# Information Theory and Probabilistic Programming

Samuel Cheng

School of ECE
University of Oklahoma

November 28, 2019

# This time

- Method of types
- Universal source coding
- Large deviation theory

## Motivation

- In previous lectures, we have introduced LLN and typical sequences. In a sense that every sequences drawn from a discrete memoryless source is typical

## Motivation

- In previous lectures, we have introduced LLN and typical sequences. In a sense that every sequences drawn from a discrete memoryless source is typical

- Take coin tossing as example again, if $Pr(Head) = 0.6$, and we throw the coin 1000 times. We expect that almost all drawn sequences with have about 600 heads. And the rest have neglible probability

## Motivation

- In previous lectures, we have introduced LLN and typical sequences. In a sense that every sequences drawn from a discrete memoryless source is typical

- Take coin tossing as example again, if $Pr(Head) = 0.6$, and we throw the coin 1000 times. We expect that almost all drawn sequences with have about 600 heads. And the rest have neglible probability

- However, sometimes we are interested in the probability of getting say 400 heads, even though we know that the probability is neglible

## Motivation

- In previous lectures, we have introduced LLN and typical sequences. In a sense that every sequences drawn from a discrete memoryless source is typical

- Take coin tossing as example again, if $Pr(Head) = 0.6$, and we throw the coin 1000 times. We expect that almost all drawn sequences with have about 600 heads. And the rest have neglible probability

- However, sometimes we are interested in the probability of getting say 400 heads, even though we know that the probability is neglible $\rightarrow$ method of types

## Motivation

By the end of the class, we will be able to solve the following nontrivial puzzle

- Tom throws a unbiased dice for 10,000 times and adds all values

## Motivation

By the end of the class, we will be able to solve the following nontrivial puzzle

- Tom throws a unbiased dice for 10,000 times and adds all values
- For whatever reason, he is not happy until the sum is at least 40,000. If not, he will just throw the dice again for 10,000

## Motivation

By the end of the class, we will be able to solve the following nontrivial puzzle

- Tom throws a unbiased dice for 10,000 times and adds all values
- For whatever reason, he is not happy until the sum is at least 40,000. If not, he will just throw the dice again for 10,000
- Now, by the time he eventually got a sequence with sum at least 40,000, *approximately how many ones in the sequence?*

## Type class

Continue with the coin-tossing example

- Recall that the probability of getting a particular sequence with 600 heads is

    $$0.6^{600}0.4^{400}$$

## Type class

Continue with the coin-tossing example

- Recall that the probability of getting a particular sequence with 600 heads is

$$0.6^{600}0.4^{400} = 2^{-1000(-0.6 \log 0.6 - 0.4 \log 0.4)}$$

## Type class

Continue with the coin-tossing example

- Recall that the probability of getting a particular sequence with 600 heads is

$$0.6^{600}0.4^{400} = 2^{-1000(-0.6 \log 0.6 - 0.4 \log 0.4)} = 2^{-NH(X)}$$

## Type class

Continue with the coin-tossing example

- Recall that the probability of getting a particular sequence with 600 heads is

$$0.6^{600}0.4^{400} = 2^{-1000(-0.6\log 0.6 - 0.4\log 0.4)} = 2^{-NH(X)}$$

- How about the probability of getting a particular sequence with 400 heads? It is

$$0.6^{400}0.4^{600} = 2^{-1000(-0.4\log 0.6 - 0.6\log 0.4)}$$

## Type class

Continue with the coin-tossing example

- Recall that the probability of getting a particular sequence with 600 heads is

$$0.6^{600}0.4^{400} = 2^{-1000(-0.6 \log 0.6 - 0.4 \log 0.4)} = 2^{-NH(X)}$$

- How about the probability of getting a particular sequence with 400 heads? It is

$$0.6^{400}0.4^{600} = 2^{-1000(-0.4 \log 0.6 - 0.6 \log 0.4)}$$
$$= 2^{-1000(-0.4 \log 0.4 - 0.6 \log 0.6 + 0.4 \log \frac{0.4}{0.6} + 0.6 \log \frac{0.6}{0.4})}$$

## Type class

Continue with the coin-tossing example

- Recall that the probability of getting a particular sequence with 600 heads is

$$0.6^{600}0.4^{400} = 2^{-1000(-0.6\log 0.6 - 0.4\log 0.4)} = 2^{-NH(X)}$$

- How about the probability of getting a particular sequence with 400 heads? It is

$$\begin{aligned}
0.6^{400}0.4^{600} &= 2^{-1000(-0.4\log 0.6 - 0.6\log 0.4)} \\
&= 2^{-1000(-0.4\log 0.4 - 0.6\log 0.6 + 0.4\log\frac{0.4}{0.6} + 0.6\log\frac{0.6}{0.4})} \\
&= 2^{-N(H(X)+KL((0.4,0.6)||(0.6,0.4)))}
\end{aligned}$$

## Type class

Continue with the coin-tossing example

- Recall that the probability of getting a particular sequence with 600 heads is

$$0.6^{600}0.4^{400} = 2^{-1000(-0.6\log 0.6 - 0.4 \log 0.4)} = 2^{-NH(X)}$$

- How about the probability of getting a particular sequence with 400 heads? It is

$$0.6^{400}0.4^{600} = 2^{-1000(-0.4\log 0.6 - 0.6\log 0.4)}$$
$$= 2^{-1000(-0.4\log 0.4 - 0.6\log 0.6 + 0.4\log\frac{0.4}{0.6} + 0.6\log\frac{0.6}{0.4})}$$
$$= 2^{-N(H(X) + KL((0.4,0.6)||(0.6,0.4)))}$$

- Every sequence with 400 heads has the same probability. And in general, sequences with the same fraction of outcomes have same probability and we can put them into the same (type) class

# Type class

- For convenience, let us denote the number of $a$ in the sequence $x^N$ as $\mathcal{N}(a|x^N)$

## Type class

- For convenience, let us denote the number of $a$ in the sequence $x^N$ as $\mathcal{N}(a|x^N)$

- Then for any valid distribution of $X$, $p(x)$, we will define a type class $T(p_X)$ as the set containing all sequences such that $\frac{\mathcal{N}(a|x^N)}{N} \approx p(a)$, $\forall a \in \mathcal{X}$

## Type class

- For convenience, let us denote the number of $a$ in the sequence $x^N$ as $\mathcal{N}(a|x^N)$

- Then for any valid distribution of $X$, $p(x)$, we will define a type class $T(p_X)$ as the set containing all sequences such that $\frac{\mathcal{N}(a|x^N)}{N} \approx p(a)$, $\forall a \in \mathcal{X}$

- Let us reserve $q(x)$ as the true distribution of $x$ (i.e., $q(Head) = 0.6$ and $q(Tail) = 0.4$). And in general, we expect all sequences drawn from the source should belongs to $T(q)$ asymptotically

## Type class

- For convenience, let us denote the number of $a$ in the sequence $x^N$ as $\mathcal{N}(a|x^N)$
- Then for any valid distribution of $X$, $p(x)$, we will define a type class $T(p_X)$ as the set containing all sequences such that $\frac{\mathcal{N}(a|x^N)}{N} \approx p(a)$, $\forall a \in \mathcal{X}$
- Let us reserve $q(x)$ as the true distribution of $x$ (i.e., $q(Head) = 0.6$ and $q(Tail) = 0.4$). And in general, we expect all sequences drawn from the source should belongs to $T(q)$ asymptotically
- Let's also refer $p_{x^N}$ as the empirical distribution of $x^N$. That is $p_{x^N}(a) = \frac{\mathcal{N}(a|x^N)}{N}$. So $T(p_{x^N})$ is the type class containing $x^N$

# Example

Let $\mathcal{X} \in \{1, 2, 3\}$ and $x^N = 11321$

- $p_{x^N}(1) = \frac{3}{5}$,

## Example

Let $\mathcal{X} \in \{1, 2, 3\}$ and $x^N = 11321$

- $p_{x^N}(1) = \frac{3}{5}$, $p_{x^N}(2) = \frac{1}{5}$, $p_{x^N}(3) = \frac{1}{5}$

## Example

Let $\mathcal{X} \in \{1, 2, 3\}$ and $x^N = 11321$

- $p_{x^N}(1) = \frac{3}{5}$, $p_{x^N}(2) = \frac{1}{5}$, $p_{x^N}(3) = \frac{1}{5}$
- $T(p_{x^N}) = \{11123, 11132, 11231, 11321, \cdots\}$ containing all sequences with three 1's, one 2, and one 3

## Example

Let $\mathcal{X} \in \{1, 2, 3\}$ and $x^N = 11321$

- $p_{x^N}(1) = \frac{3}{5}$, $p_{x^N}(2) = \frac{1}{5}$, $p_{x^N}(3) = \frac{1}{5}$
- $T(p_{x^N}) = \{11123, 11132, 11231, 11321, \cdots\}$ containing all sequences with three 1's, one 2, and one 3
- $|T(p_{x^N})| = \frac{5!}{3!1!1!} = 20$.

## Example

Let $\mathcal{X} \in \{1, 2, 3\}$ and $x^N = 11321$

- $p_{x^N}(1) = \frac{3}{5}$, $p_{x^N}(2) = \frac{1}{5}$, $p_{x^N}(3) = \frac{1}{5}$
- $T(p_{x^N}) = \{11123, 11132, 11231, 11321, \cdots\}$ containing all sequences with three 1's, one 2, and one 3
- $|T(p_{x^N})| = \frac{5!}{3!1!1!} = 20$. In general,

$$|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))! \cdots}$$

## Example

Let $\mathcal{X} \in \{1, 2, 3\}$ and $x^N = 11321$

- $p_{x^N}(1) = \frac{3}{5}$, $p_{x^N}(2) = \frac{1}{5}$, $p_{x^N}(3) = \frac{1}{5}$
- $T(p_{x^N}) = \{11123, 11132, 11231, 11321, \cdots\}$ containing all sequences with three 1's, one 2, and one 3
- $|T(p_{x^N})| = \frac{5!}{3!1!1!} = 20$. In general,

$$|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))! \cdots}$$

Actually we don't care too much what $|T(p)|$ is exactly. We will provide bounds for $|T(p)|$ as we come back later on

## Example

Let $\mathcal{X} \in \{1, 2, 3\}$ and $x^N = 11321$

- $p_{x^N}(1) = \frac{3}{5}$, $p_{x^N}(2) = \frac{1}{5}$, $p_{x^N}(3) = \frac{1}{5}$
- $T(p_{x^N}) = \{11123, 11132, 11231, 11321, \cdots\}$ containing all sequences with three 1's, one 2, and one 3
- $|T(p_{x^N})| = \frac{5!}{3!1!1!} = 20$. In general,

$$|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))! \cdots}$$

  Actually we don't care too much what $|T(p)|$ is exactly. We will provide bounds for $|T(p)|$ as we come back later on

- And for any sequence $\mathbf{y}$ in $T(p_{x^N})$, $p(\mathbf{y}) = q(1)^3 q(2) q(3)$, where $q(\cdot)$ is the true distribution

# Type sequence probability

Even though we have seen that in the coin toss example, let's restate it more formally.

### Theorem 1

If $x^N \in T(p)$ and $q(\cdot)$ is the true distribution of $X$, the probability of getting $x^N$ from sampling $q(\cdot)$ for $N$ times, as denoted as $q^N(x^N)$, is given by

$$2^{-N(H(p)+KL(p||q))}$$

# Type sequence probability

Even though we have seen that in the coin toss example, let's restate it more formally.

### Theorem 1

If $x^N \in T(p)$ and $q(\cdot)$ is the true distribution of $X$, the probability of getting $x^N$ from sampling $q(\cdot)$ for $N$ times, as denoted as $q^N(x^N)$, is given by

$$2^{-N(H(p)+KL(p||q))}$$

### Proof

$$q^N(x^N) = \prod_{i=1}^{N} q(x_i) = 2^{\sum_{i=1}^{N} \log q(x_i)}$$

# Type sequence probability

Even though we have seen that in the coin toss example, let's restate it more formally.

## Theorem 1

If $x^N \in T(p)$ and $q(\cdot)$ is the true distribution of $X$, the probability of getting $x^N$ from sampling $q(\cdot)$ for $N$ times, as denoted as $q^N(x^N)$, is given by

$$2^{-N(H(p)+KL(p||q))}$$

## Proof

$$q^N(x^N) = \prod_{i=1}^{N} q(x_i) = 2^{\sum_{i=1}^{N} \log q(x_i)} = 2^{\sum_{a \in \mathcal{X}} \mathcal{N}(a|x^N) \log q(a)}$$

# Type sequence probability

Even though we have seen that in the coin toss example, let's restate it more formally.

### Theorem 1

If $x^N \in T(p)$ and $q(\cdot)$ is the true distribution of $X$, the probability of getting $x^N$ from sampling $q(\cdot)$ for $N$ times, as denoted as $q^N(x^N)$, is given by

$$2^{-N(H(p)+KL(p||q))}$$

### Proof

$$q^N(x^N) = \prod_{i=1}^{N} q(x_i) = 2^{\sum_{i=1}^{N} \log q(x_i)} = 2^{\sum_{a \in \mathcal{X}} \mathcal{N}(a|x^N) \log q(a)}$$

$$= 2^{-N \sum_{a \in \mathcal{X}} -p_{x_N}(a) \log q(a)}$$

# Type sequence probability

Even though we have seen that in the coin toss example, let's restate it more formally.

## Theorem 1

If $x^N \in T(p)$ and $q(\cdot)$ is the true distribution of $X$, the probability of getting $x^N$ from sampling $q(\cdot)$ for $N$ times, as denoted as $q^N(x^N)$, is given by

$$2^{-N(H(p)+KL(p||q))}$$

## Proof

$$q^N(x^N) = \prod_{i=1}^{N} q(x_i) = 2^{\sum_{i=1}^{N} \log q(x_i)} = 2^{\sum_{a \in \mathcal{X}} \mathcal{N}(a|x^N) \log q(a)}$$

$$= 2^{-N \sum_{a \in \mathcal{X}} -p_{x_N}(a) \log q(a)} = 2^{-N\left(-\sum_{a \in \mathcal{X}} p(a) \log p(a) - \sum_{a \in \mathcal{X}} p(a) \log \frac{p(a)}{q(a)}\right)}$$

## Type sequence probability

Even though we have seen that in the coin toss example, let's restate it more formally.

### Theorem 1

If $x^N \in T(p)$ and $q(\cdot)$ is the true distribution of $X$, the probability of getting $x^N$ from sampling $q(\cdot)$ for $N$ times, as denoted as $q^N(x^N)$, is given by

$$2^{-N(H(p)+KL(p||q))}$$

### Proof

$$q^N(x^N) = \prod_{i=1}^{N} q(x_i) = 2^{\sum_{i=1}^{N} \log q(x_i)} = 2^{\sum_{a \in \mathcal{X}} \mathcal{N}(a|x^N) \log q(a)}$$

$$= 2^{-N \sum_{a \in \mathcal{X}} -p_{x_N}(a) \log q(a)} = 2^{-N\left(-\sum_{a \in \mathcal{X}} p(a) \log p(a) - \sum_{a \in \mathcal{X}} p(a) \log \frac{p(a)}{q(a)}\right)}$$

$$= 2^{-N(H(p)+KL(p||q))}$$

# Probability of a sequence in the "typical" class

If $x^N \in T(q)$, where $q(\cdot)$ is the true distribution of $X$, then

$$q^N(x^N) = 2^{-NH(q)} = 2^{-NH(X)}$$

# Probability of a sequence in the "typical" class

If $x^N \in T(q)$, where $q(\cdot)$ is the true distribution of $X$, then

$$q^N(x^N) = 2^{-NH(q)} = 2^{-NH(X)}$$

### Remarks

- Note that the probability is exactly equal to $2^{-NH(X)}$

# Probability of a sequence in the "typical" class

If $x^N \in T(q)$, where $q(\cdot)$ is the true distribution of $X$, then

$$q^N(x^N) = 2^{-NH(q)} = 2^{-NH(X)}$$

### Remarks

- Note that the probability is exactly equal to $2^{-NH(X)}$
- Recall that this is the probability of a typical sequence supposed to be. Therefore, any $x^N$ in $T(q)$ is a typical sequence ($T(q) \subset A_\epsilon^N(X)$)

# Set of all empirical distribution $\mathcal{P}_N(X)$

Denote $\mathcal{P}_N(X)$ as the set of all empirical distribution of $X$ in a length-$N$ sequence

# Set of all empirical distribution $\mathcal{P}_N(X)$

Denote $\mathcal{P}_N(X)$ as the set of all empirical distribution of $X$ in a length-$N$ sequence

## Example

If $X \in \{0, 1\}$,
$$\mathcal{P}_N(X) = \left\{ (p_X(0), p_X(1)) : \left( \frac{0}{N}, \frac{N}{N} \right), \left( \frac{1}{N}, \frac{N-1}{N} \right), \cdots, \left( \frac{N}{N}, \frac{0}{N} \right) \right\}$$

Note that $|\mathcal{P}_N(X)| = N + 1$

# Set of all empirical distribution $\mathcal{P}_N(X)$

Denote $\mathcal{P}_N(X)$ as the set of all empirical distribution of $X$ in a length-$N$ sequence

## Example

If $X \in \{0, 1\}$,
$$\mathcal{P}_N(X) = \left\{ (p_X(0), p_X(1)) : \left( \frac{0}{N}, \frac{N}{N} \right), \left( \frac{1}{N}, \frac{N-1}{N} \right), \cdots, \left( \frac{N}{N}, \frac{0}{N} \right) \right\}$$

Note that $|\mathcal{P}_N(X)| = N + 1$

- Since a type is uniquely characterized by a distribution of $X$ in a length-$N$ sequence

# Set of all empirical distribution $\mathcal{P}_N(X)$

Denote $\mathcal{P}_N(X)$ as the set of all empirical distribution of $X$ in a length-$N$ sequence

### Example

If $X \in \{0, 1\}$,
$$\mathcal{P}_N(X) = \left\{ (p_X(0), p_X(1)) : \left( \frac{0}{N}, \frac{N}{N} \right), \left( \frac{1}{N}, \frac{N-1}{N} \right), \cdots, \left( \frac{N}{N}, \frac{0}{N} \right) \right\}$$

Note that $|\mathcal{P}_N(X)| = N + 1$

- Since a type is uniquely characterized by a distribution of $X$ in a length-$N$ sequence
- Each element $p$ of $\mathcal{P}_N(X)$ corresponds a type $T(p)$

# Set of all empirical distribution $\mathcal{P}_N(X)$

Denote $\mathcal{P}_N(X)$ as the set of all empirical distribution of $X$ in a length-$N$ sequence

## Example

If $X \in \{0, 1\}$,
$$\mathcal{P}_N(X) = \left\{ (p_X(0), p_X(1)) : \left( \frac{0}{N}, \frac{N}{N} \right), \left( \frac{1}{N}, \frac{N-1}{N} \right), \cdots, \left( \frac{N}{N}, \frac{0}{N} \right) \right\}$$

Note that $|\mathcal{P}_N(X)| = N + 1$

- Since a type is uniquely characterized by a distribution of $X$ in a length-$N$ sequence
- Each element $p$ of $\mathcal{P}_N(X)$ corresponds a type $T(p)$
- Number of types is $|\mathcal{P}_N(X)|$

## Number of types

It is not too difficult to count the exact number of types. But in practice, we don't quite bother with it as long as we know that the number is relatively "small"

### Theorem 2

$$|\mathcal{P}_N(X)| \leq (N+1)^{|\mathcal{X}|}$$

# Number of types

It is not too difficult to count the exact number of types. But in practice, we don't quite bother with it as long as we know that the number is relatively "small"

### Theorem 2

$$|\mathcal{P}_N(X)| \leq (N+1)^{|\mathcal{X}|}$$

### Proof

Note that each type is specified by the empirical probability of each outcome of $X$. And the possible values of the empirical probabilities are $\frac{0}{N}, \frac{1}{N}, \cdots, \frac{N}{N}$ ($N+1$ of them).

# Number of types

It is not too difficult to count the exact number of types. But in practice, we don't quite bother with it as long as we know that the number is relatively "small"

## Theorem 2

$$|\mathcal{P}_N(X)| \leq (N+1)^{|\mathcal{X}|}$$

## Proof

Note that each type is specified by the empirical probability of each outcome of $X$. And the possible values of the empirical probabilities are $\frac{0}{N}$, $\frac{1}{N}$, $\cdots$, $\frac{N}{N}$ ($N+1$ of them). Since there are $|\mathcal{X}|$ elements, the number of types is bounded by $(N+1)^{|\mathcal{X}|}$

## Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

### Theorem 3

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \le |T(p)| \le 2^{NH(p)}$$

# Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

**Theorem 3**

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

**Proof**

Let's assume $p(\cdot)$ is the actual distribution of $X$ here

$$1 \geq \sum_{x^N \in T(p)} p^N(x^N)$$

# Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

### Theorem 3

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

### Proof

Let's assume $p(\cdot)$ is the actual distribution of $X$ here

$$1 \geq \sum_{x^N \in T(p)} p^N(x^N) = \sum_{x^N \in T(p)} 2^{-NH(p)} = |T(p)| 2^{-NH(p)}$$

# Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

**Theorem 3**

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

**Proof**

Let's assume $p(\cdot)$ is the actual distribution of $X$ here

$$1 \geq \sum_{x^N \in T(p)} p^N(x^N) = \sum_{x^N \in T(p)} 2^{-NH(p)} = |T(p)| 2^{-NH(p)}$$

$$1 = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(\hat{p}))$$

# Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

## Theorem 3

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

## Proof

Let's assume $p(\cdot)$ is the actual distribution of $X$ here

$$1 \geq \sum_{x^N \in T(p)} p^N(x^N) = \sum_{x^N \in T(p)} 2^{-NH(p)} = |T(p)| 2^{-NH(p)}$$

$$1 = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(\hat{p})) \leq \sum_{\hat{p} \in \mathcal{P}_N} \max_{\tilde{p}} Pr(T(\tilde{p}))$$

# Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

## Theorem 3

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

## Proof

Let's assume $p(\cdot)$ is the actual distribution of $X$ here

$$1 \geq \sum_{x^N \in T(p)} p^N(x^N) = \sum_{x^N \in T(p)} 2^{-NH(p)} = |T(p)| 2^{-NH(p)}$$

$$1 = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(\hat{p})) \leq \sum_{\hat{p} \in \mathcal{P}_N} \max_{\tilde{p}} Pr(T(\tilde{p})) = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(p))$$

# Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

### Theorem 3

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

### Proof

Let's assume $p(\cdot)$ is the actual distribution of $X$ here

$$1 \geq \sum_{x^N \in T(p)} p^N(x^N) = \sum_{x^N \in T(p)} 2^{-NH(p)} = |T(p)| 2^{-NH(p)}$$

$$1 = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(\hat{p})) \leq \sum_{\hat{p} \in \mathcal{P}_N} \max_{\tilde{p}} Pr(T(\tilde{p})) = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(p)) \leq (N+1)^{|\mathcal{X}|} Pr(T(p))$$

# Size of a type class

Recall that $|T(p)| = \frac{N!}{(Np(x_1))!(Np(x_2))!(Np(x_3))!\cdots}$ but the following bounds are much more useful in practice

## Theorem 3

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

## Proof

Let's assume $p(\cdot)$ is the actual distribution of $X$ here

$$1 \geq \sum_{x^N \in T(p)} p^N(x^N) = \sum_{x^N \in T(p)} 2^{-NH(p)} = |T(p)| 2^{-NH(p)}$$

$$1 = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(\hat{p})) \leq \sum_{\hat{p} \in \mathcal{P}_N} \max_{\tilde{p}} Pr(T(\tilde{p})) = \sum_{\hat{p} \in \mathcal{P}_N} Pr(T(p)) \leq (N+1)^{|\mathcal{X}|} Pr(T(p))$$

$$= (N+1)^{|\mathcal{X}|} |T(p)| 2^{-NH(p)}$$

# Probability of a type class

**Theorem 4**

Let the true distribution of $X$ is $q(\cdot)$, then

$$\frac{2^{-N(KL(p||q))}}{(N+1)^{|\mathcal{X}|}} \leq Pr(T(p)) \leq 2^{-N(KL(p||q))}$$

# Probability of a type class

### Theorem 4

Let the true distribution of $X$ is $q(\cdot)$, then

$$\frac{2^{-N(KL(p||q))}}{(N+1)^{|\mathcal{X}|}} \leq Pr(T(p)) \leq 2^{-N(KL(p||q))}$$

### Proof

From Theorem 1, each sequence in $T(p)$ has probability $2^{-N(H(p)+KL(p||q))}$ and since $\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$ from Theorem 3,

$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} 2^{-N(H(p)+KL(p||q))} \leq Pr(T(p)) \leq 2^{NH(p)} 2^{-N(H(p)+KL(p||q))}$$

# Summary of type

- Type class $T(p)$ contains all sequences with empirical distribution of $p$. That is,

$$T(p) = \left\{ x^N : \frac{\mathcal{N}(a|x^N)}{N} = p(a) \right\}$$

## Summary of type

- Type class $T(p)$ contains all sequences with empirical distribution of $p$. That is,
$$T(p) = \left\{ x^N : \frac{\mathcal{N}(a|x^N)}{N} = p(a) \right\}$$

- All sequences in the type class $T(p)$ has the same probability ($q(\cdot)$ is the true distribution)
$$q^N(x^N) = 2^{-N(H(p)+KL(p||q))}$$

## Summary of type

- Type class $T(p)$ contains all sequences with empirical distribution of $p$. That is,
$$T(p) = \left\{ x^N : \frac{\mathcal{N}(a|x^N)}{N} = p(a) \right\}$$

- All sequences in the type class $T(p)$ has the same probability ($q(\cdot)$ is the true distribution)
$$q^N(x^N) = 2^{-N(H(p)+KL(p||q)}$$

- There are about $2^{NH(p)}$ sequences in $T(p)$
$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

## Summary of type

- Type class $T(p)$ contains all sequences with empirical distribution of $p$. That is,
$$T(p) = \left\{ x^N : \frac{\mathcal{N}(a|x^N)}{N} = p(a) \right\}$$

- All sequences in the type class $T(p)$ has the same probability ($q(\cdot)$ is the true distribution)
$$q^N(x^N) = 2^{-N(H(p)+KL(p||q))}$$

- There are about $2^{NH(p)}$ sequences in $T(p)$
$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

- Probability of getting a sequence in $T(p)$ is about $2^{-N(KL(p||q))}$. More precisely,
$$\frac{2^{-N(KL(p||q))}}{(N+1)^{|\mathcal{X}|}} \leq Pr(T(p)) \leq 2^{-N(KL(p||q))}$$

## Summary of type

- Type class $T(p)$ contains all sequences with empirical distribution of $p$. That is,
$$T(p) = \left\{ x^N : \frac{\mathcal{N}(a|x^N)}{N} = p(a) \right\}$$

- All sequences in the type class $T(p)$ has the same probability ($q(\cdot)$ is the true distribution)
$$q^N(x^N) = 2^{-N(H(p)+KL(p||q))}$$

- There are about $2^{NH(p)}$ sequences in $T(p)$
$$\frac{1}{(N+1)^{|\mathcal{X}|}} 2^{NH(p)} \leq |T(p)| \leq 2^{NH(p)}$$

- Probability of getting a sequence in $T(p)$ is about $2^{-N(KL(p||q))}$. More precisely,
$$\frac{2^{-N(KL(p||q))}}{(N+1)^{|\mathcal{X}|}} \leq Pr(T(p)) \leq 2^{-N(KL(p||q))}$$

- There are $(N+1)^{|\mathcal{X}|}$ types

## Rationale

- For the compression scheme (such as Huffmann coding) that we discussed earlier in this class, one needs to know the source distribution ahead to design the encoder and decoder

## Rationale

- For the compression scheme (such as Huffmann coding) that we discussed earlier in this class, one needs to know the source distribution ahead to design the encoder and decoder

- Question: Is it possible to construct compression scheme without knowing the source distibution and still performs as good?

## Rationale

- For the compression scheme (such as Huffmann coding) that we discussed earlier in this class, one needs to know the source distribution ahead to design the encoder and decoder

- Question: Is it possible to construct compression scheme without knowing the source distibution and still performs as good?

- Answer: Yes. At least theoretically $\rightarrow$ universal source coding

## Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

## Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

### Proof

Let $R_N = R - |\mathcal{X}|\frac{\log(N+1)}{N}$, and consider the set of sequences $A = \{x^N : H(p_{x^N}) < R_N\}$ as the code book.

# Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

### Proof

Let $R_N = R - |\mathcal{X}| \frac{\log(N+1)}{N}$, and consider the set of sequences $A = \{x^N : H(p_{x^N}) < R_N\}$ as the code book. Note that the rate is $< R$ as

$$|A| = \sum_{p : H(p) < R_N} |T(p)|$$

# Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

### Proof

Let $R_N = R - |\mathcal{X}|\frac{\log(N+1)}{N}$, and consider the set of sequences
$A = \{x^N : H(p_{x^N}) < R_N\}$ as the code book. Note that the rate is $< R$ as
$$|A| = \sum_{p:H(p)<R_N} |T(p)| \leq \sum_{p:H(p)<R_N} 2^{NH(p)}$$

## Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

### Proof

Let $R_N = R - |\mathcal{X}|\frac{\log(N+1)}{N}$, and consider the set of sequences $A = \{x^N : H(p_{x^N}) < R_N\}$ as the code book. Note that the rate is $< R$ as

$$|A| = \sum_{p:H(p)<R_N} |T(p)| \leq \sum_{p:H(p)<R_N} 2^{NH(p)} < \sum_{p:H(p)<R_N} 2^{NR_N}$$

## Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

### Proof

Let $R_N = R - |\mathcal{X}|\frac{\log(N+1)}{N}$, and consider the set of sequences $A = \{x^N : H(p_{x^N}) < R_N\}$ as the code book. Note that the rate is $< R$ as

$$|A| = \sum_{p:H(p)<R_N} |T(p)| \leq \sum_{p:H(p)<R_N} 2^{NH(p)} < \sum_{p:H(p)<R_N} 2^{NR_N}$$

$$\leq (N+1)^{|\mathcal{X}|} 2^{NR_N}$$

## Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

### Proof

Let $R_N = R - |\mathcal{X}|\frac{\log(N+1)}{N}$, and consider the set of sequences $A = \{x^N : H(p_{x^N}) < R_N\}$ as the code book. Note that the rate is $< R$ as
$$|A| = \sum_{p:H(p)<R_N} |T(p)| \leq \sum_{p:H(p)<R_N} 2^{NH(p)} < \sum_{p:H(p)<R_N} 2^{NR_N}$$
$$\leq (N+1)^{|\mathcal{X}|}2^{NR_N} = 2^{N\left(R_N+|\mathcal{X}|\frac{\log(N+1)}{N}\right)} = 2^{NR}$$

## Theory of universal source coding

Given any source $Q$ with $H(Q) < R$, there exists a length-$N$ universal code of rate $R$ such that the source can be decoded losslessly as $N \to \infty$

### Proof

Let $R_N = R - |\mathcal{X}|\frac{\log(N+1)}{N}$, and consider the set of sequences
$A = \{x^N : H(p_{x^N}) < R_N\}$ as the code book. Note that the rate is $< R$ as

$$|A| = \sum_{p:H(p)<R_N} |T(p)| \leq \sum_{p:H(p)<R_N} 2^{NH(p)} < \sum_{p:H(p)<R_N} 2^{NR_N}$$

$$\leq (N+1)^{|\mathcal{X}|}2^{NR_N} = 2^{N\left(R_N+|\mathcal{X}|\frac{\log(N+1)}{N}\right)} = 2^{NR}$$

- Encoder: given input, check if input is in $A$, output index if so. Otherwise, declare failure
- Decoder: simply map index back to the sequence

# Theory of universal source coding

### Proof (con't)

Note that the probability of error $P_e$ is given by

$$P_e = \sum_{p:H(p)>R_N} Pr(T(p))$$

# Theory of universal source coding

## Proof (con't)

Note that the probability of error $P_e$ is given by

$$P_e = \sum_{p:H(p)>R_N} Pr(T(p)) \leq \sum_{p:H(p)>R_N} \max_{\tilde{p}:H(\tilde{p})>R_N} Pr(T(\tilde{p}))$$

# Theory of universal source coding

## Proof (con't)

Note that the probability of error $P_e$ is given by

$$P_e = \sum_{p:H(p)>R_N} Pr(T(p)) \leq \sum_{p:H(p)>R_N} \max_{\tilde{p}:H(\tilde{p})>R_N} Pr(T(\tilde{p}))$$

$$\leq (1+N)^{|\mathcal{X}|} 2^{-N\left(\min_{\tilde{p}:H(\tilde{p})>R_N} KL(\tilde{p}||q)\right)}$$

# Theory of universal source coding

### Proof (con't)

Note that the probability of error $P_e$ is given by

$$P_e = \sum_{p:H(p)>R_N} Pr(T(p)) \leq \sum_{p:H(p)>R_N} \max_{\tilde{p}:H(\tilde{p})>R_N} Pr(T(\tilde{p}))$$

$$\leq (1+N)^{|\mathcal{X}|} 2^{-N\left(\min_{\tilde{p}:H(\tilde{p})>R_N} KL(\tilde{p}\|q)\right)}$$

- If $H(q) < R$, as $R_N \to R$ as $N$ increases, we can find some $N_0$ such that $H(q) < R_N$ for all $N \geq N_0$

# Theory of universal source coding

### Proof (con't)

Note that the probability of error $P_e$ is given by

$$P_e = \sum_{p:H(p)>R_N} Pr(T(p)) \leq \sum_{p:H(p)>R_N} \max_{\tilde{p}:H(\tilde{p})>R_N} Pr(T(\tilde{p}))$$

$$\leq (1+N)^{|\mathcal{X}|} 2^{-N\left(\min_{\tilde{p}:H(\tilde{p})>R_N} KL(\tilde{p}\|q)\right)}$$

- If $H(q) < R$, as $R_N \to R$ as $N$ increases, we can find some $N_0$ such that $H(q) < R_N$ for all $N \geq N_0$
- Therefore, any $p$ in $\{p : H(p) > R_N\}$ cannot be the same as $q$

# Theory of universal source coding

### Proof (con't)

Note that the probability of error $P_e$ is given by

$$P_e = \sum_{p:H(p)>R_N} Pr(T(p)) \leq \sum_{p:H(p)>R_N} \max_{\tilde{p}:H(\tilde{p})>R_N} Pr(T(\tilde{p}))$$

$$\leq (1+N)^{|\mathcal{X}|} 2^{-N\left(\min_{\tilde{p}:H(\tilde{p})>R_N} KL(\tilde{p}||q)\right)}$$

- If $H(q) < R$, as $R_N \to R$ as $N$ increases, we can find some $N_0$ such that $H(q) < R_N$ for all $N \geq N_0$
- Therefore, any $p$ in $\{p : H(p) > R_N\}$ cannot be the same as $q$
- $\Rightarrow \min_{\tilde{p}:H(\tilde{p})>R_N} KL(\tilde{p}||q) > 0$ for $N \geq N_0$

# Theory of universal source coding

## Proof (con't)

Note that the probability of error $P_e$ is given by

$$P_e = \sum_{p:H(p)>R_N} Pr(T(p)) \leq \sum_{p:H(p)>R_N} \max_{\tilde{p}:H(\tilde{p})>R_N} Pr(T(\tilde{p}))$$

$$\leq (1+N)^{|\mathcal{X}|} 2^{-N\left(\min_{\tilde{p}:H(\tilde{p})>R_N} KL(\tilde{p}||q)\right)}$$

- If $H(q) < R$, as $R_N \to R$ as $N$ increases, we can find some $N_0$ such that $H(q) < R_N$ for all $N \geq N_0$
- Therefore, any $p$ in $\{p : H(p) > R_N\}$ cannot be the same as $q$
- $\Rightarrow \min_{\tilde{p}:H(\tilde{p})>R_N} KL(\tilde{p}||q) > 0$ for $N \geq N_0$
- Hence, $P_e \to 0$ as $N \to \infty$

# Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
  - Construct a dictionary including all previously seen segments

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
    - First parse segment into segments that haven't seen before $\Rightarrow$

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
  - Construct a dictionary including all previously seen segments
  - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
  - First parse segment into segments that haven't seen before $\Rightarrow$
    1
    1

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
  - Construct a dictionary including all previously seen segments
  - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
  - First parse segment into segments that haven't seen before ⇒
    1, 0
    1  2

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
  - Construct a dictionary including all previously seen segments
  - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
  - First parse segment into segments that haven't seen before $\Rightarrow$
    $\underset{1}{1}, \underset{2}{0}, \underset{3}{11}$

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
    - First parse segment into segments that haven't seen before $\Rightarrow$
      $\overset{1}{1}, \overset{2}{0}, \overset{3}{11}, \overset{4}{01}$

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
    - First parse segment into segments that haven't seen before $\Rightarrow$
      $$\underset{1}{1}, \underset{2}{0}, \underset{3}{11}, \underset{4}{01}, \underset{5}{110}$$

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
    - First parse segment into segments that haven't seen before $\Rightarrow$
      $$\underset{1}{1}, \underset{2}{0}, \underset{3}{11}, \underset{4}{01}, \underset{5}{110}, \underset{6}{111}$$

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
    - First parse segment into segments that haven't seen before $\Rightarrow$
      $\underset{1}{1}, \underset{2}{0}, \underset{3}{11}, \underset{4}{01}, \underset{5}{110}, \underset{6}{111}, \underset{7}{10}$

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
  - Construct a dictionary including all previously seen segments
  - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
  - First parse segment into segments that haven't seen before $\Rightarrow$
    $$\overset{1}{1}, \overset{2}{0}, \overset{3}{11}, \overset{4}{01}, \overset{5}{110}, \overset{6}{111}, \overset{7}{10}, \overset{8}{111}$$
  - Encode each segment into representation containing a pair of numbers:

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
    - First parse segment into segments that haven't seen before $\Rightarrow$
      $\underset{1}{1}, \underset{2}{0}, \underset{3}{11}, \underset{4}{01}, \underset{5}{110}, \underset{6}{111}, \underset{7}{10}, \underset{8}{111}$
    - Encode each segment into representation containing a pair of numbers: 1) index of segment (excluding the last bit) in the dictionary;

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
    - Construct a dictionary including all previously seen segments
    - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
    - First parse segment into segments that haven't seen before $\Rightarrow$
      $\underset{1}{1}, \underset{2}{0}, \underset{3}{11}, \underset{4}{01}, \underset{5}{110}, \underset{6}{111}, \underset{7}{10}, \underset{8}{111}$
    - Encode each segment into representation containing a pair of numbers: 1) index of segment (excluding the last bit) in the dictionary; 2) the last bit

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
  - Construct a dictionary including all previously seen segments
  - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
  - First parse segment into segments that haven't seen before $\Rightarrow$
    $\overset{1}{1}, \overset{2}{0}, \overset{3}{11}, \overset{4}{01}, \overset{5}{110}, \overset{6}{111}, \overset{7}{10}, \overset{8}{111}$
  - Encode each segment into representation containing a pair of numbers: 1) index of segment (excluding the last bit) in the dictionary; 2) the last bit $\Rightarrow (0,1), (0,0), (1,1), (2,1), (3,0), (3,1), (1,0), (6,\varnothing)$

## Lempel-Ziv coding

- Its variants are widely used by compression tools almost everywhere (zip, pkzip, tiff, etc.)
- Main ideas
  - Construct a dictionary including all previously seen segments
  - Bits needed to send a new segment can be reduced taking advantage known segment in the dictionary
- Example: let's compress 10110111011110111
  - First parse segment into segments that haven't seen before $\Rightarrow$
    $\underset{1}{1}, \underset{2}{0}, \underset{3}{11}, \underset{4}{01}, \underset{5}{110}, \underset{6}{111}, \underset{7}{10}, \underset{8}{111}$
  - Encode each segment into representation containing a pair of numbers: 1) index of segment (excluding the last bit) in the dictionary; 2) the last bit $\Rightarrow (0,1), (0,0), (1,1), (2,1), (3,0), (3,1), (1,0), (6,\varnothing)$
  - Encode representation to bit stream. Note that as the dictionary grows, number of bits needed to store the index increases $\Rightarrow$
    01000111010110011100010110

# Lempel-Ziv decoding

- Decode bitstream back to representation
  $0100011101011001110010110 \Rightarrow$
  $(0, 1), (0, 0), (1, 1), (2, 1), (3, 0), (3, 1), (1, 0), (6, \varnothing)$
- Build dictionary and decode

# Lempel-Ziv decoding

- Decode bitstream back to representation
  0100011101011001110010110 $\Rightarrow$
  $(0,1), (0,0), (1,1), (2,1), (3,0), (3,1), (1,0), (6,\varnothing)$
- Build dictionary and decode

$$1$$
$$1$$

$\Rightarrow 1$

# Lempel-Ziv decoding

- Decode bitstream back to representation
  01000111010110011100010110 $\Rightarrow$
  $(0,1), (0,0), (1,1), (2,1), (3,0), (3,1), (1,0), (6,\varnothing)$
- Build dictionary and decode

$$
\begin{array}{cc}
1 & 2 \\
1 & 0
\end{array}
$$

$\Rightarrow 10$

# Lempel-Ziv decoding

- Decode bitstream back to representation
  0100011101011001110010110 $\Rightarrow$
  $(0,1),(0,0),(1,1),(2,1),(3,0),(3,1),(1,0),(6,\varnothing)$
- Build dictionary and decode

$$
\begin{array}{ccc}
1 & 2 & 3 \\
1 & 0 & 11
\end{array}
$$

$\Rightarrow 1011$

## Lempel-Ziv decoding

- Decode bitstream back to representation
  $01000111010110011100101110 \Rightarrow$
  $(0,1), (0,0), (1,1), (2,1), (3,0), (3,1), (1,0), (6,\varnothing)$
- Build dictionary and decode

$$
\begin{array}{cccc}
1 & 2 & 3 & 4 \\
1 & 0 & 11 & 01
\end{array}
$$

$\Rightarrow 101101$

# Lempel-Ziv decoding

- Decode bitstream back to representation
  $0100011101011001110010110 \Rightarrow$
  $(0,1),(0,0),(1,1),(2,1),(3,0),(3,1),(1,0),(6,\varnothing)$
- Build dictionary and decode

$$
\begin{array}{ccccc}
1 & 2 & 3 & 4 & 5 \\
1 & 0 & 11 & 01 & 110
\end{array}
$$

$\Rightarrow 101101110$

# Lempel-Ziv decoding

- Decode bitstream back to representation
  $0100011101011001110010110 \Rightarrow$
  $(0,1),(0,0),(1,1),(2,1),(3,0),(3,1),(1,0),(6,\varnothing)$
- Build dictionary and decode

$$
\begin{array}{cccccc}
1 & 2 & 3 & 4 & 5 & 6 \\
1 & 0 & 11 & 01 & 110 & 111
\end{array}
$$

$\Rightarrow 101101110111$

## Lempel-Ziv decoding

- Decode bitstream back to representation
  0100011101011001110010110 $\Rightarrow$
  $(0,1),(0,0),(1,1),(2,1),(3,0),(3,1),(1,0),(6,\varnothing)$
- Build dictionary and decode

| 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|
| 1 | 0 | 11 | 01 | 110 | 111 | 10 |

$\Rightarrow$ 10110111011110

# Lempel-Ziv decoding

- Decode bitstream back to representation
  $0100011101011001110010110 \Rightarrow$
  $(0,1),(0,0),(1,1),(2,1),(3,0),(3,1),(1,0),(6,\varnothing)$
- Build dictionary and decode

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|
| 1 | 0 | 11 | 01 | 110 | 111 | 10 | 111 |

$\Rightarrow 10110111011110111$

# Motivation

- Let's revisit some coin tossing example. Say if a coin is fair, and we toss if for 1000 times, we know that we will almost always get 500 heads. So getting, say, 400 heads has neglible probability

## Motivation

- Let's revisit some coin tossing example. Say if a coin is fair, and we toss if for 1000 times, we know that we will almost always get 500 heads. So getting, say, 400 heads has neglible probability

- However, if we insist finding the probability of getting 400 heads, from discussion up to now, we know that it is just

$$Pr(T((0.4, 0.6))) \approx 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))}$$

# Motivation

- Let's revisit some coin tossing example. Say if a coin is fair, and we toss if for 1000 times, we know that we will almost always get 500 heads. So getting, say, 400 heads has neglible probability

- However, if we insist finding the probability of getting 400 heads, from discussion up to now, we know that it is just

$$Pr(T((0.4, 0.6))) \approx 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))}$$

- Now, what if we are interested in the probability of a more general case? Say what is the probability of getting $> 300$ and $< 400$ heads?

## Sanov's Theorem

Let $\mathcal{E} = \{p : 0.3 \leq p(Head) \leq 0.4\}$ and $q(\cdot) = (0.5, 0.5)$ is the true distribution, then

$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_{1000})$$

## Sanov's Theorem

Let $\mathcal{E} = \{p : 0.3 \leq p(Head) \leq 0.4\}$ and $q(\cdot) = (0.5, 0.5)$ is the true distribution, then

$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_{1000}) = \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} Pr(T(p))$$

## Sanov's Theorem

Let $\mathcal{E} = \{p : 0.3 \leq p(Head) \leq 0.4\}$ and $q(\cdot) = (0.5, 0.5)$ is the true distribution, then

$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_{1000}) = \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} Pr(T(p)) \approx \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} 2^{-1000(KL(p||q))}$$

## Sanov's Theorem

Let $\mathcal{E} = \{p : 0.3 \leq p(Head) \leq 0.4\}$ and $q(\cdot) = (0.5, 0.5)$ is the true distribution, then

$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_{1000}) = \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} Pr(T(p)) \approx \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} 2^{-1000(KL(p||q))}$$

$$= 2^{-1000(KL((0.4, 0.6)||(0.5, 0.5)))} + 2^{-1000(KL((0.399, 0.601)||(0.5, 0.5)))} +$$

$$2^{-1000(KL((0.398, 0.602)||(0.5, 0.5)))} + \cdots + 2^{-1000(KL((0.3, 0.7)||(0.5, 0.5)))}$$

## Sanov's Theorem

Let $\mathcal{E} = \{p : 0.3 \leq p(Head) \leq 0.4\}$ and $q(\cdot) = (0.5, 0.5)$ is the true distribution, then

$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_{1000}) = \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} Pr(T(p)) \approx \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} 2^{-1000(KL(p||q))}$$

$$= 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))} + 2^{-1000(KL((0.399,0.601)||(0.5,0.5)))} +$$

$$2^{-1000(KL((0.398,0.602)||(0.5,0.5)))} + \cdots + 2^{-1000(KL((0.3,0.7)||(0.5,0.5)))}$$

$$\leq |\mathcal{P}_{1000}| 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))}$$

## Sanov's Theorem

Let $\mathcal{E} = \{p : 0.3 \leq p(Head) \leq 0.4\}$ and $q(\cdot) = (0.5, 0.5)$ is the true distribution, then

$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_{1000}) = \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} Pr(T(p)) \approx \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} 2^{-1000(KL(p||q))}$$

$$= 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))} + 2^{-1000(KL((0.399,0.601)||(0.5,0.5)))} +$$

$$2^{-1000(KL((0.398,0.602)||(0.5,0.5)))} + \cdots + 2^{-1000(KL((0.3,0.7)||(0.5,0.5)))}$$

$$\leq |\mathcal{P}_{1000}| 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))}$$

### Sanov's Theorem

Let $X_1, X_2, \cdots, X_N$ be i.i.d. $\sim q(\cdot)$ and $\mathcal{E}$ be a set of distribution. Then
$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_N) \leq (N+1)^{|\mathcal{X}|} 2^{-N(KL(p^*||q))},$$

where $p^* = \arg\min_{p \in \mathcal{E}} KL(p||q)$.

## Sanov's Theorem

Let $\mathcal{E} = \{p : 0.3 \leq p(Head) \leq 0.4\}$ and $q(\cdot) = (0.5, 0.5)$ is the true distribution, then

$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_{1000}) = \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} Pr(T(p)) \approx \sum_{p \in \mathcal{E} \cap \mathcal{P}_{1000}} 2^{-1000(KL(p||q))}$$

$$= 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))} + 2^{-1000(KL((0.399,0.601)||(0.5,0.5)))} +$$

$$2^{-1000(KL((0.398,0.602)||(0.5,0.5)))} + \cdots + 2^{-1000(KL((0.3,0.7)||(0.5,0.5)))}$$

$$\leq |\mathcal{P}_{1000}| 2^{-1000(KL((0.4,0.6)||(0.5,0.5)))}$$

### Sanov's Theorem

Let $X_1, X_2, \cdots, X_N$ be i.i.d. $\sim q(\cdot)$ and $\mathcal{E}$ be a set of distribution. Then
$$Pr(\mathcal{E}) = Pr(\mathcal{E} \cap \mathcal{P}_N) \leq (N+1)^{|\mathcal{X}|} 2^{-N(KL(p^*||q))},$$

where $p^* = \arg\min_{p \in \mathcal{E}} KL(p||q)$. Moreover, given a rather weak condition (closure of interior of $\mathcal{E}$ is $\mathcal{E}$ itself), we have

$$\frac{1}{N} \log Pr(\mathcal{E}) \to -KL(p^*||q)$$

# Conditional limit theorem

- The first part of Sanov's Theorem is easy to show as similar to the example. However, the second half will need some more math background (mostly mathematical analysis) to understand the proof and so we will skip it here

# Conditional limit theorem

- The first part of Sanov's Theorem is easy to show as similar to the example. However, the second half will need some more math background (mostly mathematical analysis) to understand the proof and so we will skip it here

- The latter part of Sanov's Theorem suggests that the probability of getting $\mathcal{E}$ is the same as the probability of getting $T(p^*)$

## Conditional limit theorem

- The first part of Sanov's Theorem is easy to show as similar to the example. However, the second half will need some more math background (mostly mathematical analysis) to understand the proof and so we will skip it here

- The latter part of Sanov's Theorem suggests that the probability of getting $\mathcal{E}$ is the same as the probability of getting $T(p^*)$

- It turns out that we can claim something stronger. We will state the theorem below without proof

# Conditional limit theorem

- The first part of Sanov's Theorem is easy to show as similar to the example. However, the second half will need some more math background (mostly mathematical analysis) to understand the proof and so we will skip it here

- The latter part of Sanov's Theorem suggests that the probability of getting $\mathcal{E}$ is the same as the probability of getting $T(p^*)$

- It turns out that we can claim something stronger. We will state the theorem below without proof

### Conditional limit theorem

Let $\mathcal{E}$ be a closed convex subset of $\mathcal{P}$ (the set of all distributions) and $q(\cdot)$ be the true distribution which is $\notin \mathcal{E}$.

# Conditional limit theorem

- The first part of Sanov's Theorem is easy to show as similar to the example. However, the second half will need some more math background (mostly mathematical analysis) to understand the proof and so we will skip it here

- The latter part of Sanov's Theorem suggests that the probability of getting $\mathcal{E}$ is the same as the probability of getting $T(p^*)$

- It turns out that we can claim something stronger. We will state the theorem below without proof

## Conditional limit theorem

Let $\mathcal{E}$ be a closed convex subset of $\mathcal{P}$ (the set of all distributions) and $q(\cdot)$ be the true distribution which is $\notin \mathcal{E}$. If $x_1, x_2, \cdots, x_N$ are drawn from $q(\cdot)$ and we know that $p_{x_N} \in \mathcal{E}$, then

$$\frac{\mathcal{N}(a|x_N)}{N} \to p^*(a)$$

in probability as $N \to \infty$

# Examples

## Coin toss

- Let's go back to our previous example. If we throw a fair coin 1000 times and some one tells you that there are 300 to 400 heads, recall
$$\mathcal{E} = \{0.3 \leq p(Head) \leq 0.4\}$$

# Examples

## Coin toss

- Let's go back to our previous example. If we throw a fair coin 1000 times and some one tells you that there are 300 to 400 heads, recall
  $$\mathcal{E} = \{0.3 \le p(Head) \le 0.4\}$$

- Since apparently,
  $$p^* = \arg\min_{p \in \mathcal{E}} KL(p||(0.5, 0.5)) = (0.4, 0.6)$$

## Examples

---

### Coin toss

- Let's go back to our previous example. If we throw a fair coin 1000 times and some one tells you that there are 300 to 400 heads, recall
$$\mathcal{E} = \{0.3 \leq p(Head) \leq 0.4\}$$

- Since apparently,
$$p^* = \arg \min_{p \in \mathcal{E}} KL(p||(0.5, 0.5)) = (0.4, 0.6)$$

- By conditional limit theorem, knowing the the number of head is within the range, the coin behaves like a biased coin with $p(Head) = 0.4$

## Examples

### Coin toss

- Let's go back to our previous example. If we throw a fair coin 1000 times and some one tells you that there are 300 to 400 heads, recall
$$\mathcal{E} = \{0.3 \leq p(Head) \leq 0.4\}$$

- Since apparently,
$$p^* = \arg\min_{p \in \mathcal{E}} KL(p||(0.5, 0.5)) = (0.4, 0.6)$$

- By conditional limit theorem, knowing the the number of head is within the range, the coin behaves like a biased coin with $p(Head) = 0.4$

- A best bet would be there are 400 heads

## Examples

### Lower bounds

- Let say $x_1, x_2, \cdots, x_N$ are drawn from $q(\cdot)$. And we have $K$ functions $g_1(\cdot), g_2(\cdot), \cdots, g_K(\cdot)$ such that for $k = 1, \cdots, K$,
$$\sum_{i=1}^{N} g_k(x_i) p(x_i) \geq \alpha_k$$

# Examples

## Lower bounds

- Let say $x_1, x_2, \cdots, x_N$ are drawn from $q(\cdot)$. And we have $K$ functions $g_1(\cdot), g_2(\cdot), \cdots, g_K(\cdot)$ such that for $k = 1, \cdots, K$,
$$\sum_{i=1}^{N} g_k(x_i) p(x_i) \geq \alpha_k$$

- Let $\mathcal{E} = \{p : \sum_a p(a) g_k(a) \geq \alpha_k, k = 1, \cdots, K\}$

## Examples

### Lower bounds

- Let say $x_1, x_2, \cdots, x_N$ are drawn from $q(\cdot)$. And we have $K$ functions $g_1(\cdot), g_2(\cdot), \cdots, g_K(\cdot)$ such that for $k = 1, \cdots, K$,
$$\sum_{i=1}^{N} g_k(x_i) p(x_i) \geq \alpha_k$$

- Let $\mathcal{E} = \{p : \sum_a p(a) g_k(a) \geq \alpha_k, k = 1, \cdots, K\}$

- From conditional limit theorem, $\frac{\mathcal{N}(a|x^N)}{N} \to p^*(a)$, where
$$p^* = \arg\min_{p \in \mathcal{E}} KL(p\|q)$$

## Examples

### Lower bounds

- Let say $x_1, x_2, \cdots, x_N$ are drawn from $q(\cdot)$. And we have $K$ functions $g_1(\cdot), g_2(\cdot), \cdots, g_K(\cdot)$ such that for $k = 1, \cdots, K$,
$$\sum_{i=1}^{N} g_k(x_i) p(x_i) \geq \alpha_k$$

- Let $\mathcal{E} = \{ p : \sum_a p(a) g_k(a) \geq \alpha_k, k = 1, \cdots, K \}$

- From conditional limit theorem, $\frac{\mathcal{N}(a|x^N)}{N} \to p^*(a)$, where
$$p^* = \arg \min_{p \in \mathcal{E}} KL(p \| q)$$

- This is a simple constrained optimization problem and can be solved with KKT conditions. If you go through the conditions, you will find that
$$p^*(x) \propto q(x) 2^{\sum_{k=1}^{K} \lambda_k g_k(x)},$$

with $\lambda_k (\sum_a p(a) g_k(a) - \alpha_k) = 0$, $\lambda_k \geq 0$, and $\sum_a p(a) g_k(a) \geq \alpha_k$

## Examples

I think this example below gives a nice demonstration that the technique we have learned today can solve some amazing puzzle!

## Examples

I think this example below gives a nice demonstration that the technique we have learned today can solve some amazing puzzle!

### Fair dice

A fair dice is thrown 10,000 times and the sum of all outcomes is larger than 40,000, out of the 10,000 throw, how many ones do you think there are?

## Fair dice

- From the result of previous example, let $g_1(x) = x$ and $\alpha_1 = 4$, we expect

$$p^*(i) = \frac{2^{\lambda i}}{\sum_{j=1}^{6} 2^{\lambda j}}$$

for some $\lambda$

## Fair dice

- From the result of previous example, let $g_1(x) = x$ and $\alpha_1 = 4$, we expect

$$p^*(i) = \frac{2^{\lambda i}}{\sum_{j=1}^{6} 2^{\lambda j}}$$

for some $\lambda$

- $\lambda \neq 0$ since $\sum_a p(a) g_1(a) = 3.5 < 4 = \alpha_1$ if so

## Fair dice

- From the result of previous example, let $g_1(x) = x$ and $\alpha_1 = 4$, we expect

$$p^*(i) = \frac{2^{\lambda i}}{\sum_{j=1}^{6} 2^{\lambda j}}$$

for some $\lambda$

- $\lambda \neq 0$ since $\sum_a p(a)g_1(a) = 3.5 < 4 = \alpha_1$ if so
- Since $\lambda \neq 0$, by the complementary slackness constraint $\lambda_k(\sum_a p(a)g_k(a) - \alpha_k) = 0$,

$$\sum_a p(a)g_1(a) = \alpha_1 = 4$$

## Fair dice

- From the result of previous example, let $g_1(x) = x$ and $\alpha_1 = 4$, we expect

$$p^*(i) = \frac{2^{\lambda i}}{\sum_{j=1}^6 2^{\lambda j}}$$

for some $\lambda$

- $\lambda \neq 0$ since $\sum_a p(a)g_1(a) = 3.5 < 4 = \alpha_1$ if so
- Since $\lambda \neq 0$, by the complementary slackness constraint $\lambda_k(\sum_a p(a)g_k(a) - \alpha_k) = 0$,

$$\sum_a p(a)g_1(a) = \alpha_1 = 4$$

- This gives us $\lambda = 0.2519$, and thus $p^* = (0.103, 0.123, 0.146, 0.174, 0.207, 0.247)$

## Fair dice

- From the result of previous example, let $g_1(x) = x$ and $\alpha_1 = 4$, we expect

$$p^*(i) = \frac{2^{\lambda i}}{\sum_{j=1}^6 2^{\lambda j}}$$

for some $\lambda$

- $\lambda \neq 0$ since $\sum_a p(a)g_1(a) = 3.5 < 4 = \alpha_1$ if so
- Since $\lambda \neq 0$, by the complementary slackness constraint $\lambda_k(\sum_a p(a)g_k(a) - \alpha_k) = 0$,

$$\sum_a p(a)g_1(a) = \alpha_1 = 4$$

- This gives us $\lambda = 0.2519$, and thus
  $p^* = (0.103, 0.123, 0.146, 0.174, 0.207, 0.247)$
- # ones $\approx 0.103 \times 10000 = 1030$

# Normal distribution

- Univariate Normal: $\mathcal{N}(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
- Multivariate Normal: $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \frac{1}{det(2\pi\Sigma)} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \Sigma^{-1} (\mathbf{x}-\boldsymbol{\mu})}$

### Remark

*Note that $\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \Sigma) = \mathcal{N}(\boldsymbol{\mu}; \mathbf{x}, \Sigma)$. It is trivial but quite useful*

### Remark

*$\Sigma$ is known to be the covariance matrices and it has to be (symmetric) positive definite*

### Remark

*Consequently, symmetric matrices are carefully studied and understood by statisticians and information theorists (more discussion couple slides later)*

# Covariance matrices

### Definition (Covariance matrices)

Recall that for a vector random variable $\boldsymbol{X} = [X_1, X_2, \cdots, X_n]^T$, the covariance matrix $\Sigma \triangleq E[(\boldsymbol{X} - \boldsymbol{\mu})(\boldsymbol{X} - \boldsymbol{\mu})^T]$

### Remark

*Covariance matrices are always positive semi-definite since $\forall u$,*
$u^T \Sigma u = E[u^T (\boldsymbol{X} - \boldsymbol{\mu})(\boldsymbol{X} - \boldsymbol{\mu})^T u] = E[\|(\boldsymbol{X} - \boldsymbol{\mu})^T u\|^2] \geq 0$

### Remark

*In general, we usually would like to assume $\Sigma$ to be strictly positive definite. Because otherwise it means that some of its eigenvalues are zero and so in some dimension, there is actually no variation and is just constant along that dimension. Representing those dimension as random variable is troublesome since "$1/\sigma^2$" which occurs often will become infinite. Instead we can always simply strip away those dimensions to avoid complications*

# Symmetric matrices

### Lemma

$(M^T)^{-1} = (M^{-1})^T$

# Symmetric matrices

### Lemma

$(M^T)^{-1} = (M^{-1})^T$

### Proof.

$(M^{-1})^T M^T = (MM^{-1})^T = I \Rightarrow (M^{-1})^T$ is inverse of $M^T$ $\qquad \square$

### Lemma

If $M$ is symmetric, so is $M^{-1}$

# Symmetric matrices

### Lemma

$(M^T)^{-1} = (M^{-1})^T$

### Proof.

$(M^{-1})^T M^T = (MM^{-1})^T = I \Rightarrow (M^{-1})^T$ is inverse of $M^T$ $\square$

### Lemma

If $M$ is symmetric, so is $M^{-1}$

### Proof.

$(M^{-1})^T = (M^T)^{-1} = M^{-1}$ $\square$

# Hermitian matrices

- An extension of transpose operation to complex matrices is the hermitian transpose operation, which is simply the transpose and conjugate of a matrix (vector)

- We denote the hermitian transpose of $M$ as $M^\dagger \triangleq \overline{M}^T$, when $\overline{M}$ is the complex conjugate of $M$

- A matrix is Hermitian if $M^\dagger = M$. Note that a real symmetric matrix is Hermitian

# Eigenvalues of Hermitian matrices

## Lemma

*If M is Hermitian ($M^{\dagger} = M$), all eigenvalues are real*

# Eigenvalues of Hermitian matrices

### Lemma

*If $M$ is Hermitian ($M^\dagger = M$), all eigenvalues are real*

### Proof.

$$\overline{\lambda}(x^\dagger x) = (\lambda x)^\dagger x = (Mx)^\dagger x = x^\dagger M^\dagger x = x^\dagger M x = x^\dagger(\lambda x) = \lambda(x^\dagger x) \qquad \square$$

### Lemma

*If $M$ is Hermitian, eigenvectors of different eigenvalues are orthogonal*

# Eigenvalues of Hermitian matrices

### Lemma

*If M is Hermitian ($M^\dagger = M$), all eigenvalues are real*

### Proof.

$$\overline{\lambda}(x^\dagger x) = (\lambda x)^\dagger x = (Mx)^\dagger x = x^\dagger M^\dagger x = x^\dagger M x = x^\dagger (\lambda x) = \lambda(x^\dagger x) \qquad \square$$

### Lemma

*If M is Hermitian, eigenvectors of different eigenvalues are orthogonal*

### Proof.

$$\lambda_1 x_1^\dagger x_2 = (Mx_1)^\dagger x_2 = x_1^\dagger M x_2 = \lambda_2 x_1^\dagger x_2$$
$$\Rightarrow \lambda_1 \neq \lambda_2 \Rightarrow x_1^\dagger x_2 = 0$$

$\square$

# Hermitian matrices are diagonizable

### Lemma

*Hermitian matrices are diagonizable*

### Proof (*).

We will sketch the proof by construction. For any $n$-d Hermitian matrix $M$, consider an eigenvalue $\lambda$ and corresponding eigenvector $u$, without loss of generality, let's also normalize $u$ such that $\|u\| = 1$. Consider the subspace orthogonal to $u$, $U^{\perp}$, and let $v_1, \cdots, v_{n-1}$ be arbitrary orthonormal basis of $U^{\perp}$. Note that for any $k$, $Mv_k$ will be orthogonal to $u$ since

$$u^{\dagger} M v_k = u^{\dagger} M^{\dagger} v_k = (Mu)^{\dagger} v_k = \lambda u^{\dagger} v_k = 0.$$

Thus, $\left(u, v_1, \cdots, v_{n-1}\right)^{\dagger} M \left(u, v_1, \cdots, v_{n-1}\right) = \left(\begin{smallmatrix} \lambda & 0 \\ 0 & M' \end{smallmatrix}\right)$. Moreover, $M'$ is also a Hermitian matrix with one less dimension. We can apply the same process on $M'$ and "diagonalize" one more row/column. That is, $\left(\begin{smallmatrix} 1 & 0 \\ 0 & P' \end{smallmatrix}\right)^{\dagger} P^{\dagger} M P \left(\begin{smallmatrix} 1 & 0 \\ 0 & P' \end{smallmatrix}\right) = \left(\begin{smallmatrix} \lambda & 0 & \cdots \\ 0 & \lambda' & \\ & & M'' \end{smallmatrix}\right)$. We can repeat this until the entire $M$ is diagonalized $\qquad\square$

# Hermitian matrices are diagonalizable

### Remark

*We can find a orthogonal set of eigenvectors that diagonalize a Hermitian matrix. That is*

$$\left(v_1, \cdots, v_n\right)^{\dagger} M \underbrace{\left(v_1, \cdots, v_n\right)}_{V} = \begin{pmatrix} \lambda_1 & 0 & \cdots \\ 0 & \lambda_2 & \\ \vdots & & \ddots \end{pmatrix},$$

*and $V$ is unitary (orthogonal), i.e., $V^{\dagger}V = I$ and thus $V^{-1} = V^{\dagger}$. Note that $v_i \perp v_j$ if $\lambda_i \neq \lambda_j$. Otherwise, we may use Gram-Schmidt*

### Remark

*The reverse is obviously true. If a matrix can be diagonalized by a unitary matrix into a real diagonal matrix, the matrix is Hermitian*

### Remark

*Recall that real-symmetric matrices are Hermitian, thus can be diagonalized by its eigenvectors also*

# Positive definite matrices

### Definition (Positive definite)

For a Hermitian matrix $M$, it is positive definite iff $\forall x$, $x^\dagger M x > 0$

### Definition (Positive semi-definite)

For a Hermitian matrix $M$, it is positive semi-definite iff $\forall x$, $x^\dagger M x \geq 0$

### Remark

*M is positive definite (semi-definite) iff all its eigenvalue is larger (larger or equal to) 0*

# Positive definite matrices

### Definition (Positive definite)

For a Hermitian matrix $M$, it is positive definite iff $\forall x$, $x^\dagger M x > 0$

### Definition (Positive semi-definite)

For a Hermitian matrix $M$, it is positive semi-definite iff $\forall x$, $x^\dagger M x \geq 0$

### Remark

*$M$ is positive definite (semi-definite) iff all its eigenvalue is larger (larger or equal to) 0*

### Proof.

$\Rightarrow$: assume positive definite but some eigenvalue $< 0$, WLOG, let $\lambda_1 < 0$, then $v_1^\dagger M v_1 = \lambda_1 < 0$ contradicts that $M$ is positive definite

$\Leftarrow$: If $\forall k, \lambda_k > 0$, for any $x$,

$x^\dagger M x = (V^\dagger x)^\dagger \begin{pmatrix} \lambda_1 & 0 \\ & \ddots \\ 0 & \end{pmatrix} V^\dagger x = \sum_i \lambda_i (V^\dagger x)_i^2 > 0$ $\qquad \square$

# Eigenvectors and eigenvalues of covariance matrices

- WLOG, let's assume $\mathbf{X} = [X_1, X_2, \cdots, X_n]^T$ is zero mean. So the covariance matrix $\Sigma_X = E[\mathbf{X}\mathbf{X}^T]$

# Eigenvectors and eigenvalues of covariance matrices

- WLOG, let's assume $\mathbf{X} = [X_1, X_2, \cdots, X_n]^T$ is zero mean. So the covariance matrix $\Sigma_X = E[\mathbf{X}\mathbf{X}^T]$
- Covariance matrices are real symmetric (hence Hermitian) and so can be diagonalized by its eigenvectors. That is,
  - $P^T \Sigma_X P = D$, where $P = [u_1, u_2, \cdots, u_n]$ with $u_k$ being eigenvectors of $\Sigma$ and $D$ is a diagonal matrix with eigenvalues $\lambda_1, \lambda_2, \cdots, \lambda_n$ as the diagonal elements

## Eigenvectors and eigenvalues of covariance matrices

- WLOG, let's assume $\mathbf{X} = [X_1, X_2, \cdots, X_n]^T$ is zero mean. So the covariance matrix $\Sigma_X = E[\mathbf{X}\mathbf{X}^T]$
- Covariance matrices are real symmetric (hence Hermitian) and so can be diagonalized by its eigenvectors. That is,
  - $P^T \Sigma_X P = D$, where $P = [u_1, u_2, \cdots, u_n]$ with $u_k$ being eigenvectors of $\Sigma$ and $D$ is a diagonal matrix with eigenvalues $\lambda_1, \lambda_2, \cdots, \lambda_n$ as the diagonal elements
- Let $\mathbf{Y} = P^T \mathbf{X}$, note that the covariance matrix of $\mathbf{Y}$

$$\Sigma_Y = E[\mathbf{Y}\mathbf{Y}^T] = E[P^T \mathbf{X}\mathbf{X}^T P] = P^T E[\mathbf{X}\mathbf{X}^T]P = P^T \Sigma_X P = D$$

is diagonalized

## Eigenvectors and eigenvalues of covariance matrices

- WLOG, let's assume $\mathbf{X} = [X_1, X_2, \cdots, X_n]^T$ is zero mean. So the covariance matrix $\Sigma_X = E[\mathbf{X}\mathbf{X}^T]$
- Covariance matrices are real symmetric (hence Hermitian) and so can be diagonalized by its eigenvectors. That is,
  - $P^T \Sigma_X P = D$, where $P = [u_1, u_2, \cdots, u_n]$ with $u_k$ being eigenvectors of $\Sigma$ and $D$ is a diagonal matrix with eigenvalues $\lambda_1, \lambda_2, \cdots, \lambda_n$ as the diagonal elements
- Let $\mathbf{Y} = P^T \mathbf{X}$, note that the covariance matrix of $\mathbf{Y}$

$$\Sigma_Y = E[\mathbf{Y}\mathbf{Y}^T] = E[P^T\mathbf{X}\mathbf{X}^T P] = P^T E[\mathbf{X}\mathbf{X}^T]P = P^T \Sigma_X P = D$$

  is diagonalized

  - So the variance of $Y_k$ is simply $\lambda_k$
  - $E[Y_i Y_j] = 0$ for $i \neq j$. That is, $Y_i \perp\!\!\!\perp Y_j$ for $i \neq j$
  - Note that $\mathbf{Y} = P^T\mathbf{X}$ is just principal component analysis (PCA)

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$

---

[1]$tr(AB) = \sum_i \sum_j a_{i,j} b_{j,i} = \sum_j \sum_i b_{j,i} a_{i,j} = tr(BA)$

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0

---

[1]$tr(AB) = \sum_i \sum_j a_{i,j} b_{j,i} = \sum_j \sum_i b_{j,i} a_{i,j} = tr(BA)$

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
  - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$

---

[1] $tr(AB) = \sum_i \sum_j a_{i,j} b_{j,i} = \sum_j \sum_i b_{j,i} a_{i,j} = tr(BA)$

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T \mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T \Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
    - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
    - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
      $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})])$

---

[1] $tr(AB) = \sum_i \sum_j a_{i,j} b_{j,i} = \sum_j \sum_i b_{j,i} a_{i,j} = tr(BA)$

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$

- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
  - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
    $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[tr((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$

---

[1]$tr(AB) = \sum_i \sum_j a_{i,j}b_{j,i} = \sum_j \sum_i b_{j,i}a_{i,j} = tr(BA)$

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
  - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
    $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[tr((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
    $= E[tr((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T])$

---

[1]$tr(AB) = \sum_i \sum_j a_{i,j} b_{j,i} = \sum_j \sum_i b_{j,i} a_{i,j} = tr(BA)$

## Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
    - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
    - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
    $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[tr((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
    $= E[tr((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^{n} \lambda_i$
    - Similarly, if we "reconstruct" $\mathbf{X}$ as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of
    $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = tr(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T])$

---

[1] $tr(AB) = \sum_i \sum_j a_{i,j}b_{j,i} = \sum_j \sum_i b_{j,i}a_{i,j} = tr(BA)$

## Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$

- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$

  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0

  - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
    $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[tr((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
    $= E[tr((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^{n} \lambda_i$

  - Similarly, if we "reconstruct" $\mathbf{X}$ as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of
    $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = tr(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) = tr(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})]P^T)$

---

[1] $tr(AB) = \sum_i \sum_j a_{i,j}b_{j,i} = \sum_j \sum_i b_{j,i}a_{i,j} = tr(BA)$

## Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
  - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
    $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[tr((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
    $= E[tr((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^{n} \lambda_i$
  - Similarly, if we "reconstruct" $\mathbf{X}$ as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of
    $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = tr(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) =$
    $tr(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})]P^T) = tr(P^T PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})])$

---

[1] $tr(AB) = \sum_i \sum_j a_{i,j} b_{j,i} = \sum_j \sum_i b_{j,i} a_{i,j} = tr(BA)$

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
  - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
    $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[tr((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
    $= E[tr((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^{n} \lambda_i$
  - Similarly, if we "reconstruct" $\mathbf{X}$ as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of
    $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = tr(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) =$
    $tr(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})]P^T) = tr(P^T P E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})]) = \sum_{i=k+1}^{n} \lambda_i$

---

[1] $tr(AB) = \sum_i \sum_j a_{i,j}b_{j,i} = \sum_j \sum_i b_{j,i}a_{i,j} = tr(BA)$

# Principal component analysis (PCA)

- Recall that $\Sigma = E[\mathbf{X}\mathbf{X}^T]$ (assume $\mathbf{X}$ is zero-mean) and $\mathbf{Y} = P^T\mathbf{X}$ with $E[\mathbf{Y}\mathbf{Y}^T] = P^T\Sigma P = D$
- Assume that the diagonal of $D$ (note that those are eigenvalues) are arranged in descending order that $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_n$
  - Generate an approximate $\hat{\mathbf{Y}}$ of $\mathbf{Y}$ by setting all components except first $k$ as 0
  - The mean square error (mse) of[1] $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
    $= tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[tr((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
    $= E[tr((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = tr(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^{n} \lambda_i$
  - Similarly, if we "reconstruct" $\mathbf{X}$ as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of
    $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = tr(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) =$
    $tr(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})]P^T) = tr(P^TPE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})]) = \sum_{i=k+1}^{n} \lambda_i$
  - Note that the eigenvectors of $\Sigma$ (columns of $P$) are known as the principal components

---

[1] $tr(AB) = \sum_i \sum_j a_{i,j}b_{j,i} = \sum_j \sum_i b_{j,i}a_{i,j} = tr(BA)$

## Practical PCA

In practice, we typically are given a dataset with samples of **X** instead of the distribution or covariance matrix of **X**. Denote the data as $\mathcal{X}$ with each row is a data point and a total of $m$ data points. Thus $\mathcal{X}$ is an $m$ by $n$ matrix

---

[2]I used the matlab notations for *ones*($\cdot$) and *mean*($\cdot$) here

[3]Note that $\hat{\Sigma}$ won't be full rank and positive definite as one would hope

## Practical PCA

In practice, we typically are given a dataset with samples of **X** instead of the distribution or covariance matrix of **X**. Denote the data as $\mathcal{X}$ with each row is a data point and a total of $m$ data points. Thus $\mathcal{X}$ is an $m$ by $n$ matrix

- Data are rarely zero-mean to begin with, but we can easily preprocess it by subtracting the mean. That is[2] $\mathcal{X} \leftarrow \mathcal{X} - ones(m, 1)mean(\mathcal{X})$

---

[2]I used the matlab notations for $ones(\cdot)$ and $mean(\cdot)$ here

[3]Note that $\hat{\Sigma}$ won't be full rank and positive definite as one would hope

## Practical PCA

In practice, we typically are given a dataset with samples of **X** instead of the distribution or covariance matrix of **X**. Denote the data as $\mathcal{X}$ with each row is a data point and a total of $m$ data points. Thus $\mathcal{X}$ is an $m$ by $n$ matrix

- Data are rarely zero-mean to begin with, but we can easily preprocess it by subtracting the mean. That is[2] $\mathcal{X} \leftarrow \mathcal{X} - ones(m, 1)mean(\mathcal{X})$
- Note that $\hat{\Sigma} \approx \frac{1}{m}\mathcal{X}^T\mathcal{X}$. We could directly compute the eigenvectors and eigenvalues of $\hat{\Sigma}$ as discussed previously. But in many cases, $m < n$ making $\hat{\Sigma}$ a bad approximate[3]

---

[2]I used the matlab notations for $ones(\cdot)$ and $mean(\cdot)$ here

[3]Note that $\hat{\Sigma}$ won't be full rank and positive definite as one would hope

## Practical PCA

In practice, we typically are given a dataset with samples of **X** instead of the distribution or covariance matrix of **X**. Denote the data as $\mathcal{X}$ with each row is a data point and a total of $m$ data points. Thus $\mathcal{X}$ is an $m$ by $n$ matrix

- Data are rarely zero-mean to begin with, but we can easily preprocess it by subtracting the mean. That is[2] $\mathcal{X} \leftarrow \mathcal{X} - ones(m, 1)mean(\mathcal{X})$
- Note that $\hat{\Sigma} \approx \frac{1}{m}\mathcal{X}^T\mathcal{X}$. We could directly compute the eigenvectors and eigenvalues of $\hat{\Sigma}$ as discussed previously. But in many cases, $m < n$ making $\hat{\Sigma}$ a bad approximate[3]
    - A more common approach is to decompose $\mathcal{X}$ with singular value decomposition (SVD) instead

---

[2]I used the matlab notations for $ones(\cdot)$ and $mean(\cdot)$ here

[3]Note that $\hat{\Sigma}$ won't be full rank and positive definite as one would hope

# Singular value decomposition (SVD)



$M = U \cdot \Sigma \cdot V^*$

- Every matrix $M$ can be decomposed as $M = UDV^\dagger$, where $D$ is diagonal and $U, V$ are unitary. The diagonal terms in $\Sigma$ are known to be the <span style="color:red">singular values</span>

# Singular value decomposition (SVD)



$$M = U \cdot \Sigma \cdot V^*$$

- Every matrix $M$ can be decomposed as $M = UDV^\dagger$, where $D$ is diagonal and $U, V$ are unitary. The diagonal terms in $\Sigma$ are known to be the <span style="color:red">singular values</span>
- For real matrix $M$, we can write $M = UDV^T$ instead. $U, V$ are now "real unitary" or orthogonal

# Singular value decomposition (SVD)



$$M = U \cdot \Sigma \cdot V^*$$

- Every matrix $M$ can be decomposed as $M = UDV^\dagger$, where $D$ is diagonal and $U, V$ are unitary. The diagonal terms in $\Sigma$ are known to be the <span style="color:red">singular values</span>

- For real matrix $M$, we can write $M = UDV^T$ instead. $U, V$ are now "real unitary" or orthogonal

  - Note that $M^T M = VD^T U^T UDV^T = VD^2 V^T$. Therefore, $V$ are really eigenvectors of $M^T M$ with eigenvalues equal to the square of the singular values

# Singular value decomposition (SVD)



$$M = U \cdot \Sigma \cdot V^*$$

- Every matrix $M$ can be decomposed as $M = UDV^\dagger$, where $D$ is diagonal and $U, V$ are unitary. The diagonal terms in $\Sigma$ are known to be the <span style="color:red">singular values</span>

- For real matrix $M$, we can write $M = UDV^T$ instead. $U, V$ are now "real unitary" or orthogonal
  - Note that $M^T M = VD^T U^T UDV^T = VD^2 V^T$. Therefore, $V$ are really eigenvectors of $M^T M$ with eigenvalues equal to the square of the singular values
  - Similar, we have $MM^T = UD^2 U^T$

## PCA with SVD

So from previous slides, instead of first estimating the covariance matrix and then diagonalize it. We should directly decompose the data $\mathcal{X}$ with SVD instead. The process is summarized below

## PCA with SVD

So from previous slides, instead of first estimating the covariance matrix and then diagonalize it. We should directly decompose the data $\mathcal{X}$ with SVD instead. The process is summarized below

- Estimate mean from data and subtract mean from that

## PCA with SVD

So from previous slides, instead of first estimating the covariance matrix and then diagonalize it. We should directly decompose the data $\mathcal{X}$ with SVD instead. The process is summarized below

- Estimate mean from data and subtract mean from that
- Decomposed the mean subtracted data with SVD. We get $\mathcal{X} = UDV^T$

## PCA with SVD

So from previous slides, instead of first estimating the covariance matrix and then diagonalize it. We should directly decompose the data $\mathcal{X}$ with SVD instead. The process is summarized below

- Estimate mean from data and subtract mean from that
- Decomposed the mean subtracted data with SVD. We get $\mathcal{X} = UDV^T$
- Note that column of $V$ are now the principal components, and we can transform a data column as $V^T x$. The entire data set can be transformed as $\mathcal{Y} = \mathcal{X}V$

## PCA with SVD

So from previous slides, instead of first estimating the covariance matrix and then diagonalize it. We should directly decompose the data $\mathcal{X}$ with SVD instead. The process is summarized below

- Estimate mean from data and subtract mean from that
- Decomposed the mean subtracted data with SVD. We get $\mathcal{X} = UDV^T$
- Note that column of $V$ are now the principal components, and we can transform a data column as $V^T x$. The entire data set can be transformed as $\mathcal{Y} = \mathcal{X}V$
    - The first few columns of $\mathcal{Y}$ will contain most "information" regarding the original $\mathcal{X}$

## PCA with SVD

So from previous slides, instead of first estimating the covariance matrix and then diagonalize it. We should directly decompose the data $\mathcal{X}$ with SVD instead. The process is summarized below

- Estimate mean from data and subtract mean from that
- Decomposed the mean subtracted data with SVD. We get $\mathcal{X} = UDV^T$
- Note that column of $V$ are now the principal components, and we can transform a data column as $V^T x$. The entire data set can be transformed as $\mathcal{Y} = \mathcal{X}V$
    - The first few columns of $\mathcal{Y}$ will contain most "information" regarding the original $\mathcal{X}$
    - For example, they can be taken as features for recognition or one can omit other columns besides the first few for "compression" as discussed earlier

# Marginalization of normal distribution

- Consider $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{Z}}, \Sigma_{\mathbf{Z}})$ and let say $\mathbf{X}$ is a segment of $\mathbf{Z}$. That is, $\mathbf{Z} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$ for some $\mathbf{Y}$. Then how should $\mathbf{X}$ behave?

## Marginalization of normal distribution

- Consider $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu_Z}, \Sigma_{\mathbf{Z}})$ and let say $\mathbf{X}$ is a segment of $\mathbf{Z}$. That is, $\mathbf{Z} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$ for some $\mathbf{Y}$. Then how should $\mathbf{X}$ behave?

- We can find the pdf of $\mathbf{X}$ by just marginalizing that of $\mathbf{Z}$. That is

$$
\begin{aligned}
p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{y}) d\mathbf{y} \\
&= \frac{1}{\sqrt{\det(2\pi\Sigma)}} \int \exp\left( -\frac{1}{2} \begin{pmatrix} \mathbf{x} - \boldsymbol{\mu_X} \\ \mathbf{y} - \boldsymbol{\mu_Y} \end{pmatrix}^T \Sigma^{-1} \begin{pmatrix} \mathbf{x} - \boldsymbol{\mu_X} \\ \mathbf{y} - \boldsymbol{\mu_Y} \end{pmatrix} \right) d\mathbf{y}
\end{aligned}
$$

# Marginalization of normal distribution

- Denote $\Sigma^{-1}$ as $\Lambda$ (also known as the precision matrix). And partition both $\Sigma$ and $\Lambda$ into $\Sigma = \begin{pmatrix} \Sigma_{\mathbf{XX}} & \Sigma_{\mathbf{XY}} \\ \Sigma_{\mathbf{YX}} & \Sigma_{\mathbf{YY}} \end{pmatrix}$ and $\Lambda = \begin{pmatrix} \Lambda_{\mathbf{XX}} & \Lambda_{\mathbf{XY}} \\ \Lambda_{\mathbf{YX}} & \Lambda_{\mathbf{YY}} \end{pmatrix}$

## Marginalization of normal distribution

- Denote $\Sigma^{-1}$ as $\Lambda$ (also known as the precision matrix). And partition both $\Sigma$ and $\Lambda$ into $\Sigma = \begin{pmatrix} \Sigma_{\mathbf{XX}} & \Sigma_{\mathbf{XY}} \\ \Sigma_{\mathbf{YX}} & \Sigma_{\mathbf{YY}} \end{pmatrix}$ and $\Lambda = \begin{pmatrix} \Lambda_{\mathbf{XX}} & \Lambda_{\mathbf{XY}} \\ \Lambda_{\mathbf{YX}} & \Lambda_{\mathbf{YY}} \end{pmatrix}$

- Then we have

$$
\begin{aligned}
p(\mathbf{x}) &= \frac{1}{\sqrt{\det(2\pi\Sigma)}} \int \exp\left( -\frac{1}{2} \left[ (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}})^T \Lambda_{\mathbf{XX}} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}}) \right.\right. \\
&\quad + (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}})^T \Lambda_{\mathbf{YX}} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}}) + (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}})^T \Lambda_{\mathbf{XY}} (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}) \\
&\quad \left.\left. + (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}})^T \Lambda_{\mathbf{YY}} (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}) \right] \right) d\mathbf{y} \\
&= \frac{e^{-\frac{(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}})^T \Lambda_{\mathbf{XX}} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}})}{2}}}{\sqrt{\det(2\pi\Sigma)}} \int \exp\left( -\frac{1}{2} \left[ (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}})^T \Lambda_{\mathbf{YX}} (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}}) \right.\right. \\
&\quad \left.\left. + (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}})^T \Lambda_{\mathbf{XY}} (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}) + (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}})^T \Lambda_{\mathbf{YY}} (\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}) \right] \right) d\mathbf{y}
\end{aligned}
$$

## Marginalization of normal distribution

To proceed, let's apply the completing square trick on
$(\mathbf{y} - \boldsymbol{\mu_Y})^T \Lambda_{\mathbf{YX}}(\mathbf{x} - \boldsymbol{\mu_X}) + (\mathbf{x} - \boldsymbol{\mu_X})^T \Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu_Y}) + (\mathbf{y} - \boldsymbol{\mu_Y})^T \Lambda_{\mathbf{YY}}(\mathbf{y} - \boldsymbol{\mu_Y})$.
For the ease of exposition, let us denote $\tilde{\mathbf{x}}$ as $\mathbf{x} - \boldsymbol{\mu_X}$ and $\tilde{\mathbf{y}}$ as $\mathbf{y} - \boldsymbol{\mu_Y}$. We
have

## Marginalization of normal distribution

To proceed, let's apply the completing square trick on
$(\mathbf{y} - \boldsymbol{\mu_Y})^T \Lambda_{\mathbf{YX}} (\mathbf{x} - \boldsymbol{\mu_X}) + (\mathbf{x} - \boldsymbol{\mu_X})^T \Lambda_{\mathbf{XY}} (\mathbf{y} - \boldsymbol{\mu_Y}) + (\mathbf{y} - \boldsymbol{\mu_Y})^T \Lambda_{\mathbf{YY}} (\mathbf{y} - \boldsymbol{\mu_Y})$.
For the ease of exposition, let us denote $\tilde{\mathbf{x}}$ as $\mathbf{x} - \boldsymbol{\mu_X}$ and $\tilde{\mathbf{y}}$ as $\mathbf{y} - \boldsymbol{\mu_Y}$. We
have

$$
\begin{aligned}
&\tilde{\mathbf{y}}^T \Lambda_{\mathbf{YX}} \tilde{\mathbf{x}} + \tilde{\mathbf{x}}^T \Lambda_{\mathbf{XY}} \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \Lambda_{\mathbf{YY}} \tilde{\mathbf{y}} \\
=&(\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1} \Lambda_{\mathbf{YX}} \tilde{\mathbf{x}})^T \Lambda_{\mathbf{YY}} (\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1} \Lambda_{\mathbf{YX}} \tilde{\mathbf{x}}) - \tilde{\mathbf{x}}^T \Lambda_{\mathbf{XY}} \Lambda_{\mathbf{YY}}^{-1} \Lambda_{\mathbf{YX}} \tilde{\mathbf{x}},
\end{aligned}
$$

where we use the fact that $\Lambda = \Sigma^{-1}$ is symmetric and so $\Lambda_{\mathbf{XY}} = \Lambda_{\mathbf{YX}}$

## Marginalization of normal distribution

$$p(\mathbf{x}) = \frac{e^{-\frac{\bar{\mathbf{x}}^T(\Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\bar{\mathbf{x}}}{2}}}{\sqrt{\det(2\pi\Sigma)}} \int e^{-\frac{(\bar{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\bar{\mathbf{x}})^T\Lambda_{\mathbf{YY}}(\bar{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\bar{\mathbf{x}})}{2}} d\mathbf{y}$$

# Marginalization of normal distribution

$$p(\mathbf{x}) = \frac{e^{-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}}}{\sqrt{\det(2\pi\Sigma)}} \int e^{-\frac{(\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})^T \Lambda_{\mathbf{YY}}(\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})}{2}} \, d\mathbf{y}$$

$$= \frac{\sqrt{\det(2\pi\Lambda_{\mathbf{YY}}^{-1})}}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}\right)$$

## Marginalization of normal distribution

$$p(\mathbf{x}) = \frac{e^{-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}}}{\sqrt{\det(2\pi\Sigma)}} \int e^{-\frac{(\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})^T \Lambda_{\mathbf{YY}}(\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})}{2}} d\mathbf{y}$$

$$= \frac{\sqrt{\det(2\pi\Lambda_{\mathbf{YY}}^{-1})}}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}\right)$$

$$\stackrel{(a)}{=} \frac{\sqrt{\det(2\pi\Lambda_{\mathbf{YY}}^{-1})}}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{\tilde{\mathbf{x}}^T\Sigma_{\mathbf{XX}}^{-1}\tilde{\mathbf{x}}}{2}\right)$$

# Marginalization of normal distribution

$$
\begin{aligned}
p(\mathbf{x}) &= \frac{e^{-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}}-\Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}}}{\sqrt{\det(2\pi\Sigma)}} \int e^{-\frac{(\tilde{\mathbf{y}}+\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})^T\Lambda_{\mathbf{YY}}(\tilde{\mathbf{y}}+\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})}{2}} \, d\mathbf{y} \\
&= \frac{\sqrt{\det(2\pi\Lambda_{\mathbf{YY}}^{-1})}}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}}-\Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}\right) \\
&\overset{(a)}{=} \frac{\sqrt{\det(2\pi\Lambda_{\mathbf{YY}}^{-1})}}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{\tilde{\mathbf{x}}^T\Sigma_{\mathbf{XX}}^{-1}\tilde{\mathbf{x}}}{2}\right) \\
&\overset{(b)}{=} \frac{1}{\sqrt{\det(2\pi\Sigma_{\mathbf{XX}})}} \exp\left(-\frac{\tilde{\mathbf{x}}^T\Sigma_{\mathbf{XX}}^{-1}\tilde{\mathbf{x}}}{2}\right)
\end{aligned}
$$

# Marginalization of normal distribution

$$p(\mathbf{x}) = \frac{e^{-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}}}{\sqrt{\det(2\pi\Sigma)}} \int e^{-\frac{(\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})^T\Lambda_{\mathbf{YY}}(\tilde{\mathbf{y}} + \Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}\tilde{\mathbf{x}})}{2}} d\mathbf{y}$$

$$= \frac{\sqrt{\det(2\pi\Lambda_{\mathbf{YY}}^{-1})}}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{\tilde{\mathbf{x}}^T(\Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}})\tilde{\mathbf{x}}}{2}\right)$$

$$\stackrel{(a)}{=} \frac{\sqrt{\det(2\pi\Lambda_{\mathbf{YY}}^{-1})}}{\sqrt{\det(2\pi\Sigma)}} \exp\left(-\frac{\tilde{\mathbf{x}}^T\Sigma_{\mathbf{XX}}^{-1}\tilde{\mathbf{x}}}{2}\right)$$

$$\stackrel{(b)}{=} \frac{1}{\sqrt{\det(2\pi\Sigma_{\mathbf{XX}})}} \exp\left(-\frac{\tilde{\mathbf{x}}^T\Sigma_{\mathbf{XX}}^{-1}\tilde{\mathbf{x}}}{2}\right)$$

$$= \frac{1}{\sqrt{\det(2\pi\Sigma_{\mathbf{XX}})}} \exp\left(-\frac{(\mathbf{x} - \mu_{\mathbf{X}})^T\Sigma_{\mathbf{XX}}^{-1}(\mathbf{x} - \mu_{\mathbf{X}})}{2}\right),$$

where (a) and (b) will be shown next

# (a) $\Sigma_{\mathbf{XX}}^{-1} = \Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}$

## Lemma

Assume $\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} \tilde{A} & \tilde{B} \\ \tilde{C} & \tilde{D} \end{pmatrix}$, then $A^{-1} = \tilde{A} - \tilde{B}\tilde{D}^{-1}\tilde{C}$

## Proof.

Note that $\begin{pmatrix} A & B \\ C & D \end{pmatrix}\begin{pmatrix} \tilde{A} & \tilde{B} \\ \tilde{C} & \tilde{D} \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & I \end{pmatrix}$. Thus $A\tilde{A} + B\tilde{C} = I$ and

$A\tilde{B} + B\tilde{D} = 0$. So

$A(\tilde{A} - \tilde{B}\tilde{D}^{-1}\tilde{C}) = A\tilde{A} - (A\tilde{B})\tilde{D}^{-1}\tilde{C} = A\tilde{A} + B\tilde{D}\tilde{D}^{-1}\tilde{C} = A\tilde{A} + B\tilde{C} = I$ $\quad\square$

# (b) $\det(a\Sigma) = \det(a\Sigma_{\mathbf{YY}})\det(a\Lambda_{\mathbf{XX}}^{-1})$

## Lemma

Assume $\begin{pmatrix} A & B \\ C & D \end{pmatrix}^{-1} = \begin{pmatrix} \tilde{A} & \tilde{B} \\ \tilde{C} & \tilde{D} \end{pmatrix}$, then $det\begin{pmatrix} A & B \\ C & D \end{pmatrix} = det(D)det(\tilde{A}^{-1})$

## Proof.

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} = \begin{pmatrix} I & 0 \\ 0 & D \end{pmatrix} \begin{pmatrix} A & B \\ D^{-1}C & I \end{pmatrix}$$

$$= \begin{pmatrix} I & 0 \\ 0 & D \end{pmatrix} \begin{pmatrix} I & B \\ 0 & I \end{pmatrix} \begin{pmatrix} A - BD^{-1}C & 0 \\ D^{-1}C & I \end{pmatrix}$$

$$\Rightarrow det\begin{pmatrix} A & B \\ C & D \end{pmatrix} = det(D)det(A - BD^{-1}C) = det(D)det(\tilde{A}^{-1}) \qquad \square$$

# Review

- ML: $\hat{x} = \arg\max_x p(x|\hat{\theta}), \hat{\theta} = \arg\max_\theta p(o|\theta)$
- MAP: $\hat{x} = \arg\max_x p(x|\hat{\theta}), \hat{\theta} = \arg\max_\theta p(\theta|o)$
- Bayesian: $\hat{x} = \sum_\theta p(\theta|o) \sum_x x p(x|\theta)$
- For zero-mean $\mathbf{X}$, $\Sigma_X = E[\mathbf{X}\mathbf{X}^T]$ and say we have $P^T \Sigma_X P = D$. The transformed $\mathbf{Y} = P^T \mathbf{X}$ are independent to each other
  - Note that the transform is just principal component analysis
- Marginalization of a normal distribution is still a normal distribution
- (a) $\Sigma_{\mathbf{XX}}^{-1} = \Lambda_{\mathbf{XX}} - \Lambda_{\mathbf{XY}}\Lambda_{\mathbf{YY}}^{-1}\Lambda_{\mathbf{YX}}$
- (b) $\det(a\Sigma) = \det(a\Sigma_{\mathbf{YY}})\det(a\Lambda_{\mathbf{XX}}^{-1})$ for any constant $a$

## Conditioning of normal distribution

- Consider the same $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{Z}}, \Sigma_{\mathbf{Z}})$ and $\mathbf{Z} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$. What will $\mathbf{X}$ be like if $\mathbf{Y}$ is observed to be $\mathbf{y}$?

## Conditioning of normal distribution

- Consider the same $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{Z}}, \Sigma_{\mathbf{Z}})$ and $\mathbf{Z} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$. What will $\mathbf{X}$ be like if $\mathbf{Y}$ is observed to be $\mathbf{y}$?

- Basically, we want to find $p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}, \mathbf{y})/p(\mathbf{y})$

## Conditioning of normal distribution

- Consider the same $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu_Z}, \Sigma_{\mathbf{Z}})$ and $\mathbf{Z} = \begin{pmatrix} \mathbf{X} \\ \mathbf{Y} \end{pmatrix}$. What will $\mathbf{X}$ be like if $\mathbf{Y}$ is observed to be $\mathbf{y}$?

- Basically, we want to find $p(\mathbf{x}|\mathbf{y}) = p(\mathbf{x}, \mathbf{y})/p(\mathbf{y})$

- From previous result, we have $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu_Y}, \Sigma_{\mathbf{YY}})$. Therefore,

$$p(\mathbf{x}|\mathbf{y}) \propto \exp\left(-\frac{1}{2}\left[\begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \end{pmatrix}^T \Sigma^{-1} \begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \end{pmatrix} - \tilde{\mathbf{y}}^T \Sigma_{\mathbf{YY}}^{-1} \tilde{\mathbf{y}}\right]\right)$$

$$\propto \exp\left(-\frac{1}{2}[\tilde{\mathbf{x}}^T \Lambda_{\mathbf{XX}} \tilde{\mathbf{x}} + \tilde{\mathbf{x}}^T \Lambda_{\mathbf{XY}} \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \Lambda_{\mathbf{YX}} \tilde{\mathbf{x}}]\right),$$

where we use $\tilde{\mathbf{x}}$ and $\tilde{\mathbf{y}}$ as shorthands of $\mathbf{x} - \boldsymbol{\mu_X}$ and $\mathbf{y} - \boldsymbol{\mu_Y}$ as before

# Conditioning of normal distribution

- Completing the square for $\tilde{\mathbf{x}}$, we have

$$p(\mathbf{x}|\mathbf{y}) \propto \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})^T\Lambda_{\mathbf{XX}}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})\right)$$

$$= \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))^T\Lambda_{\mathbf{XX}}\right.$$

$$\left.(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))\right)$$

## Conditioning of normal distribution

- Completing the square for $\tilde{\mathbf{x}}$, we have

$$p(\mathbf{x}|\mathbf{y}) \propto \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})^T\Lambda_{\mathbf{XX}}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})\right)$$

$$= \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \mu_{\mathbf{Y}}))^T\Lambda_{\mathbf{XX}}\right.$$
$$\left.(\mathbf{x} - \mu_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \mu_{\mathbf{Y}}))\right)$$

- Therefore $\mathbf{X}|\mathbf{y}$ is Gaussian distributed with mean $\mu_{\mathbf{X}} - \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \mu_{\mathbf{Y}})$ and covariance $\Lambda_{\mathbf{XX}}^{-1}$

## Conditioning of normal distribution

- Completing the square for $\tilde{\mathbf{x}}$, we have

$$p(\mathbf{x}|\mathbf{y}) \propto \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})^T\Lambda_{\mathbf{XX}}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})\right)$$

$$= \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \mu_{\mathbf{Y}}))^T\Lambda_{\mathbf{XX}}\right.$$

$$\left.(\mathbf{x} - \mu_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \mu_{\mathbf{Y}})))\right)$$

- Therefore $\mathbf{X}|\mathbf{y}$ is Gaussian distributed with mean $\mu_{\mathbf{X}} - \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \mu_{\mathbf{Y}})$ and covariance $\Lambda_{\mathbf{XX}}^{-1}$
- Note that since $\Lambda_{\mathbf{XX}}\Sigma_{\mathbf{XY}} + \Lambda_{\mathbf{XY}}\Sigma_{\mathbf{YY}} = 0 \Rightarrow \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}} = -\Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}$ and from (a), we have

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\mu_{\mathbf{X}} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \mu_{\mathbf{Y}}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$$

## Interpretation of conditioning

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{X}} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$$

- When the observation of $\mathbf{Y}$ is exactly the mean, the conditioned mean does not change

## Interpretation of conditioning

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu_X} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \boldsymbol{\mu_Y}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$$

- When the observation of $\mathbf{Y}$ is exactly the mean, the conditioned mean does not change
- Otherwise, it needs to be modified and the size of the adjustment decreases with $\Sigma_{\mathbf{YY}}$, the variance of $\mathbf{Y}$ for the 1-D case.
  - The observation is less reliable with the increase of $\Sigma_{\mathbf{YY}}$. The adjustment is finally scaled by $\Sigma_{\mathbf{XY}}$, which translates the variation of $\mathbf{Y}$ to the variation of $\mathbf{X}$

## Interpretation of conditioning

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu_X} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \boldsymbol{\mu_Y}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$$

- When the observation of **Y** is exactly the mean, the conditioned mean does not change
- Otherwise, it needs to be modified and the size of the adjustment decreases with $\Sigma_{\mathbf{YY}}$, the variance of **Y** for the 1-D case.
    - The observation is less reliable with the increase of $\Sigma_{\mathbf{YY}}$. The adjustment is finally scaled by $\Sigma_{\mathbf{XY}}$, which translates the variation of **Y** to the variation of **X**
    - In particular, if **X** and **Y** are negatively correlated, the sign of the adjustment will be reversed

## Interpretation of conditioning

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu_X} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \boldsymbol{\mu_Y}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$$

- When the observation of $\mathbf{Y}$ is exactly the mean, the conditioned mean does not change
- Otherwise, it needs to be modified and the size of the adjustment decreases with $\Sigma_{\mathbf{YY}}$, the variance of $\mathbf{Y}$ for the 1-D case.
  - The observation is less reliable with the increase of $\Sigma_{\mathbf{YY}}$. The adjustment is finally scaled by $\Sigma_{\mathbf{XY}}$, which translates the variation of $\mathbf{Y}$ to the variation of $\mathbf{X}$
  - In particular, if $\mathbf{X}$ and $\mathbf{Y}$ are negatively correlated, the sign of the adjustment will be reversed
- As for the variance of the conditioned variable, it always decreases and the decrease is larger if $\Sigma_{\mathbf{YY}}$ is smaller and $\Sigma_{\mathbf{XY}}$ is larger ($\mathbf{X}$ and $\mathbf{Y}$ are more correlated)

## Uncorrelated implies independence

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu_X} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \boldsymbol{\mu_Y}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$$

If **X** and **Y** are uncorrelated, $\Sigma_{\mathbf{XY}} = 0$. Then

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu_X}, \Sigma_{\mathbf{XX}})$$

Note that the statistics of **X** does not change with respect to **y** and so **X** is independent of **Y**

# $X \perp\!\!\!\perp Y | Z$ if $\rho_{XZ}\rho_{YZ} = \rho_{XY}$

### Corollary

*Given multivariate Gaussian variables $X$, $Y$ and $Z$, we have $X$ and $Y$ are conditionally independent given $Z$ if $\rho_{XZ}\rho_{YZ} = \rho_{XY}$, where $\rho_{XZ} = \frac{E[(X-E(X))(Z-E(Z))]}{\sqrt{E[(X-E(X))^2]E[(Z-E(Z))^2]}}$ is the correlation coefficent between $X$ and $Z$. Similarly, $\rho_{YZ}$ and $\rho_{XY}$ are the correlation coefficients between $Y$ and $Z$, and $X$ and $Y$, respectively.*

# $X \perp\!\!\!\perp Y | Z$ if $\rho_{XZ}\rho_{YZ} = \rho_{XY}$

### Proof.

- From the definition of correlation coefficient,
$$\Sigma = \begin{pmatrix} \sigma_{XX} & \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} \\ \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sigma_{YY} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} \\ \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} & \sigma_{ZZ} \end{pmatrix}$$

# $X \perp\!\!\!\perp Y | Z$ if $\rho_{XZ}\rho_{YZ} = \rho_{XY}$

**Proof.**

- From the definition of correlation coefficient,
$$\Sigma = \begin{pmatrix} \sigma_{XX} & \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} \\ \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sigma_{YY} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} \\ \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} & \sigma_{ZZ} \end{pmatrix}$$

- Then from the conditioning result, we have

$$\Sigma_{\left(\begin{smallmatrix}X\\Y\end{smallmatrix}\right)|z} = \begin{pmatrix} \sigma_{XX} & \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} \\ \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sigma_{YY} \end{pmatrix}$$

$$- \begin{pmatrix} \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} \end{pmatrix} \sigma_{ZZ}^{-1} \begin{pmatrix} \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} \\ \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} \end{pmatrix}$$

$$= \begin{pmatrix} \sigma_{XX}(1-\rho_{XZ}^2) & \sqrt{\sigma_{XX}\sigma_{YY}}(\rho_{XY} - \rho_{XZ}\rho_{YZ}) \\ \sqrt{\sigma_{XX}\sigma_{YY}}(\rho_{XY} - \rho_{XZ}\rho_{YZ}) & \sigma_{YY}(1-\rho_{YZ}^2) \end{pmatrix}$$

# $X \perp\!\!\!\perp Y | Z$ if $\rho_{XZ}\rho_{YZ} = \rho_{XY}$

### Proof.

- From the definition of correlation coefficient,

$$\Sigma = \begin{pmatrix} \sigma_{XX} & \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} \\ \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sigma_{YY} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} \\ \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} & \sigma_{ZZ} \end{pmatrix}$$

- Then from the conditioning result, we have

$$\begin{aligned} \Sigma_{\binom{X}{Y}|z} &= \begin{pmatrix} \sigma_{XX} & \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} \\ \sqrt{\sigma_{XX}\sigma_{YY}}\rho_{XY} & \sigma_{YY} \end{pmatrix} \\ &\qquad - \begin{pmatrix} \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} & \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} \end{pmatrix} \sigma_{ZZ}^{-1} \begin{pmatrix} \sqrt{\sigma_{XX}\sigma_{ZZ}}\rho_{XZ} \\ \sqrt{\sigma_{YY}\sigma_{ZZ}}\rho_{YZ} \end{pmatrix} \\ &= \begin{pmatrix} \sigma_{XX}(1-\rho_{XZ}^2) & \sqrt{\sigma_{XX}\sigma_{YY}}(\rho_{XY}-\rho_{XZ}\rho_{YZ}) \\ \sqrt{\sigma_{XX}\sigma_{YY}}(\rho_{XY}-\rho_{XZ}\rho_{YZ}) & \sigma_{YY}(1-\rho_{YZ}^2) \end{pmatrix} \end{aligned}$$

- Therefore, $X$ and $Y$ are uncorrelated given $Z$ when the off-diagonal is zero and this gives us $\rho_{XY} = \rho_{XZ}\rho_{YZ}$. Since for Gaussian variables, uncorrelatedness implies independence. This concludes the proof.

## Gaussian Process

- Consider a 1-D discrete-time signal, and say the signal is joint Gaussian and two points are conditional independent given points in the middle

- If the variance is stationary and say the correlation coefficent between two adjacent points is $\rho$, further assume that the variance is normalized to 1. WLOG, then

$$\Sigma = \begin{pmatrix} 1 & \rho & \rho^2 & \cdots & \\ \rho & 1 & \rho & \rho^2 & \cdots \\ \rho^2 & \rho & 1 & \rho & \cdots \\ & & \cdots & & \end{pmatrix}$$

# Product of normal distributions

- Assume that we tries to recover some vector parameter **x**, which is subject to multivariate Gaussian noise

# Product of normal distributions

- Assume that we tries to recover some vector parameter **x**, which is subject to multivariate Gaussian noise
- Say we made two measurements $\mathbf{y}_1$ and $\mathbf{y}_2$, where $\mathbf{Y}_1 \sim \mathcal{N}(\mathbf{x}, \Sigma_{\mathbf{Y}_1})$ and $\mathbf{Y}_2 \sim \mathcal{N}(\mathbf{x}, \Sigma_{\mathbf{Y}_2})$. Note that even though both measurements have mean **x**, they have different covariance
  - This variation, for instance, can be due to environment change between the two measurements

## Product of normal distributions

- Assume that we tries to recover some vector parameter $\mathbf{x}$, which is subject to multivariate Gaussian noise
- Say we made two measurements $\mathbf{y}_1$ and $\mathbf{y}_2$, where $\mathbf{Y}_1 \sim \mathcal{N}(\mathbf{x}, \Sigma_{\mathbf{Y}_1})$ and $\mathbf{Y}_2 \sim \mathcal{N}(\mathbf{x}, \Sigma_{\mathbf{Y}_2})$. Note that even though both measurements have mean $\mathbf{x}$, they have different covariance
  - This variation, for instance, can be due to environment change between the two measurements
- Now, if we want to compute the overall likelihood, $p(\mathbf{y}_1, \mathbf{y}_2|\mathbf{x})$. Assuming that $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are conditionally independent given $\mathbf{X}$, we have

$$p(\mathbf{y}_1, \mathbf{y}_2|\mathbf{x}) = p(\mathbf{y}_1|\mathbf{x})p(\mathbf{y}_2|\mathbf{x})$$
$$= \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}).$$

## Product of normal distributions

- Assume that we tries to recover some vector parameter **x**, which is subject to multivariate Gaussian noise
- Say we made two measurements $\mathbf{y}_1$ and $\mathbf{y}_2$, where $\mathbf{Y}_1 \sim \mathcal{N}(\mathbf{x}, \Sigma_{\mathbf{Y}_1})$ and $\mathbf{Y}_2 \sim \mathcal{N}(\mathbf{x}, \Sigma_{\mathbf{Y}_2})$. Note that even though both measurements have mean **x**, they have different covariance
    - This variation, for instance, can be due to environment change between the two measurements
- Now, if we want to compute the overall likelihood, $p(\mathbf{y}_1, \mathbf{y}_2|\mathbf{x})$. Assuming that $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are conditionally independent given **X**, we have

$$p(\mathbf{y}_1, \mathbf{y}_2|\mathbf{x}) = p(\mathbf{y}_1|\mathbf{x})p(\mathbf{y}_2|\mathbf{x})$$
$$= \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}).$$

- Essentially, we just need to compute the product of two Gaussian pdfs. Such computation is very useful and it occurs often when one needs to perform inference

# Product of normal distributions

As in previous cases, the product turns out to be normal also. However, unlike them, the product is not a pdf and so it does not normalize to 1. So we have to compute both the scaling factor and the exponent explicitly. Let us start with the exponent.

$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})$$

# Product of normal distributions

As in previous cases, the product turns out to be normal also. However, unlike them, the product is not a pdf and so it does not normalize to 1. So we have to compute both the scaling factor and the exponent explicitly. Let us start with the exponent.

$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})$$

$$\propto \exp\left(-\frac{1}{2}[(\mathbf{x} - \mathbf{y}_1)^T \Lambda_{\mathbf{Y}_1}(\mathbf{x} - \mathbf{y}_1) + (\mathbf{x} - \mathbf{y}_2)^T \Lambda_{\mathbf{Y}_2}(\mathbf{x} - \mathbf{y}_2)]\right)$$

## Product of normal distributions

As in previous cases, the product turns out to be normal also. However, unlike them, the product is not a pdf and so it does not normalize to 1. So we have to compute both the scaling factor and the exponent explicitly. Let us start with the exponent.

$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1}) \mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})$$

$$\propto \exp\left(-\frac{1}{2}[(\mathbf{x}-\mathbf{y}_1)^T \Lambda_{\mathbf{Y}_1}(\mathbf{x}-\mathbf{y}_1) + (\mathbf{x}-\mathbf{y}_2)^T \Lambda_{\mathbf{Y}_2}(\mathbf{x}-\mathbf{y}_2)]\right)$$

$$\propto \exp\left(-\frac{1}{2}[\mathbf{x}^T(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})\mathbf{x} - (\mathbf{y}_2^T \Lambda_{\mathbf{Y}_2} + \mathbf{y}_1^T \Lambda_{\mathbf{Y}_1})\mathbf{x} - \mathbf{x}^T(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1)]\right)$$

# Product of normal distributions

As in previous cases, the product turns out to be normal also. However, unlike them, the product is not a pdf and so it does not normalize to 1. So we have to compute both the scaling factor and the exponent explicitly. Let us start with the exponent.

$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})$$

$$\propto \exp\left(-\frac{1}{2}[(\mathbf{x}-\mathbf{y}_1)^T\Lambda_{\mathbf{Y}_1}(\mathbf{x}-\mathbf{y}_1) + (\mathbf{x}-\mathbf{y}_2)^T\Lambda_{\mathbf{Y}_2}(\mathbf{x}-\mathbf{y}_2)]\right)$$

$$\propto \exp\left(-\frac{1}{2}[\mathbf{x}^T(\Lambda_{\mathbf{Y}_1}+\Lambda_{\mathbf{Y}_2})\mathbf{x} - (\mathbf{y}_2^T\Lambda_{\mathbf{Y}_2}+\mathbf{y}_1^T\Lambda_{\mathbf{Y}_1})\mathbf{x} - \mathbf{x}^T(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2+\Lambda_{\mathbf{Y}_1}\mathbf{y}_1)]\right)$$

$$\propto e^{-\frac{1}{2}[(\mathbf{x}-(\Lambda_{\mathbf{Y}_1}+\Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2+\Lambda_{\mathbf{Y}_1}\mathbf{y}_1))^T(\Lambda_{\mathbf{Y}_1}+\Lambda_{\mathbf{Y}_2})(\mathbf{x}-(\Lambda_{\mathbf{Y}_1}+\Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2+\Lambda_{\mathbf{Y}_1}\mathbf{y}_1))]}$$

# Product of normal distributions

As in previous cases, the product turns out to be normal also. However, unlike them, the product is not a pdf and so it does not normalize to 1. So we have to compute both the scaling factor and the exponent explicitly. Let us start with the exponent.
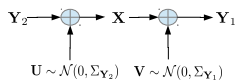
$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1}) \mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})$$

$$\propto \exp\left(-\frac{1}{2}[(\mathbf{x} - \mathbf{y}_1)^T \Lambda_{\mathbf{Y}_1}(\mathbf{x} - \mathbf{y}_1) + (\mathbf{x} - \mathbf{y}_2)^T \Lambda_{\mathbf{Y}_2}(\mathbf{x} - \mathbf{y}_2)]\right)$$

$$\propto \exp\left(-\frac{1}{2}[\mathbf{x}^T(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})\mathbf{x} - (\mathbf{y}_2^T \Lambda_{\mathbf{Y}_2} + \mathbf{y}_1^T \Lambda_{\mathbf{Y}_1})\mathbf{x} - \mathbf{x}^T(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1)]\right)$$

$$\propto e^{-\frac{1}{2}[(\mathbf{x} - (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1))^T(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})(\mathbf{x} - (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1))]}$$

$$\propto \mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})$$

Therefore,

$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1}) \mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})$$

$$= K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}) \mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})$$

for some scaling factor $K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})$ independent of $\mathbf{x}$

# Product of normal distributions

- One can compute the scaling factor $K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})$ directly

## Product of normal distributions

- One can compute the scaling factor $K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})$ directly
- However, it is much easier to take advantage for the following setup when $\mathbf{Y}_1 \perp\!\!\!\perp \mathbf{Y}_2 | \mathbf{X}$ as shown below

## Product of normal distributions

- One can compute the scaling factor $K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})$ directly
- However, it is much easier to take advantage for the following setup when $\mathbf{Y}_1 \perp\!\!\!\perp \mathbf{Y}_2 | \mathbf{X}$ as shown below

$$\mathbf{Y}_2 \longrightarrow \bigoplus \longrightarrow \mathbf{X} \longrightarrow \bigoplus \longrightarrow \mathbf{Y}_1$$

$$\mathbf{U} \sim \mathcal{N}(0, \Sigma_{\mathbf{Y}_2}) \quad \mathbf{V} \sim \mathcal{N}(0, \Sigma_{\mathbf{Y}_1})$$

- Since $\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) = \mathcal{N}(\mathbf{x}; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2})$ and $\mathbf{Y}_1 \perp\!\!\!\perp \mathbf{Y}_2 | \mathbf{X}$, we have

$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) = \underbrace{\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})}_{p(y_1|x)=p(y_1|x,y_2)}\underbrace{\mathcal{N}(\mathbf{x}; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2})}_{p(x|y_2)} = p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)$$

# Product of normal distributions

- Then, marginalizing **x** out from $p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)$, we have



$p(\mathbf{y}_1|\mathbf{y}_2) = \int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x}$. However, from the figure,

$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$$
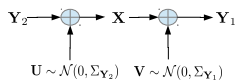
# Product of normal distributions

- Then, marginalizing $\mathbf{x}$ out from $p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)$, we have



$p(\mathbf{y}_1|\mathbf{y}_2) = \int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x}$. However, from the figure,

$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$$

- On the other hand,

$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = \int \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})d\mathbf{x}$$

$$= \int K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})d\mathbf{x}$$

$$= K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}).$$

## Product of normal distributions

- Then, marginalizing **x** out from $p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)$, we have

$$\mathbf{Y}_2 \longrightarrow \bigoplus \longrightarrow \mathbf{X} \longrightarrow \bigoplus \longrightarrow \mathbf{Y}_1$$

$$\mathbf{U} \sim \mathcal{N}(0, \Sigma_{\mathbf{Y}_2}) \quad \mathbf{V} \sim \mathcal{N}(0, \Sigma_{\mathbf{Y}_1})$$

$p(\mathbf{y}_1|\mathbf{y}_2) = \int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x}$. However, from the figure,

$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$$

- On the other hand,

$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = \int \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})d\mathbf{x}$$

$$= \int K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})d\mathbf{x}$$

$$= K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}).$$

- Thus we have $K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$

## Product of normal distributions

- Then, marginalizing $\mathbf{x}$ out from $p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)$, we have



$p(\mathbf{y}_1|\mathbf{y}_2) = \int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x}$. However, from the figure,

$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$$

- On the other hand,

$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = \int \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})d\mathbf{x}$$

$$= \int K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})d\mathbf{x}$$

$$= K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}).$$

- Thus we have $K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$ and so

$$\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})$$

$$= \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})$$

## Product of normal distributions

Let us try to interpret the product as the overall likelihood after making two observations. Consider the simpler case when $\mathbf{X}$, $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are all scaler

## Product of normal distributions

Let us try to interpret the product as the overall likelihood after making two observations. Consider the simpler case when $\mathbf{X}$, $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are all scaler

- The mean considering both observations,
  $(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y)$, is essential a weighted average of observations $\mathbf{y}_2$ and $\mathbf{y}_1$
    - The weight is higher when the precision $\Lambda_{\mathbf{Y}_2}$ or $\Lambda_{\mathbf{Y}_1}$ is larger

## Product of normal distributions

Let us try to interpret the product as the overall likelihood after making two observations. Consider the simpler case when $\mathbf{X}$, $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are all scaler

- The mean considering both observations,
  $(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y)$, is essential a weighted average of observations $\mathbf{y}_2$ and $\mathbf{y}_1$
    - The weight is higher when the precision $\Lambda_{\mathbf{Y}_2}$ or $\Lambda_{\mathbf{Y}_1}$ is larger
- The overall variance $(\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1}$ is always smaller than the individual variance $\Sigma_{\mathbf{Y}_2}$ and $\Sigma_{\mathbf{Y}_1}$
    - We are more certain with $\mathbf{x}$ after considering both $\mathbf{y}_1$ and $\mathbf{y}_2$

## Product of normal distributions

Let us try to interpret the product as the overall likelihood after making two observations. Consider the simpler case when $\mathbf{X}$, $\mathbf{Y}_1$ and $\mathbf{Y}_2$ are all scaler

- The mean considering both observations, $(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y)$, is essential a weighted average of observations $\mathbf{y}_2$ and $\mathbf{y}_1$
    - The weight is higher when the precision $\Lambda_{\mathbf{Y}_2}$ or $\Lambda_{\mathbf{Y}_1}$ is larger
- The overall variance $(\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1}$ is always smaller than the individual variance $\Sigma_{\mathbf{Y}_2}$ and $\Sigma_{\mathbf{Y}_1}$
    - We are more certain with $\mathbf{x}$ after considering both $\mathbf{y}_1$ and $\mathbf{y}_2$
- The scaling factor, $\mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$, can be interpreted as how much one can believe on the overall likelihood.
    - The value is reasonable since when the two observations are far away with respect to the overall variance $\Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1}$, the likelihood will become less reliable
    - The scaling factor is especially useful when we deal with mixture of Gaussian to be discussed next

## Review

- PCA (assume zero mean)
  - Via eigen-decomposition
    1. $\Sigma \approx \frac{1}{m} \mathcal{X}^T \mathcal{X}$
    2. $P^T \Sigma P = D$
    3. $Y = P^T X$
  - Via SVD
    1. $U^T \mathcal{X} V = D$
    2. $Y = V^T X$
- Marginalization of a normal distribution is still a normal distribution
- Conditioning of normal distribution:
  $\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu_X} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \boldsymbol{\mu_Y}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$
- Product of normal distribution:
  $\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) =$
  $\mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})$

## Division of normal distributions

- To compute $\frac{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)}{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)}$, note that from the product formula earlier

$$\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)\mathcal{N}(\mathbf{x};(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),(\Lambda_1-\Lambda_2)^{-1})$$
$$=\mathcal{N}(\boldsymbol{\mu}_2;(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),\Lambda_2^{-1}+(\Lambda_1-\Lambda_2)^{-1})\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)$$

## Division of normal distributions

- To compute $\frac{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)}{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)}$, note that from the product formula earlier

$$\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)\mathcal{N}(\mathbf{x};(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),(\Lambda_1-\Lambda_2)^{-1})$$
$$=\mathcal{N}(\boldsymbol{\mu}_2;(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),\Lambda_2^{-1}+(\Lambda_1-\Lambda_2)^{-1})\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)$$

- Therefore,

$$\frac{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)}{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)} = \frac{\mathcal{N}(\mathbf{x};(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),(\Lambda_1-\Lambda_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_2;(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),\Lambda_2^{-1}+(\Lambda_1-\Lambda_2)^{-1})}$$
$$= \frac{\mathcal{N}(\mathbf{x};\boldsymbol{\mu},(\Lambda_1-\Lambda_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_2;\boldsymbol{\mu},\Lambda_2^{-1}+(\Lambda_1-\Lambda_2)^{-1})},$$

where $\boldsymbol{\mu} = (\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2)$

## Division of normal distributions

- To compute $\frac{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)}{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)}$, note that from the product formula earlier

$$\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)\mathcal{N}(\mathbf{x};(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),(\Lambda_1-\Lambda_2)^{-1})$$
$$=\mathcal{N}(\boldsymbol{\mu}_2;(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),\Lambda_2^{-1}+(\Lambda_1-\Lambda_2)^{-1})\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)$$

- Therefore,

$$\frac{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_1,\Sigma_1)}{\mathcal{N}(\mathbf{x};\boldsymbol{\mu}_2,\Sigma_2)} = \frac{\mathcal{N}(\mathbf{x};(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),(\Lambda_1-\Lambda_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_2;(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2),\Lambda_2^{-1}+(\Lambda_1-\Lambda_2)^{-1})}$$
$$= \frac{\mathcal{N}(\mathbf{x};\boldsymbol{\mu},(\Lambda_1-\Lambda_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_2;\boldsymbol{\mu},\Lambda_2^{-1}+(\Lambda_1-\Lambda_2)^{-1})},$$

where $\boldsymbol{\mu}=(\Lambda_1-\Lambda_2)^{-1}(\Lambda_1\boldsymbol{\mu}_1-\Lambda_2\boldsymbol{\mu}_2)$

- Note that the final pdf will be Gaussian-like if $\Lambda_1 \succeq \Lambda_2$. Otherwise, one can still write out the pdf using the precision matrix. But the covariance matrix will not be defined (Try plot some pdfs out yourselves)

## Mixture of Gaussians

Consider an electrical system that outputs signal of different statistics when it is on and off

- When the system is on, the output signal $S$ behaves like $\mathcal{N}(5, 1)$. When the system is off is off, $S$ behaves like $\mathcal{N}(0, 1)$

## Mixture of Gaussians

Consider an electrical system that outputs signal of different statistics when it is on and off

- When the system is on, the output signal $S$ behaves like $\mathcal{N}(5, 1)$. When the system is off is off, $S$ behaves like $\mathcal{N}(0, 1)$
- If someone measuring the signal does not know the status of the system but only knows that the system is on 40% of the time, then to the observer, the signal $S$ behaves like a mixture of Gaussians

## Mixture of Gaussians

Consider an electrical system that outputs signal of different statistics when it is on and off

- When the system is on, the output signal $S$ behaves like $\mathcal{N}(5,1)$. When the system is off is off, $S$ behaves like $\mathcal{N}(0,1)$
- If someone measuring the signal does not know the status of the system but only knows that the system is on 40% of the time, then to the observer, the signal $S$ behaves like a mixture of Gaussians
- The pdf of $S$ will be $0.4\mathcal{N}(s;5,1) + 0.6\mathcal{N}(s;0,1)$ as shown below

## Mixture of Gaussians

- A main limitation of normal distribution is that it is unimodal

## Mixture of Gaussians

- A main limitation of normal distribution is that it is unimodal
- Mixture of Gaussian distribution allows multimodal and can virtually model any pdfs. But there is a computational cost for this gain

## Mixture of Gaussians

- A main limitation of normal distribution is that it is unimodal
- Mixture of Gaussian distribution allows multimodal and can virtually model any pdfs. But there is a computational cost for this gain
- Let us illustrate this with the following example:
  - Consider two mixtures of Gaussian likelihood of $x$ given two observations $y_1$ and $y_2$ as follows:

$$p(y_1|x) = 0.6\mathcal{N}(x; 0, 1) + 0.4\mathcal{N}(x; 5, 1);$$
$$p(y_2|x) = 0.5\mathcal{N}(x; -2, 1) + 0.5\mathcal{N}(x; 4, 1).$$

What is the overall likelihood, $p(y_1, y_2|x)$?

## Mixture of Gaussians

- A main limitation of normal distribution is that it is unimodal
- Mixture of Gaussian distribution allows multimodal and can virtually model any pdfs. But there is a computational cost for this gain
- Let us illustrate this with the following example:
  - Consider two mixtures of Gaussian likelihood of $x$ given two observations $y_1$ and $y_2$ as follows:

  $$p(y_1|x) = 0.6\mathcal{N}(x; 0, 1) + 0.4\mathcal{N}(x; 5, 1);$$
  $$p(y_2|x) = 0.5\mathcal{N}(x; -2, 1) + 0.5\mathcal{N}(x; 4, 1).$$

  What is the overall likelihood, $p(y_1, y_2|x)$?

- As usual, it is reasonable to assume the observations to be conditionally independent given $x$. Then,

$$
\begin{aligned}
p(y_1, y_2|x) &= p(y_1|x)p(y_2|x) \\
&= (0.6\mathcal{N}(x; 0, 1) + 0.4\mathcal{N}(x; 5, 1))(0.5\mathcal{N}(x; -2, 1) + 0.5\mathcal{N}(x; 4, 1)) \\
&= 0.3\mathcal{N}(x; 0, 1)\mathcal{N}(x; -2, 1) + 0.2\mathcal{N}(x; 5, 1)\mathcal{N}(x; -2, 1) \\
&\quad + 0.3\mathcal{N}(x; 0, 1)\mathcal{N}(x; 4, 1) + 0.2\mathcal{N}(x; 5, 1)\mathcal{N}(x; 4, 1)
\end{aligned}
$$

## Explosion of Gaussians

- The last step involves computing products of Gaussians but we have learned it in previous sections. Using the previous result,

$$p(y_1, y_2|x) = 0.3\mathcal{N}(-2; 0, 2)\mathcal{N}(x; -1, 0.5) + 0.2\mathcal{N}(-2; 5, 2)\mathcal{N}(x; 1.5, 0.5)$$
$$+ 0.3\mathcal{N}(4; 0, 2)\mathcal{N}(x; 2, 0.5) + 0.2\mathcal{N}(4; 5, 2)\mathcal{N}(x; 4.5, 0.5).$$

# Explosion of Gaussians

- The last step involves computing products of Gaussians but we have learned it in previous sections. Using the previous result,

$$p(y_1, y_2 | x) = 0.3 \mathcal{N}(-2; 0, 2) \mathcal{N}(x; -1, 0.5) + 0.2 \mathcal{N}(-2; 5, 2) \mathcal{N}(x; 1.5, 0.5)$$
$$+ 0.3 \mathcal{N}(4; 0, 2) \mathcal{N}(x; 2, 0.5) + 0.2 \mathcal{N}(4; 5, 2) \mathcal{N}(x; 4.5, 0.5).$$

So we have the overall likelihood is a mixture of four Gaussians

## Explosion of Gaussians

- The last step involves computing products of Gaussians but we have learned it in previous sections. Using the previous result,

$$p(y_1, y_2|x) = 0.3\mathcal{N}(-2; 0, 2)\mathcal{N}(x; -1, 0.5) + 0.2\mathcal{N}(-2; 5, 2)\mathcal{N}(x; 1.5, 0.5)$$
$$+ 0.3\mathcal{N}(4; 0, 2)\mathcal{N}(x; 2, 0.5) + 0.2\mathcal{N}(4; 5, 2)\mathcal{N}(x; 4.5, 0.5).$$

  So we have the overall likelihood is a mixture of four Gaussians

- Let's repeat our discussion but with $n$ observations instead. The overall likelihood will be a mixture of $2^n$ Gaussians!

  - Therefore, the computation will quickly become intractable as the number of observations increases

## Explosion of Gaussians

- The last step involves computing products of Gaussians but we have learned it in previous sections. Using the previous result,

$$p(y_1, y_2|x) = 0.3\mathcal{N}(-2; 0, 2)\mathcal{N}(x; -1, 0.5) + 0.2\mathcal{N}(-2; 5, 2)\mathcal{N}(x; 1.5, 0.5)$$
$$+ 0.3\mathcal{N}(4; 0, 2)\mathcal{N}(x; 2, 0.5) + 0.2\mathcal{N}(4; 5, 2)\mathcal{N}(x; 4.5, 0.5).$$

  So we have the overall likelihood is a mixture of four Gaussians

- Let's repeat our discussion but with $n$ observations instead. The overall likelihood will be a mixture of $2^n$ Gaussians!
  - Therefore, the computation will quickly become intractable as the number of observations increases
  - Fortunately, in reality, some of the Gaussians in the mixture tend to have a very small weight

## Reduce number of components in Gaussian mixtures

- For instance, in our previous numerical example, if we continue our numerical computation for the two observation example, we have

$$p(y_1, y_2|x) = 0.4163 \mathcal{N}(x; -1, 0.5) + 3.5234 \times 10^{-6} \mathcal{N}(x; 1.5, 0.5)$$
$$+ 0.0202 \mathcal{N}(x; 2, 0.5) + 0.5734 \mathcal{N}(x; 4.5, 0.5).$$

## Reduce number of components in Gaussian mixtures

- For instance, in our previous numerical example, if we continue our numerical computation for the two observation example, we have

$$p(y_1, y_2|x) = 0.4163\mathcal{N}(x; -1, 0.5) + 3.5234 \times 10^{-6}\mathcal{N}(x; 1.5, 0.5)$$
$$+ 0.0202\mathcal{N}(x; 2, 0.5) + 0.5734\mathcal{N}(x; 4.5, 0.5).$$

- We can see that the weight for the component at mean 1.5 is very small. And the component at mean 2 has a rather small weight also.

## Reduce number of components in Gaussian mixtures

- For instance, in our previous numerical example, if we continue our numerical computation for the two observation example, we have

$$p(y_1, y_2|x) = 0.4163\mathcal{N}(x; -1, 0.5) + 3.5234 \times 10^{-6}\mathcal{N}(x; 1.5, 0.5)$$
$$+ 0.0202\mathcal{N}(x; 2, 0.5) + 0.5734\mathcal{N}(x; 4.5, 0.5).$$

- We can see that the weight for the component at mean 1.5 is very small. And the component at mean 2 has a rather small weight also.
- Even with the four Gaussian components, the overall likelihood is essentially just a bimodal distribution as shown in the figure below

# Reduce number of components in Gaussian mixtures

- Therefore, we may approximate $p(y_1, y_2|x)$ with only two of its original component as $0.4163/(0.4163 + 0.5734)\mathcal{N}(x; -1, 0.5) + 0.5734/(0.4163 + 0.5734)\mathcal{N}(x; 4.5, 0.5) = 0.4206\mathcal{N}(x; -1, 0.5) + 0.5794\mathcal{N}(x; 4.5, 0.5)$

# Reduce number of components in Gaussian mixtures

- Therefore, we may approximate $p(y_1, y_2|x)$ with only two of its original component as $0.4163/(0.4163 + 0.5734)\mathcal{N}(x; -1, 0.5) + 0.5734/(0.4163 + 0.5734)\mathcal{N}(x; 4.5, 0.5) = 0.4206\mathcal{N}(x; -1, 0.5) + 0.5794\mathcal{N}(x; 4.5, 0.5)$

- However, it is not always a good approximation strategy just to dump away the small components in a Gaussian mixture

## Another example

Consider

$$p(x) = 0.1\mathcal{N}(x; -0.2, 1) + 0.1\mathcal{N}(x; -0.1, 1) + 0.1\mathcal{N}(x; 0, 1) + 0.1\mathcal{N}(x; 0.1, 1)$$
$$+ 0.1\mathcal{N}(x; 0.2, 1) + 0.5\mathcal{N}(x; 5, 1).$$

## Another example

Consider

$$p(x) = 0.1\mathcal{N}(x; -0.2, 1) + 0.1\mathcal{N}(x; -0.1, 1) + 0.1\mathcal{N}(x; 0, 1) + 0.1\mathcal{N}(x; 0.1, 1)$$
$$+ 0.1\mathcal{N}(x; 0.2, 1) + 0.5\mathcal{N}(x; 5, 1).$$

- Let say we want to reduce $p(x)$ to only a mixture of two Gaussians. It is tempting to just dumping four smallest one and renormalized the weight. For example, if we choose to remove the first four components, we have

$$\hat{p}(x) = 1/6\mathcal{N}(x; 0.2, 1) + 5/6\mathcal{N}(x; 5, 1)$$

## Another example

Consider

$$p(x) = 0.1\mathcal{N}(x; -0.2, 1) + 0.1\mathcal{N}(x; -0.1, 1) + 0.1\mathcal{N}(x; 0, 1) + 0.1\mathcal{N}(x; 0.1, 1)$$
$$+ 0.1\mathcal{N}(x; 0.2, 1) + 0.5\mathcal{N}(x; 5, 1).$$

- Let say we want to reduce $p(x)$ to only a mixture of two Gaussians. It is tempting to just dumping four smallest one and renormalized the weight. For example, if we choose to remove the first four components, we have

$$\hat{p}(x) = 1/6\mathcal{N}(x; 0.2, 1) + 5/6\mathcal{N}(x; 5, 1)$$

- The approximation $\hat{p}(x)$ is significantly different from $p(x)$ as shown below

## Merging components

- The problem is that while the first five components are all relatively small compared to the last one, they are all quite similar and their combined contribution is comparable to the latter

## Merging components

- The problem is that while the first five components are all relatively small compared to the last one, they are all quite similar and their combined contribution is comparable to the latter
- Actually the first five components are so similar that their combined contribution can be accurately modeled as one Gaussian

# Merging components

- The problem is that while the first five components are all relatively small compared to the last one, they are all quite similar and their combined contribution is comparable to the latter
- Actually the first five components are so similar that their combined contribution can be accurately modeled as one Gaussian
- So rather than discarding the components, one can get a much more accurate approximation by merging them. The approximation is illustrated as $\tilde{p}(x)$ in the figure below

# Merging components

To successfully obtain such approximation $\tilde{p}(x)$, we have to answer two questions:

- which components to merge?
- how to merge them?

## Which Components to Merge?

It is reasonable to pick similar components to merge. The question is how do will gauge the similarity between two components.

## Which Components to Merge?

It is reasonable to pick similar components to merge. The question is how do will gauge the similarity between two components.

- Consider two pdfs $p(\mathbf{x})$ and $q(\mathbf{x})$, note that we can define an inner product of $p(\mathbf{x})$ and $q(\mathbf{x})$ by

$$\langle p(\mathbf{x}), q(\mathbf{x}) \rangle = \int p(\mathbf{x}) q(\mathbf{x}) d\mathbf{x}$$

## Which Components to Merge?

It is reasonable to pick similar components to merge. The question is how do will gauge the similarity between two components.

- Consider two pdfs $p(\mathbf{x})$ and $q(\mathbf{x})$, note that we can define an inner product of $p(\mathbf{x})$ and $q(\mathbf{x})$ by

$$\langle p(\mathbf{x}), q(\mathbf{x})\rangle = \int p(\mathbf{x})q(\mathbf{x})d\mathbf{x}$$

- Note that the inner product is well defined and $\langle p(\mathbf{x}), p(\mathbf{x})\rangle \geq 0$

## Which Components to Merge?

It is reasonable to pick similar components to merge. The question is how do will gauge the similarity between two components.

- Consider two pdfs $p(\mathbf{x})$ and $q(\mathbf{x})$, note that we can define an inner product of $p(\mathbf{x})$ and $q(\mathbf{x})$ by

$$\langle p(\mathbf{x}), q(\mathbf{x}) \rangle = \int p(\mathbf{x}) q(\mathbf{x}) d\mathbf{x}$$

- Note that the inner product is well defined and $\langle p(\mathbf{x}), p(\mathbf{x}) \rangle \geq 0$
- By Cauchy-Schwartz inequality,

$$\frac{\langle p(\mathbf{x}), q(\mathbf{x}) \rangle}{\sqrt{\langle p(\mathbf{x}), p(\mathbf{x}) \rangle \langle q(\mathbf{x}), q(\mathbf{x}) \rangle}} = \frac{\int p(\mathbf{x}) q(\mathbf{x}) d\mathbf{x}}{\sqrt{\int p(\mathbf{x})^2 d\mathbf{x} \int q(\mathbf{x})^2 d\mathbf{x}}} \leq 1$$

## Which Components to Merge?

It is reasonable to pick similar components to merge. The question is how do will gauge the similarity between two components.

- Consider two pdfs $p(\mathbf{x})$ and $q(\mathbf{x})$, note that we can define an inner product of $p(\mathbf{x})$ and $q(\mathbf{x})$ by

$$\langle p(\mathbf{x}), q(\mathbf{x}) \rangle = \int p(\mathbf{x}) q(\mathbf{x}) d\mathbf{x}$$

- Note that the inner product is well defined and $\langle p(\mathbf{x}), p(\mathbf{x}) \rangle \geq 0$
- By Cauchy-Schwartz inequality,

$$\frac{\langle p(\mathbf{x}), q(\mathbf{x}) \rangle}{\sqrt{\langle p(\mathbf{x}), p(\mathbf{x}) \rangle \langle q(\mathbf{x}), q(\mathbf{x}) \rangle}} = \frac{\int p(\mathbf{x}) q(\mathbf{x}) d\mathbf{x}}{\sqrt{\int p(\mathbf{x})^2 d\mathbf{x} \int q(\mathbf{x})^2 d\mathbf{x}}} \leq 1$$

- The inner product maximizes $(= 1)$ when $p(\mathbf{x}) = q(\mathbf{x})$. This suggests a very reasonable similarity measure between two pdfs

# Similarity measure

- Let's define

$$Sim(p(\mathbf{x}), q(\mathbf{x})) \triangleq \frac{\int p(\mathbf{x})q(\mathbf{x})d\mathbf{x}}{\sqrt{\int p(\mathbf{x})^2 d\mathbf{x} \int q(\mathbf{x})^2 d\mathbf{x}}}$$

## Similarity measure

- Let's define

$$Sim(p(\mathbf{x}), q(\mathbf{x})) \triangleq \frac{\int p(\mathbf{x})q(\mathbf{x})d\mathbf{x}}{\sqrt{\int p(\mathbf{x})^2 d\mathbf{x} \int q(\mathbf{x})^2 d\mathbf{x}}}$$

- In particular, if $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_p, \Sigma_p)$ and $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_q, \Sigma_q)$, we have (please verify)

$$Sim(\mathcal{N}(\boldsymbol{\mu}_p, \Sigma_p), \mathcal{N}(\boldsymbol{\mu}_q, \Sigma_q)) = \frac{\mathcal{N}(\boldsymbol{\mu}_p; \boldsymbol{\mu}_q, \Sigma_p + \Sigma_q)}{\sqrt{\mathcal{N}(0; 0, 2\Sigma_p)\mathcal{N}(0; 0, 2\Sigma_q)}},$$

which can be computed very easily and is equal to one only when means and covariances are the same

## How to Merge Components?

Say we have $n$ components $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$, $\mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2)$, $\cdots$, $\mathcal{N}(\boldsymbol{\mu}_n, \Sigma_n)$ with weights $w_1, w_2, \cdots, w_n$. What should the combined component be like?

## How to Merge Components?

Say we have $n$ components $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1), \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2), \cdots, \mathcal{N}(\boldsymbol{\mu}_n, \Sigma_n)$ with weights $w_1, w_2, \cdots, w_n$. What should the combined component be like?

- Combined component weight should equal to total weight $\sum_{i=1}^n w_i$

## How to Merge Components?

Say we have $n$ components $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1), \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2), \cdots, \mathcal{N}(\boldsymbol{\mu}_n, \Sigma_n)$ with weights $w_1, w_2, \cdots, w_n$. What should the combined component be like?

- Combined component weight should equal to total weight $\sum_{i=1}^{n} w_i$
- Combined mean will simply be $\sum_{i=1}^{n} \hat{w}_i \boldsymbol{\mu}_i$, where $\hat{w}_i = \frac{w_i}{\sum_{i=1}^{n} w_i}$

## How to Merge Components?

Say we have $n$ components $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1), \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2), \cdots, \mathcal{N}(\boldsymbol{\mu}_n, \Sigma_n)$ with weights $w_1, w_2, \cdots, w_n$. What should the combined component be like?

- Combined component weight should equal to total weight $\sum_{i=1}^n w_i$
- Combined mean will simply be $\sum_{i=1}^n \hat{w}_i \boldsymbol{\mu}_i$, where $\hat{w}_i = \frac{w_i}{\sum_{i=1}^n w_i}$
- For combined covariance, it may be tempting to approximate it as $\sum_{i=1}^n \hat{w}_i \Sigma_i$.

# How to Merge Components?

Say we have $n$ components $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1), \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2), \cdots, \mathcal{N}(\boldsymbol{\mu}_n, \Sigma_n)$ with weights $w_1, w_2, \cdots, w_n$. What should the combined component be like?

- Combined component weight should equal to total weight $\sum_{i=1}^{n} w_i$
- Combined mean will simply be $\sum_{i=1}^{n} \hat{w}_i \boldsymbol{\mu}_i$, where $\hat{w}_i = \frac{w_i}{\sum_{i=1}^{n} w_i}$
- For combined covariance, it may be tempting to approximate it as $\sum_{i=1}^{n} \hat{w}_i \Sigma_i$.
    - However, it is an underestimate

## How to Merge Components?

Say we have $n$ components $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1)$, $\mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2)$, $\cdots$, $\mathcal{N}(\boldsymbol{\mu}_n, \Sigma_n)$ with weights $w_1, w_2, \cdots, w_n$. What should the combined component be like?

- Combined component weight should equal to total weight $\sum_{i=1}^{n} w_i$
- Combined mean will simply be $\sum_{i=1}^{n} \hat{w}_i \boldsymbol{\mu}_i$, where $\hat{w}_i = \frac{w_i}{\sum_{i=1}^{n} w_i}$
- For combined covariance, it may be tempting to approximate it as $\sum_{i=1}^{n} \hat{w}_i \Sigma_i$.
    - However, it is an underestimate
    - Because the weighted sum only counted the contribution of variation among each component, it did not take into account the variation due to different means across components.

## How to Merge Components?

Say we have $n$ components $\mathcal{N}(\boldsymbol{\mu}_1, \Sigma_1), \mathcal{N}(\boldsymbol{\mu}_2, \Sigma_2), \cdots, \mathcal{N}(\boldsymbol{\mu}_n, \Sigma_n)$ with weights $w_1, w_2, \cdots, w_n$. What should the combined component be like?

- Combined component weight should equal to total weight $\sum_{i=1}^{n} w_i$
- Combined mean will simply be $\sum_{i=1}^{n} \hat{w}_i \boldsymbol{\mu}_i$, where $\hat{w}_i = \frac{w_i}{\sum_{i=1}^{n} w_i}$
- For combined covariance, it may be tempting to approximate it as $\sum_{i=1}^{n} \hat{w}_i \Sigma_i$.
  - However, it is an underestimate
  - Because the weighted sum only counted the contribution of variation among each component, it did not take into account the variation due to different means across components.
  - Instead, let's denote **X** as the variable sampled from the mixture. That is, $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_i, \Sigma_i)$ with probability $\hat{w}_i$. Then, we have (please verify)

$$\Sigma = E[\mathbf{X}\mathbf{X}^T] - E[\mathbf{X}]E[\mathbf{X}]^T$$
$$= \sum_{i=1}^{n} \hat{w}_i(\Sigma_i + \boldsymbol{\mu}_i\boldsymbol{\mu}_i^T) - \sum_{i=1}^{n}\sum_{j=1}^{n} \hat{w}_i\hat{w}_j\boldsymbol{\mu}_i\boldsymbol{\mu}_j^T.$$

Now, go back to our previous numerical example

- Recall that $p(x) = 0.1\mathcal{N}(x; -0.2, 1) + 0.1\mathcal{N}(x; -0.1, 1) +$
  $0.1\mathcal{N}(x; 0, 1) + 0.1\mathcal{N}(x; 0.1, 1) + 0.1\mathcal{N}(x; 0.2, 1) + 0.5\mathcal{N}(x; 5, 1)$

Now, go back to our previous numerical example

- Recall that $p(x) = 0.1\mathcal{N}(x; -0.2, 1) + 0.1\mathcal{N}(x; -0.1, 1) + 0.1\mathcal{N}(x; 0, 1) + 0.1\mathcal{N}(x; 0.1, 1) + 0.1\mathcal{N}(x; 0.2, 1) + 0.5\mathcal{N}(x; 5, 1)$

- If we merge the five smallest components (one can easily check that they are also more similar to each other than to the last component), we have $\tilde{p}(x) = 0.5\mathcal{N}(x; 0, 1.02) + 0.5\mathcal{N}(x; 5, 1)$ as shown again below. The approximate pdf is virtually indistinguishable from the original

## Review multivariate normal

- Marginalization of a normal distribution is still a normal distribution
- Conditioning of normal distribution:
  $\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\mu_{\mathbf{X}} + \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}(\mathbf{y} - \mu_{\mathbf{Y}}), \Sigma_{\mathbf{XX}} - \Sigma_{\mathbf{XY}}\Sigma_{\mathbf{YY}}^{-1}\Sigma_{\mathbf{YX}})$
- Product of normal distribution:
  $\mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) =$
  $\mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}y), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})$
- Division of normal distribution:

$$\frac{\mathcal{N}(\mathbf{x}; \mu_1, \Sigma_1)}{\mathcal{N}(\mathbf{x}; \mu_2, \Sigma_2)} = \frac{\mathcal{N}(\mathbf{x}; \mu, (\Lambda_1 - \Lambda_2)^{-1})}{\mathcal{N}(\mu_2; \mu, \Lambda_2^{-1} + (\Lambda_1 - \Lambda_2)^{-1})},$$

  where $\mu = (\Lambda_1 - \Lambda_2)^{-1}(\Lambda_1\mu_1 - \Lambda_2\mu_2)$

- Similarity measure

$$Sim(\mathcal{N}(\mu_p, \Sigma_p), \mathcal{N}(\mu_q, \Sigma_q)) = \frac{\mathcal{N}(\mu_p; \mu_q, \Sigma_p + \Sigma_q)}{\sqrt{\mathcal{N}(0; 0, 2\Sigma_p)\mathcal{N}(0; 0, 2\Sigma_q)}},$$

# Bernoulli distribution

- Consider someone flips a biased coin. The probability of the outcome is described by the Bernoulli distribution. Denote $X = 1$ for a head and $X = 0$ for a tail. Let $Pr(X = 1) = p$.

# Bernoulli distribution

- Consider someone flips a biased coin. The probability of the outcome is described by the Bernoulli distribution. Denote $X = 1$ for a head and $X = 0$ for a tail. Let $Pr(X = 1) = p$. Then the Bernoulli distribution is simply

$$Bern(x|p) = \begin{cases} p, & x = 1 \\ 1 - p, & x = 0 \end{cases}$$

## Bernoulli distribution

- Consider someone flips a biased coin. The probability of the outcome is described by the Bernoulli distribution. Denote $X = 1$ for a head and $X = 0$ for a tail. Let $Pr(X = 1) = p$. Then the Bernoulli distribution is simply

$$Bern(x|p) = \begin{cases} p, & x = 1 \\ 1 - p, & x = 0 \end{cases}$$

- More concisely, we can write it as

$$Bern(x|p) = p^x(1 - p)^{1-x},$$

## Bernoulli distribution

- Consider someone flips a biased coin. The probability of the outcome is described by the Bernoulli distribution. Denote $X = 1$ for a head and $X = 0$ for a tail. Let $Pr(X = 1) = p$. Then the Bernoulli distribution is simply

$$
Bern(x|p) = \begin{cases} p, & x = 1 \\ 1 - p, & x = 0 \end{cases}
$$

- More concisely, we can write it as

$$
Bern(x|p) = p^x (1 - p)^{1-x},
$$

- The mean and variance are

$$
E[X] = p \cdot 1 + (1 - p) \cdot 0 = p
$$

$$
Var[X] = p \cdot (1 - p)^2 + (1 - p) \cdot p^2 = p(1 - p)
$$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1 - p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1-p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1-p)^{N-x}$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1 - p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1 - p)^{N-x}$
  $= Np \sum_{x=1}^{N} \frac{(N-1)!}{(x-1)!(N-x)!} p^{x-1} (1 - p)^{N-x}$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1-p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1-p)^{N-x}$
  $= Np \sum_{x=1}^{N} \frac{(N-1)!}{(x-1)!(N-x)!} p^{x-1} (1-p)^{N-x} = Np \sum_{x=0}^{N-1} Bin(x|p, N-1)$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1-p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1-p)^{N-x}$
    $= Np \sum_{x=1}^{N} \frac{(N-1)!}{(x-1)!(N-x)!} p^{x-1} (1-p)^{N-x} = Np \sum_{x=0}^{N-1} Bin(x|p, N-1)$
    $= Np$

## Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1-p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1-p)^{N-x}$
    $= Np \sum_{x=1}^{N} \frac{(N-1)!}{(x-1)!(N-x)!} p^{x-1}(1-p)^{N-x} = Np \sum_{x=0}^{N-1} Bin(x|p, N-1)$
    $= Np$
  - Similar, $E[X(X-1)] = \sum_{x=2}^{N} \frac{N!}{(x-2)!(N-x)!} p^x (1-p)^{N-x}$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1-p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1-p)^{N-x}$
    $= Np \sum_{x=1}^{N} \frac{(N-1)!}{(x-1)!(N-x)!} p^{x-1}(1-p)^{N-x} = Np \sum_{x=0}^{N-1} Bin(x|p, N-1)$
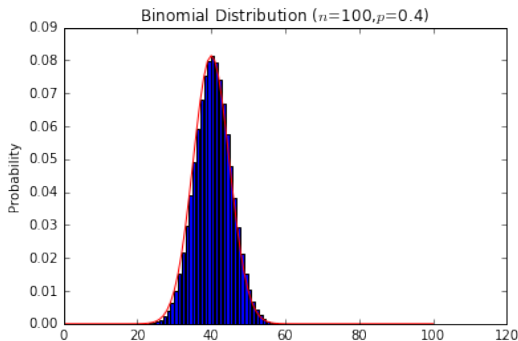    $= Np$
  - Similar, $E[X(X-1)] = \sum_{x=2}^{N} \frac{N!}{(x-2)!(N-x)!} p^x (1-p)^{N-x}$
    $= N(N-1)p^2 \sum_{x=0}^{N-2} Bin(x|p, N-2) = N(N-1)p^2$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1-p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1-p)^{N-x}$
    $= Np \sum_{x=1}^{N} \frac{(N-1)!}{(x-1)!(N-x)!} p^{x-1}(1-p)^{N-x} = Np \sum_{x=0}^{N-1} Bin(x|p, N-1)$
    $= Np$
  - Similar, $E[X(X-1)] = \sum_{x=2}^{N} \frac{N!}{(x-2)!(N-x)!} p^x (1-p)^{N-x}$
    $= N(N-1)p^2 \sum_{x=0}^{N-2} Bin(x|p, N-2) = N(N-1)p^2$
  - Therefore, $Var[X] = E[X^2] - E[X]^2$

# Binomial distribution ($N$ trials)

- Repeat the experiment for $N$ times, the probability of the outcome will now be described by the binomial distribution. Note that $x$ is now the number of obtained heads, we have

$$Bin(x|p, N) = \binom{N}{x} p^x (1-p)^{N-x},$$

- Mean and variances are given by
  - $E[X] = \sum_{x=0}^{N} Bin(x|p)x = \sum_{x=1}^{N} \frac{N!}{(x-1)!(N-x)!} p^x (1-p)^{N-x}$
    $= Np \sum_{x=1}^{N} \frac{(N-1)!}{(x-1)!(N-x)!} p^{x-1} (1-p)^{N-x} = Np \sum_{x=0}^{N-1} Bin(x|p, N-1)$
    $= Np$
  - Similar, $E[X(X-1)] = \sum_{x=2}^{N} \frac{N!}{(x-2)!(N-x)!} p^x (1-p)^{N-x}$
    $= N(N-1)p^2 \sum_{x=0}^{N-2} Bin(x|p, N-2) = N(N-1)p^2$
  - Therefore, $Var[X] = E[X^2] - E[X]^2 = E[X(X-1)] + E[X] - E[X]^2 = N(N-1)p^2 + Np - (Np)^2 = Np(1-p)$

## Binomial distribution

As shown below, the binomial distribution can be model well with a normal distribution $\mathcal{N}(Np, Np(1-p))$ for large $N$



Binomial Distribution ($n=100, p=0.4$)

The binomial distribution is shown in blue and an approximation by normal distribution is shown in red

## Conjugate prior

- Note that both Bernoulli and binomial distributions have the form $p^u(1-p)^v$

## Conjugate prior

- Note that both Bernoulli and binomial distributions have the form $p^u(1-p)^v$
- To estimate $p$, recall that the ML estimator will try to compute

$$\hat{p} = \arg\max_p p(u,v|p) = \arg\max_p p^u(1-p)^v$$

## Conjugate prior

- Note that both Bernoulli and binomial distributions have the form $p^u(1-p)^v$
- To estimate $p$, recall that the ML estimator will try to compute

$$\hat{p} = \arg \max_p p(u, v|p) = \arg \max_p p^u(1-p)^v$$

- Now if we would like to use the MAP estimator instead, we need to introduce a prior $p(p)$ and solve instead

$$\hat{p} = \arg \max_p p(u, v|p)p(p) = \arg \max_p p^u(1-p)^v p(p)$$

## Conjugate prior

- Note that both Bernoulli and binomial distributions have the form $p^u(1-p)^v$
- To estimate $p$, recall that the ML estimator will try to compute

$$\hat{p} = \arg\max_p p(u,v|p) = \arg\max_p p^u(1-p)^v$$

- Now if we would like to use the MAP estimator instead, we need to introduce a prior $p(p)$ and solve instead

$$\hat{p} = \arg\max_p p(u,v|p)p(p) = \arg\max_p p^u(1-p)^v p(p)$$

- It is very difficult to determine the prior unanimously. Actually it can be controversial just to determine the form of it

## Conjugate prior

- Note that both Bernoulli and binomial distributions have the form $p^u(1-p)^v$
- To estimate $p$, recall that the ML estimator will try to compute

$$\hat{p} = \arg\max_p p(u, v|p) = \arg\max_p p^u(1-p)^v$$

- Now if we would like to use the MAP estimator instead, we need to introduce a prior $p(p)$ and solve instead

$$\hat{p} = \arg\max_p p(u, v|p)p(p) = \arg\max_p p^u(1-p)^v p(p)$$

- It is very difficult to determine the prior unanimously. Actually it can be controversial just to determine the form of it
- However, if we select $p(p)$ of a form $p(p) \propto p^a(1-p)^b$, then the resulting posterior distribution with the same form as before. This choice is often chosen for practical purposes, and a prior with same "form" as its likelihood (and thus posterior) is known as the conjugate prior

## Beta distribution

- The conjugate prior of both Bernoulli and binomial distributions is the beta distribution. Its pdf is given by

$$Beta(x|a, b) = \frac{x^{a-1}(1-x)^{b-1}}{B(a, b)},$$

where $X \in [0, 1]$ and $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$

## Beta distribution

- The conjugate prior of both Bernoulli and binomial distributions is the beta distribution. Its pdf is given by

$$Beta(x|a, b) = \frac{x^{a-1}(1 - x)^{b-1}}{B(a, b)},$$

where $X \in [0, 1]$ and $B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$



- Note that with $a = b = 1$, $Beta(x|1, 1) = 1$. It is the same as no prior

# Gamma function

Note that $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} \, dx$

- $\Gamma(1) = \int_0^\infty e^{-x} \, dx = -e^{-x}\big|_0^\infty = 1$

# Gamma function

Note that $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} \, dx$

- $\Gamma(1) = \int_0^\infty e^{-x} \, dx = -e^{-x} \big|_0^\infty = 1$
- For $z > 1$, we have $\Gamma(z) = (z-1)\Gamma(z-1)$

# Gamma function

Note that $\Gamma(z) = \displaystyle\int_0^\infty x^{z-1} e^{-x}\, dx$

- $\Gamma(1) = \displaystyle\int_0^\infty e^{-x}\, dx = -e^{-x}\big|_0^\infty = 1$
- For $z > 1$, we have $\Gamma(z) = (z-1)\Gamma(z-1)$

### Proof.

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x}\, dx$$

# Gamma function

Note that $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x}\, dx$

- $\Gamma(1) = \int_0^\infty e^{-x}\, dx = -e^{-x}\big|_0^\infty = 1$
- For $z > 1$, we have $\Gamma(z) = (z-1)\Gamma(z-1)$

### Proof.

$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x}\, dx = -\int_0^\infty x^{z-1} de^{-x}$

# Gamma function

Note that $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x}\, dx$

- $\Gamma(1) = \int_0^\infty e^{-x}\, dx = -e^{-x}|_0^\infty = 1$
- For $z > 1$, we have $\Gamma(z) = (z-1)\Gamma(z-1)$

Proof.

$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x}\, dx = -\int_0^\infty x^{z-1} de^{-x}$

$= -x^{z-1} e^{-x}|_0^\infty + (z-1) \int_0^\infty x^{z-2} e^{-x}\, dx$

# Gamma function

Note that $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} \, dx$

- $\Gamma(1) = \int_0^\infty e^{-x} \, dx = -e^{-x}|_0^\infty = 1$
- For $z > 1$, we have $\Gamma(z) = (z-1)\Gamma(z-1)$

**Proof.**

$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} \, dx = -\int_0^\infty x^{z-1} de^{-x}$

$= -x^{z-1} e^{-x}|_0^\infty + (z-1)\int_0^\infty x^{z-2} e^{-x} \, dx$

$= (z-1)\int_0^\infty x^{z-2} e^{-x} \, dx = (z-1)\Gamma(z-1)$ □

# Gamma function

Note that $\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} \, dx$

- $\Gamma(1) = \int_0^\infty e^{-x} \, dx = -e^{-x}|_0^\infty = 1$
- For $z > 1$, we have $\Gamma(z) = (z-1)\Gamma(z-1)$

**Proof.**

$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} \, dx = -\int_0^\infty x^{z-1} de^{-x}$

$= -x^{z-1} e^{-x}|_0^\infty + (z-1)\int_0^\infty x^{z-2} e^{-x} \, dx$

$= (z-1)\int_0^\infty x^{z-2} e^{-x} \, dx = (z-1)\Gamma(z-1)$  $\square$

- Therefore, for integer $z > 1$, $\Gamma(z) = (z-1)!$

## Mode of beta distribution

The mode is the peak of a distribution. Recall that
$Beta(x|a, b) = \frac{x^{a-1}(1-x)^{b-1}}{B(a,b)}$. Set

$$\frac{\partial Beta(x|a, b)}{\partial x} = \frac{(a-1)x^{a-2}(1-x)^{b-1} - (b-1)x^{a-1}(1-x)^{b-2}}{B(a, b)} = 0,$$

we have $(a-1)(1-x) = (b-1)x \Rightarrow x = \frac{a-1}{a+b-2}$

# Mean and variance of Beta distribution

Note that $\int_{x=0}^{1} p(x|a, b) = 1 \Rightarrow \int_{x=0}^{1} x^{a-1}(1-x)^{b-1} = B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$.
This gives us a handy trick to manipulate beta distribution

## Mean and variance of Beta distribution

Note that $\int_{x=0}^{1} p(x|a, b) = 1 \Rightarrow \int_{x=0}^{1} x^{a-1}(1-x)^{b-1} = B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$.
This gives us a handy trick to manipulate beta distribution

$$E[X] = \int_{x=0}^{1} xBeta(x|a, b)dx = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{x=0}^{1} x^a (1-x)^{b-1} dx$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} = \frac{a}{a+b}$$

## Mean and variance of Beta distribution

Note that $\int_{x=0}^{1} p(x|a, b) = 1 \Rightarrow \int_{x=0}^{1} x^{a-1}(1-x)^{b-1} = B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$.
This gives us a handy trick to manipulate beta distribution

$$E[X] = \int_{x=0}^{1} x Beta(x|a, b) dx = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \int_{x=0}^{1} x^a (1 - x)^{b-1} dx$$

$$= \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a + 1)\Gamma(b)}{\Gamma(a + b + 1)} = \frac{a}{a + b}$$

Similarly, $E[X^2] = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{x=0}^{1} x^{a+1}(1 - x)^{b-1} dx$

## Mean and variance of Beta distribution

Note that $\int_{x=0}^{1} p(x|a,b) = 1 \Rightarrow \int_{x=0}^{1} x^{a-1}(1-x)^{b-1} = B(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$.
This gives us a handy trick to manipulate beta distribution

$$E[X] = \int_{x=0}^{1} xBeta(x|a,b)dx = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{x=0}^{1} x^{a}(1-x)^{b-1}dx$$

$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} = \frac{a}{a+b}$$

Similarly, $E[X^2] = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{x=0}^{1} x^{a+1}(1-x)^{b-1}dx = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)} =$
$\frac{a(a+1)}{(a+b)(a+b+1)}$.

## Mean and variance of Beta distribution

Note that $\int_{x=0}^{1} p(x|a, b) = 1 \Rightarrow \int_{x=0}^{1} x^{a-1}(1-x)^{b-1} = B(a, b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$.
This gives us a handy trick to manipulate beta distribution

$$E[X] = \int_{x=0}^{1} xBeta(x|a, b)dx = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{x=0}^{1} x^{a}(1-x)^{b-1}dx$$
$$= \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+1)\Gamma(b)}{\Gamma(a+b+1)} = \frac{a}{a+b}$$

Similarly, $E[X^2] = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \int_{x=0}^{1} x^{a+1}(1-x)^{b-1}dx = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)} = \frac{a(a+1)}{(a+b)(a+b+1)}$. Thus,

$$Var[X] = E[X^2] - E[X]^2 = \frac{a(a+1)}{(a+b)(a+b+1)} - \frac{a^2}{(a+b)^2}$$
$$= \frac{a(a+1)(a+b) - a^2(a+b+1)}{(a+b)^2(a+b+1)} = \frac{ab}{(a+b)^2(a+b+1)}$$

## Posterior estimate of probability $p$

Consider the coin flipping example again. Let say the prior probability[4] of the coin is beta distributed with parameters $a$ and $b$. And we flip the coin once to get outcome $x$.

---

[4]Note that this can be very confusing at the beginning. Beware that we are talking about the distribution of the probability of some outcome

## Posterior estimate of probability $p$

Consider the coin flipping example again. Let say the prior probability[4] of the coin is beta distributed with parameters $a$ and $b$. And we flip the coin once to get outcome $x$. Upon observing $x$, we can estimate $p$ by

$$p(p|x, a, b)$$

---

[4]Note that this can be very confusing at the beginning. Beware that we are talking about the distribution of the probability of some outcome

## Posterior estimate of probability $p$

Consider the coin flipping example again. Let say the prior probability[4] of the coin is beta distributed with parameters $a$ and $b$. And we flip the coin once to get outcome $x$. Upon observing $x$, we can estimate $p$ by

$$p(p|x, a, b) = Const1 \cdot Beta(p|a, b) Bern(x|p)$$

---

[4]Note that this can be very confusing at the beginning. Beware that we are talking about the distribution of the probability of some outcome

# Posterior estimate of probability $p$

Consider the coin flipping example again. Let say the prior probability[4] of the coin is beta distributed with parameters $a$ and $b$. And we flip the coin once to get outcome $x$. Upon observing $x$, we can estimate $p$ by

$$p(p|x, a, b) = Const1 \cdot Beta(p|a, b)Bern(x|p)$$
$$= Const2 \cdot p^{a-1+x}(1-p)^{b-1+1-x}$$
$$= Beta(p|\tilde{a}, \tilde{b})$$

So the posterior probability distribution is also beta distributed and the parameters just changed to $\tilde{a} \leftarrow a + x$ and $\tilde{b} \leftarrow b + 1 - x$

---

[4]Note that this can be very confusing at the beginning. Beware that we are talking about the distribution of the probability of some outcome

## Posterior estimate of probability p

Let say we continue our example and we flip the coin by $N$ times and obtain $x$ head. So instead of the Bernoulli likelihood, we have a binomial likelihood. Like the last slide, we have the same beta prior with parameters $a$ and $b$.

## Posterior estimate of probability $p$

Let say we continue our example and we flip the coin by $N$ times and obtain $x$ head. So instead of the Bernoulli likelihood, we have a binomial likelihood. Like the last slide, we have the same beta prior with parameters $a$ and $b$. After the experiment $x$, we can update the distribution of our estimated $p$ by

$$
\begin{aligned}
p(p|x, a, b) =& Const1 \cdot Beta(p|a, b)Bin(x|p, N) \\
=& Const2 \cdot p^{a-1+x}(1-p)^{b-1+N-x} \\
=& Beta(p|\tilde{a}, \tilde{b})
\end{aligned}
$$

## Posterior estimate of probability $p$

Let say we continue our example and we flip the coin by $N$ times and obtain $x$ head. So instead of the Bernoulli likelihood, we have a binomial likelihood. Like the last slide, we have the same beta prior with parameters $a$ and $b$. After the experiment $x$, we can update the distribution of our estimated $p$ by

$$
\begin{aligned}
p(p|x, a, b) =& Const1 \cdot Beta(p|a, b)Bin(x|p, N) \\
=& Const2 \cdot p^{a-1+x}(1-p)^{b-1+N-x} \\
=& Beta(p|\tilde{a}, \tilde{b})
\end{aligned}
$$

Again, the posterior distribution is still beta but with parameters updated to $\tilde{a} \leftarrow a + x$ and $\tilde{b} \leftarrow b + N - x$

## Prior and regularization

- One major reason of introducing prior is for the sake of "regularizing" the answer
- Another coin example
  - Fall back to high school, assume that we flip a coin for 10 times and got 3 heads. We want to estimate the chance of getting heads

## Prior and regularization

- One major reason of introducing prior is for the sake of "regularizing" the answer
- Another coin example
    - Fall back to high school, assume that we flip a coin for 10 times and got 3 heads. We want to estimate the chance of getting heads
    - 3/10, right?
    - And if I asked you chance of getting another head in the future, you will say the chance of getting another head is 3/10

## Prior and regularization

- One major reason of introducing prior is for the sake of "regularizing" the answer
- Another coin example
  - Fall back to high school, assume that we flip a coin for 10 times and got 3 heads. We want to estimate the chance of getting heads
  - 3/10, right?
  - And if I asked you chance of getting another head in the future, you will say the chance of getting another head is 3/10
  - Now, if I actually flip the coin for 10 times and got no head, what do you expect the chance of getting a head next time?

## Prior and regularization

- One major reason of introducing prior is for the sake of "regularizing" the answer
- Another coin example
    - Fall back to high school, assume that we flip a coin for 10 times and got 3 heads. We want to estimate the chance of getting heads
    - 3/10, right?
    - And if I asked you chance of getting another head in the future, you will say the chance of getting another head is 3/10
    - Now, if I actually flip the coin for 10 times and got no head, what do you expect the chance of getting a head next time?
    - 0? Okay, the estimate is a bit extreme. We know that it is very difficult to make a coin that always gives a tail

## Prior and regularization

- One major reason of introducing prior is for the sake of "regularizing" the answer
- Another coin example
  - Fall back to high school, assume that we flip a coin for 10 times and got 3 heads. We want to estimate the chance of getting heads
  - 3/10, right?
  - And if I asked you chance of getting another head in the future, you will say the chance of getting another head is 3/10
  - Now, if I actually flip the coin for 10 times and got no head, what do you expect the chance of getting a head next time?
  - 0? Okay, the estimate is a bit extreme. We know that it is very difficult to make a coin that always gives a tail
  - How about we first assumed that we actually flipped two times and got 1 head before we did experiment? We will estimate 1/12 instead of 0/10

# Prior and regularization

- We can verify that this is exactly what we got for a Beta prior with $a = 2$ and $b = 2$.

## Prior and regularization

- We can verify that this is exactly what we got for a Beta prior with $a = 2$ and $b = 2$. Note that the posterior distribution is

  $$Beta(p|2, 2)Bin(x = 0|p, N = 10) \sim Beta(0 + a, 10 + b) = Beta(2, 12)$$

  Now, what is the MAP estimate?

## Prior and regularization

- We can verify that this is exactly what we got for a Beta prior with $a = 2$ and $b = 2$. Note that the posterior distribution is

$$Beta(p|2, 2)Bin(x = 0|p, N = 10) \sim Beta(0 + a, 10 + b) = Beta(2, 12)$$

Now, what is the MAP estimate? It should be the $p$ that maximize the posterior probability. That is the mode of $Beta(2, 12)$.

## Prior and regularization

- We can verify that this is exactly what we got for a Beta prior with $a = 2$ and $b = 2$. Note that the posterior distribution is

$$Beta(p|2,2)Bin(x = 0|p, N = 10) \sim Beta(0 + a, 10 + b) = Beta(2, 12)$$

Now, what is the MAP estimate? It should be the $p$ that maximize the posterior probability. That is the mode of $Beta(2, 12)$. Thus,

$$p_{Head}^{(MAP)} = \frac{a - 1}{a + b - 2} = \frac{1}{12}$$

## Prior and regularization

- We can verify that this is exactly what we got for a Beta prior with $a = 2$ and $b = 2$. Note that the posterior distribution is

$$Beta(p|2, 2)Bin(x = 0|p, N = 10) \sim Beta(0 + a, 10 + b) = Beta(2, 12)$$

  Now, what is the MAP estimate? It should be the $p$ that maximize the posterior probability. That is the mode of $Beta(2, 12)$. Thus,

$$p_{Head}^{(MAP)} = \frac{a - 1}{a + b - 2} = \frac{1}{12}$$

- Recall that $Beta(1, 1) = 1$ and so likelihood function is equivalent to $Beta(p|1, 1)Bin(0|p, 10) \sim Beta(1, 11)$.

## Prior and regularization

- We can verify that this is exactly what we got for a Beta prior with $a = 2$ and $b = 2$. Note that the posterior distribution is

$$Beta(p|2,2)Bin(x = 0|p, N = 10) \sim Beta(0 + a, 10 + b) = Beta(2, 12)$$

Now, what is the MAP estimate? It should be the $p$ that maximize the posterior probability. That is the mode of $Beta(2, 12)$. Thus,

$$p_{Head}^{(MAP)} = \frac{a - 1}{a + b - 2} = \frac{1}{12}$$

- Recall that $Beta(1, 1) = 1$ and so likelihood function is equivalent to $Beta(p|1, 1)Bin(0|p, 10) \sim Beta(1, 11)$. Thus the ML estimate is the mode of $Beta(1, 11) \Rightarrow p_{Head}^{(ML)} = \frac{1 - 1}{1 + 11 - 2} = \frac{0}{10} = 0$
    - This indeed is the same as our high school naïve estimate

## Bayesian estimation and regularization

- Now let's consider the Bayesian estimate. Even for the case with no prior (equivalently an uniform prior or Beta prior with $a = 1$ and $b = 1$), recall that the "posterior distribution" is $Beta(1, 11)$

## Bayesian estimation and regularization

- Now let's consider the Bayesian estimate. Even for the case with no prior (equivalently an uniform prior or Beta prior with $a = 1$ and $b = 1$), recall that the "posterior distribution" is $Beta(1, 11)$
- The Bayesian estimate should be the average $p$ summing all possibility of $p$,

## Bayesian estimation and regularization

- Now let's consider the Bayesian estimate. Even for the case with no prior (equivalently an uniform prior or Beta prior with $a = 1$ and $b = 1$), recall that the "posterior distribution" is $Beta(1, 11)$

- The Bayesian estimate should be the average $p$ summing all possibility of $p$, which is essentially just, $\int p Beta(p|1, 11) dp = E[p]$, i.e., the mean.

## Bayesian estimation and regularization

- Now let's consider the Bayesian estimate. Even for the case with no prior (equivalently an uniform prior or Beta prior with $a = 1$ and $b = 1$), recall that the "posterior distribution" is $Beta(1, 11)$

- The Bayesian estimate should be the average $p$ summing all possibility of $p$, which is essentially just, $\int p Beta(p|1, 11) dp = E[p]$, i.e., the mean. Thus

$$p_{Head}^{(Bayesian)} = \frac{a}{a + b} = \frac{1}{11}$$

## Bayesian estimation and regularization

- Now let's consider the Bayesian estimate. Even for the case with no prior (equivalently an uniform prior or Beta prior with $a = 1$ and $b = 1$), recall that the "posterior distribution" is $Beta(1, 11)$

- The Bayesian estimate should be the average $p$ summing all possibility of $p$, which is essentially just, $\int p Beta(p|1, 11) dp = E[p]$, i.e., the mean. Thus

$$p_{Head}^{(Bayesian)} = \frac{a}{a + b} = \frac{1}{11}$$

- Note that Bayesian estimation is "self-regularized" (i.e., giving less extreme results) since it inherently averages out all possible cases

# Multinomial distribution

- Binomial distribution models the probability of a binary outcome. For a random event with discrete but non-binary (more than two) outcomes, we can model the event with a multinomial distribution

## Multinomial distribution

- Binomial distribution models the probability of a binary outcome. For a random event with discrete but non-binary (more than two) outcomes, we can model the event with a multinomial distribution
- Let say the probability of each possible outcome $i$ is $p_i$. And we have conducted $N$ different experiments, let say $x_i$ is the number of times we obtain outcome $i$. Then the probability of such even is given by

## Multinomial distribution

- Binomial distribution models the probability of a binary outcome. For a random event with discrete but non-binary (more than two) outcomes, we can model the event with a multinomial distribution

- Let say the probability of each possible outcome $i$ is $p_i$. And we have conducted $N$ different experiments, let say $x_i$ is the number of times we obtain outcome $i$. Then the probability of such even is given by

$$Mult(x_1, \cdots, x_n | p_1, \cdots, p_n) = \binom{N}{x_1 x_2 \cdots x_n} p_1^{x_1} p_2^{x_2} \cdots p_n^{x_n},$$

## Multinomial distribution

- Binomial distribution models the probability of a binary outcome. For a random event with discrete but non-binary (more than two) outcomes, we can model the event with a multinomial distribution

- Let say the probability of each possible outcome $i$ is $p_i$. And we have conducted $N$ different experiments, let say $x_i$ is the number of times we obtain outcome $i$. Then the probability of such even is given by

$$Mult(x_1, \cdots, x_n | p_1, \cdots, p_n) = \binom{N}{x_1 x_2 \cdots x_n} p_1^{x_1} p_2^{x_2} \cdots p_n^{x_n},$$

- Just make sure we are in the same pace. Note that $p_1 + p_2 + \cdots + p_n = 1$ and $x_1 + x_2 + \cdots + x_n = N$

## Dirichlet distribution

- Note that the conjugate prior of multinomial distribution should take the form $x_1^{\alpha_1-1} x_2^{\alpha_2-1} \cdots x_n^{\alpha_n-1}$

## Dirichlet distribution

- Note that the conjugate prior of multinomial distribution should take the form $x_1^{\alpha_1-1} x_2^{\alpha_2-1} \cdots x_n^{\alpha_n-1}$

- It turns out that the distribution is the so-called Dirichlet distribution. Its pdf is given by

$$Dir(x_1, \cdots, x_n | \alpha_1, \cdots, \alpha_n)$$
$$= \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\cdots\Gamma(\alpha_n)} x_1^{\alpha_1-1} x_2^{\alpha_2-1} \cdots x_n^{\alpha_n-1}$$

## Dirichlet distribution

- Note that the conjugate prior of multinomial distribution should take the form $x_1^{\alpha_1-1} x_2^{\alpha_2-1} \cdots x_n^{\alpha_n-1}$

- It turns out that the distribution is the so-called Dirichlet distribution. Its pdf is given by

$$
Dir(x_1, \cdots, x_n | \alpha_1, \cdots, \alpha_n)
$$
$$
= \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\cdots\Gamma(\alpha_n)} x_1^{\alpha_1-1} x_2^{\alpha_2-1} \cdots x_n^{\alpha_n-1}
$$

- As usual since pdf should be normalized to 1, we have

$$
\int x_1^{\alpha_1-1} x_2^{\alpha_2-1} \cdots x_n^{\alpha_n-1} = \frac{\Gamma(\alpha_1)\Gamma(\alpha_2)\cdots\Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n)}
$$

# Mean, mode, variance of Dirichlet distribution

- Mean:

$$\begin{aligned} E[X_1] =& \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_n)} \int x_1^{\alpha_1} x_2^{\alpha_2-1} \cdots x_n^{\alpha_n-1} \\ =& \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_n)} \frac{\Gamma(\alpha_1 + 1)\cdots\Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n + 1)} = \frac{\alpha_1}{\alpha_1 + \cdots + \alpha_n} \end{aligned}$$

# Mean, mode, variance of Dirichlet distribution

- Mean:

$$E[X_1] = \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \int x_1^{\alpha_1} x_2^{\alpha_2 - 1} \cdots x_n^{\alpha_n - 1}$$

$$= \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \frac{\Gamma(\alpha_1 + 1) \cdots \Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n + 1)} = \frac{\alpha_1}{\alpha_1 + \cdots + \alpha_n}$$

- Similarly, $E[X_1^2] = \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \int x_1^{\alpha_1 + 1} x_2^{\alpha_2 - 1} \cdots x_n^{\alpha_n - 1} =$
  $\frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \frac{\Gamma(\alpha_1 + 2) \cdots \Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n + 2)} = \frac{(\alpha_1 + 1)\alpha_1}{(\alpha_1 + \cdots + \alpha_n + 1)(\alpha_1 + \cdots + \alpha_n)}$.

## Mean, mode, variance of Dirichlet distribution

- Mean:

$$E[X_1] = \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_n)} \int x_1^{\alpha_1} x_2^{\alpha_2 - 1} \cdots x_n^{\alpha_n - 1}$$

$$= \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_n)} \frac{\Gamma(\alpha_1 + 1)\cdots\Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n + 1)} = \frac{\alpha_1}{\alpha_1 + \cdots + \alpha_n}$$

- Similarly, $E[X_1^2] = \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_n)} \int x_1^{\alpha_1 + 1} x_2^{\alpha_2 - 1} \cdots x_n^{\alpha_n - 1} =$
  $\frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\cdots\Gamma(\alpha_n)} \frac{\Gamma(\alpha_1 + 2)\cdots\Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n + 2)} = \frac{(\alpha_1 + 1)\alpha_1}{(\alpha_1 + \cdots + \alpha_n + 1)(\alpha_1 + \cdots + \alpha_n)}$. Thus,
  $Var(X_1) = E[X_1^2] - E[X_1^2] = \frac{(\alpha_1 + 1)\alpha_1}{(\alpha_1 + \cdots + \alpha_n + 1)(\alpha_1 + \cdots + \alpha_n)} - \frac{\alpha_1^2}{(\alpha_1 + \cdots + \alpha_n)^2} =$
  $\frac{\alpha_1(\alpha_0 - \alpha_1)}{\alpha_0^2(\alpha_0 + 1)}$, where $\alpha_0 = \alpha_1 + \cdots + \alpha_n$

## Mean, mode, variance of Dirichlet distribution

- Mean:

$$E[X_1] = \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \int x_1^{\alpha_1} x_2^{\alpha_2 - 1} \cdots x_n^{\alpha_n - 1}$$

$$= \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \frac{\Gamma(\alpha_1 + 1) \cdots \Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n + 1)} = \frac{\alpha_1}{\alpha_1 + \cdots + \alpha_n}$$

- Similarly, $E[X_1^2] = \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \int x_1^{\alpha_1 + 1} x_2^{\alpha_2 - 1} \cdots x_n^{\alpha_n - 1} =$
$\frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_n)} \frac{\Gamma(\alpha_1 + 2) \cdots \Gamma(\alpha_n)}{\Gamma(\alpha_1 + \cdots + \alpha_n + 2)} = \frac{(\alpha_1 + 1)\alpha_1}{(\alpha_1 + \cdots + \alpha_n + 1)(\alpha_1 + \cdots + \alpha_n)}$. Thus,
$Var(X_1) = E[X_1^2] - E[X_1^2] = \frac{(\alpha_1 + 1)\alpha_1}{(\alpha_1 + \cdots + \alpha_n + 1)(\alpha_1 + \cdots + \alpha_n)} - \frac{\alpha_1^2}{(\alpha_1 + \cdots + \alpha_n)^2} =$
$\frac{\alpha_1(\alpha_0 - \alpha_1)}{\alpha_0^2(\alpha_0 + 1)}$, where $\alpha_0 = \alpha_1 + \cdots + \alpha_n$

- Mode: one can show that the mode of $Dir(\alpha_1, \cdots, \alpha_n)$ is

$$\frac{\alpha_i - 1}{\alpha_1 + \cdots + \alpha_n - n}.$$

We will not show it now but will leave as an exercise

# Summary of Dirichlet distribution

- Pdf:

$$Dir(\mathbf{x}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_1 + \cdots + \alpha_n)}{\Gamma(\alpha_1)\Gamma(\alpha_2) \cdots \Gamma(\alpha_n)} x_1^{\alpha_1 - 1} x_2^{\alpha_2 - 1} \cdots x_n^{\alpha_n - 1}$$

- Mean:

$$\frac{\alpha_i}{\alpha_1 + \cdots + \alpha_n}$$

- Variance:

$$\frac{\alpha_i(\alpha_0 - \alpha_i)}{\alpha_0^2(\alpha_0 + 1)}$$

- Mode:

$$\frac{\alpha_i - 1}{\alpha_1 + \cdots + \alpha_n - n}$$

# Posterior probability given Multinomial likelihood and Dirichlet prior

Upon observing $x_1, \cdots, x_n$, the posterior distribution of $p_1, \cdots, p_n$ becomes

$$p(p_1, \cdots, p_n | x_1, \cdots, x_n, \alpha_1, \cdots, \alpha_n)$$

# Posterior probability given Multinomial likelihood and Dirichlet prior

Upon observing $x_1, \cdots, x_n$, the posterior distribution of $p_1, \cdots, p_n$ becomes

$$p(p_1, \cdots, p_n | x_1, \cdots, x_n, \alpha_1, \cdots, \alpha_n)$$
$$= Const1 \cdot Dir(p_1, \cdots, p_n | \alpha_1, \cdots, \alpha_n) Mult(x_1, \cdots, x_n | p_1, \cdots, p_n)$$

# Posterior probability given Multinomial likelihood and Dirichlet prior

Upon observing $x_1, \cdots, x_n$, the posterior distribution of $p_1, \cdots, p_n$ becomes

$$p(p_1, \cdots, p_n | x_1, \cdots, x_n, \alpha_1, \cdots, \alpha_n)$$
$$= Const1 \cdot Dir(p_1, \cdots, p_n | \alpha_1, \cdots, \alpha_n) Mult(x_1, \cdots, x_n | p_1, \cdots, p_n)$$
$$= Const2 \cdot p_1^{x_1 + \alpha_1} \cdots p_n^{x_n + \alpha_n}$$
$$= Dir(p_1, \cdots, p_n | \tilde{\alpha}_1, \cdots, \tilde{\alpha}_n)$$

So the posterior distribution is Dirichlet with parameters updated to
$\tilde{\alpha}_1 \leftarrow x_1 + \alpha_1, \cdots, \tilde{\alpha}_n \leftarrow x_n + \alpha_n$

## Poisson distribution

Poisson distribution describes the number of arrival $K$ within some period. For example, one can use Poisson distribution to model the arrival process (Poisson process) of customers into a store.

## Poisson distribution

Poisson distribution describes the number of arrival $K$ within some period. For example, one can use Poisson distribution to model the arrival process (Poisson process) of customers into a store. Its pdf is given by

$$Poisson(k|\lambda T) = \frac{e^{-\lambda T}(\lambda T)^k}{k!},$$

where $k$ is a non-negative integer, $\lambda$ is rate of arrival and $T$ is the length of the observed period.

## Poisson distribution

Poisson distribution describes the number of arrival $K$ within some period. For example, one can use Poisson distribution to model the arrival process (Poisson process) of customers into a store. Its pdf is given by

$$Poisson(k|\lambda T) = \frac{e^{-\lambda T}(\lambda T)^k}{k!},$$

where $k$ is a non-negative integer, $\lambda$ is rate of arrival and $T$ is the length of the observed period. It is easy to check that (please verify)

$$Mean = \lambda T$$

$$Variance = \lambda T$$

N.B. the parameters $\lambda T$ comes as a group and so we can consider it as a single parameter

## Poisson process

Poisson process is probably the simplest random process to model event arrivals. It is based on two simple assumptions

## Poisson process

Poisson process is probably the simplest random process to model event arrivals. It is based on two simple assumptions

1. Arrival rate is invariant over time
   - That is, $\lambda$ is a constant that does not change with time

## Poisson process

Poisson process is probably the simplest random process to model event arrivals. It is based on two simple assumptions

1. Arrival rate is invariant over time
   - That is, $\lambda$ is a constant that does not change with time
2. Each arrival is independent of the other

## Poisson process

Poisson process is probably the simplest random process to model event arrivals. It is based on two simple assumptions

1. Arrival rate is invariant over time
   - That is, $\lambda$ is a constant that does not change with time
2. Each arrival is independent of the other
   - For example, even though we just have one customer coming in, the probability that the next customer to come in immediately should not decrease

## Poisson process

Poisson process is probably the simplest random process to model event arrivals. It is based on two simple assumptions

1. Arrival rate is invariant over time
   - That is, $\lambda$ is a constant that does not change with time
2. Each arrival is independent of the other
   - For example, even though we just have one customer coming in, the probability that the next customer to come in immediately should not decrease
   - It makes sense to model say customers to a department store

## Poisson process

Poisson process is probably the simplest random process to model event arrivals. It is based on two simple assumptions

1. Arrival rate is invariant over time
   - That is, $\lambda$ is a constant that does not change with time
2. Each arrival is independent of the other
   - For example, even though we just have one customer coming in, the probability that the next customer to come in immediately should not decrease
   - It makes sense to model say customers to a department store
   - It can be less perfect to model the times my car broke down. The events are likely to be related

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$.

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$.

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$.

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
  $Pr(k \text{ arrivals in } T)$

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
  $Pr(k \text{ arrivals in } T) = \binom{N}{k}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
  $Pr(k \text{ arrivals in } T) = \binom{N}{k}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$
  $= \frac{N(N-1)\cdots(N-k+1)}{k!}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
$Pr(k \text{ arrivals in } T) = \binom{N}{k}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$
$= \frac{N(N-1)\cdots(N-k+1)}{k!}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k} \approx \frac{N^k}{k!}\lambda^k\frac{T^k}{N^k}(1 - \lambda\Delta)^{N-k}$

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
  $Pr(k \text{ arrivals in } T) = \binom{N}{k}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$
  $= \frac{N(N-1)\cdots(N-k+1)}{k!}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k} \approx \frac{N^k}{k!}\lambda^k\frac{T^k}{N^k}(1 - \lambda\Delta)^{N-k}$
  $= \frac{(\lambda T)^k}{k!}(1 - \frac{\lambda T}{N})^{N-k}$

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
  $Pr(k \text{ arrivals in } T) = \binom{N}{k}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$
  $= \frac{N(N-1)\cdots(N-k+1)}{k!}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k} \approx \frac{N^k}{k!}\lambda^k\frac{T^k}{N^k}(1 - \lambda\Delta)^{N-k}$
  $= \frac{(\lambda T)^k}{k!}(1 - \frac{\lambda T}{N})^{N-k} \approx \frac{(\lambda T)^k}{k!}(1 - \frac{\lambda T}{N})^N$

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
  $Pr(k \text{ arrivals in } T) = \binom{N}{k}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$
  $= \frac{N(N-1)\cdots(N-k+1)}{k!}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k} \approx \frac{N^k}{k!}\lambda^k\frac{T^k}{N^k}(1 - \lambda\Delta)^{N-k}$
  $= \frac{(\lambda T)^k}{k!}(1 - \frac{\lambda T}{N})^{N-k} \approx \frac{(\lambda T)^k}{k!}(1 - \frac{\lambda T}{N})^N = \frac{(\lambda T)^k}{k!}\exp(-\lambda T)$,
  where we use $(1 + a/N)^N = \exp(a)$ for the last equality

## Poisson process and Poisson distribution

- Consider a period $T$ and let's the arrival rate be $\lambda$ as before. Let's partition $T$ into $N$ different very short intervals of length $\Delta$. Hence, $T = N\Delta$. We will also assume $N \to \infty$ and thus $\Delta \to 0$. The probability of getting an arrival in any interval $\Delta$ is thus $\lambda\Delta$. Moreover, since $\Delta \to 0$, the probability of getting two arrivals $\propto \Delta^2$ and is negligible compared to $\lambda\Delta$

- Then, the probability of $k$ arrivals
  $Pr(k \text{ arrivals in } T) = \binom{N}{k}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k}$
  $= \frac{N(N-1)\cdots(N-k+1)}{k!}(\lambda\Delta)^k(1 - \lambda\Delta)^{N-k} \approx \frac{N^k}{k!}\lambda^k\frac{T^k}{N^k}(1 - \lambda\Delta)^{N-k}$
  $= \frac{(\lambda T)^k}{k!}(1 - \frac{\lambda T}{N})^{N-k} \approx \frac{(\lambda T)^k}{k!}(1 - \frac{\lambda T}{N})^N = \frac{(\lambda T)^k}{k!}\exp(-\lambda T)$,
  where we use $(1 + a/N)^N = \exp(a)$ for the last equality

Note that indeed $Pr(k \text{ arrivals in } T) = Poisson(k|\lambda T)$

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$.

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$. Now,
  $Pr(\text{next event happened within in time } [t, t + \Delta])$
  $= Pr(\text{next event happened within in time } [n\Delta, (n + 1)\Delta])$

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$. Now,
  $Pr$(next event happened within in time $[t, t + \Delta]$)
  $= Pr$(next event happened within in time $[n\Delta, (n + 1)\Delta]$)
  $= Pr$(no event in first $n$ intervals)$Pr$(event happened in $n + 1$ interval)

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$. Now,
  $Pr(\text{next event happened within in time } [t, t + \Delta])$
  $= Pr(\text{next event happened within in time } [n\Delta, (n+1)\Delta])$
  $= Pr(\text{no event in first } n \text{ intervals})Pr(\text{event happened in } n+1 \text{ interval})$
  $= (1 - \lambda\Delta)^n(\lambda\Delta)$

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$. Now,
  $Pr(\text{next event happened within in time } [t, t + \Delta])$
  $= Pr(\text{next event happened within in time } [n\Delta, (n+1)\Delta])$
  $= Pr(\text{no event in first } n \text{ intervals}) Pr(\text{event happened in } n+1 \text{ interval})$
  $= (1 - \lambda\Delta)^n (\lambda\Delta)$
- Let $f_T(t)$ be the pdf of the interval time. Then,
  $f_T(t) = \frac{(1 - \lambda\Delta)^n (\lambda\Delta)}{\Delta}$

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$. Now,
  $Pr$(next event happened within in time $[t, t + \Delta]$)
  $= Pr$(next event happened within in time $[n\Delta, (n + 1)\Delta]$)
  $= Pr$(no event in first $n$ intervals)$Pr$(event happened in $n + 1$ interval)
  $= (1 - \lambda\Delta)^n(\lambda\Delta)$
- Let $f_T(t)$ be the pdf of the interval time. Then,
  $f_T(t) = \frac{(1-\lambda\Delta)^n(\lambda\Delta)}{\Delta} = \lambda(1 - \lambda\frac{t}{n})^n$

## Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$. Now,
  $Pr(\text{next event happened within in time } [t, t + \Delta])$
  $= Pr(\text{next event happened within in time } [n\Delta, (n + 1)\Delta])$
  $= Pr(\text{no event in first } n \text{ intervals})Pr(\text{event happened in } n + 1 \text{ interval})$
  $= (1 - \lambda\Delta)^n(\lambda\Delta)$
- Let $f_T(t)$ be the pdf of the interval time. Then,
  $f_T(t) = \frac{(1-\lambda\Delta)^n(\lambda\Delta)}{\Delta} = \lambda(1 - \lambda\frac{t}{n})^n = \lambda\exp(-\lambda t)$, where we use
  $(1 + a/n)^n = \exp(a)$ again for $n \to \infty$

# Interarrival time of Poisson process

Using the similar analysis, we can also easily evaluate the distribution of interarrival time, the time that the next event will happen given that an event just happened. Let $t = n\Delta$ and use the same notation as before

- Note that $t > 0$ and $\Delta \to 0$ and so $n \to \infty$. Now,
  $Pr($next event happened within in time $[t, t + \Delta])$
  $= Pr($next event happened within in time $[n\Delta, (n + 1)\Delta])$
  $= Pr($no event in first $n$ intervals$)Pr($event happened in $n + 1$ interval$)$
  $= (1 - \lambda\Delta)^n(\lambda\Delta)$
- Let $f_T(t)$ be the pdf of the interval time. Then,
  $f_T(t) = \frac{(1 - \lambda\Delta)^n(\lambda\Delta)}{\Delta} = \lambda(1 - \lambda\frac{t}{n})^n = \lambda \exp(-\lambda t)$, where we use
  $(1 + a/n)^n = \exp(a)$ again for $n \to \infty$

## Exponential distribution

$f_T(t) = \lambda \exp(-\lambda t) \triangleq Exp(t|\lambda)$ is the pdf of the exponential distribution with parameter $\lambda$. It is easy to verify that (as exercise)

- $E[T] = 1/\lambda$
- $Var(T) = 1/\lambda^2$

## Normal distribution revisit

For a univariate normal random variable, the pdf is given by

$$
\begin{aligned}
Norm(x|\mu, \sigma^2) =& \frac{1}{\sqrt{2\pi\sigma^2}} exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \\
=& \sqrt{\frac{\lambda}{2\pi}} exp\left(-\frac{\lambda(x-\mu)^2}{2}\right)
\end{aligned}
$$

with

$$
E[X|\mu, \sigma^2] = \mu,
$$

$$
E[(X-\mu)^2|\mu, \sigma^2] = \sigma^2,
$$

Recall that $\lambda = \frac{1}{\sigma^2}$ is the precision parameter that simplifies computations in many cases

# Conjugate prior of normal distribution for fixed $\sigma_2$

Consider $\sigma^2$ fixed and $\mu$ as the model parameter, then the posterior probability is given by

$$p(\mu|x; \sigma^2) \propto p(\mu, x; \sigma^2)$$

# Conjugate prior of normal distribution for fixed $\sigma_2$

Consider $\sigma^2$ fixed and $\mu$ as the model parameter, then the posterior probability is given by

$$p(\mu|x; \sigma^2) \propto p(\mu, x; \sigma^2)$$
$$= p(\mu) Norm(x|\mu; \sigma^2)$$

## Conjugate prior of normal distribution for fixed $\sigma_2$

Consider $\sigma^2$ fixed and $\mu$ as the model parameter, then the posterior probability is given by

$$p(\mu|x;\sigma^2) \propto p(\mu,x;\sigma^2)$$
$$= p(\mu)Norm(x|\mu;\sigma^2)$$
$$\propto p(\mu)exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

## Conjugate prior of normal distribution for fixed $\sigma_2$

Consider $\sigma^2$ fixed and $\mu$ as the model parameter, then the posterior probability is given by

$$p(\mu|x; \sigma^2) \propto p(\mu, x; \sigma^2)$$
$$= p(\mu)Norm(x|\mu; \sigma^2)$$
$$\propto p(\mu)exp\left(-\frac{(x - \mu)^2}{2\sigma^2}\right)$$

It is apparent that the posterior will keep the same form if $p(\mu)$ is also normal. Therefore, normal distribution is the conjugate prior of itself for fixed variance

## Posterior distribution of normal variable for fixed $\sigma^2$

Given prior $p(\mu) = Norm(\mu|\mu_0, \sigma_0^2)$ and likelihood $Norm(x|\mu; \sigma^2)$. Let's find the posterior probability,

$$p(\mu|x; \sigma^2, \mu_0, \sigma_0^2)$$
$$= Const \cdot Norm(\mu|\mu_0, \sigma_0^2) Norm(x|\mu; \sigma^2)$$

## Posterior distribution of normal variable for fixed $\sigma^2$

Given prior $p(\mu) = Norm(\mu|\mu_0, \sigma_0^2)$ and likelihood $Norm(x|\mu; \sigma^2)$. Let's find the posterior probability,

$$
\begin{aligned}
&p(\mu|x; \sigma^2, \mu_0, \sigma_0^2) \\
=&Const \cdot Norm(\mu|\mu_0, \sigma_0^2)Norm(x|\mu; \sigma^2) \\
=&Const2 \cdot exp\left(-\frac{(x-\mu)^2}{2\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) \\
=&Norm\left(\mu; \tilde{\mu}, \tilde{\sigma}^2\right),
\end{aligned}
$$

## Posterior distribution of normal variable for fixed $\sigma^2$

Given prior $p(\mu) = Norm(\mu|\mu_0, \sigma_0^2)$ and likelihood $Norm(x|\mu; \sigma^2)$. Let's find the posterior probability,

$$
\begin{aligned}
&p(\mu|x; \sigma^2, \mu_0, \sigma_0^2) \\
=& Const \cdot Norm(\mu|\mu_0, \sigma_0^2) Norm(x|\mu; \sigma^2) \\
=& Const2 \cdot exp\left(-\frac{(x - \mu)^2}{2\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}\right) \\
=& Norm\left(\mu; \tilde{\mu}, \tilde{\sigma}^2\right),
\end{aligned}
$$

where $\tilde{\mu} = \frac{\sigma_0^2 x + \mu_0 \sigma^2}{\sigma_0^2 + \sigma^2}$ and $\tilde{\sigma}^2 = \frac{\sigma_0^2 \sigma^2}{\sigma_0^2 + \sigma^2}$. Alternatively, $\tilde{\lambda} = \lambda_0 + \lambda$ and $\tilde{\mu} = \frac{\lambda}{\tilde{\lambda}} x + \frac{\lambda_0}{\tilde{\lambda}} \mu_0$. Note that we have already came across the more general expression when we studied product of multivariate normal distribution

# Conjugate prior of normal distribution for fixed $\mu$

Consider $\mu$ fixed and $\lambda$ as the model parameter

$$p(x|\lambda; \mu) \propto p(x, \lambda; \mu) = p(\lambda) Norm(x|\lambda; \mu)$$

# Conjugate prior of normal distribution for fixed $\mu$

Consider $\mu$ fixed and $\lambda$ as the model parameter

$$p(x|\lambda; \mu) \propto p(x, \lambda; \mu) = p(\lambda) Norm(x|\lambda; \mu)$$
$$\propto p(\lambda) \sqrt{\lambda} \exp\left(-\frac{\lambda(x - \mu)^2}{2}\right)$$

## Conjugate prior of normal distribution for fixed $\mu$

Consider $\mu$ fixed and $\lambda$ as the model parameter

$$p(x|\lambda;\mu) \propto p(x,\lambda;\mu) = p(\lambda)Norm(x|\lambda;\mu)$$
$$\propto p(\lambda)\sqrt{\lambda}\exp\left(-\frac{\lambda(x-\mu)^2}{2}\right)$$

More generally, when we have $N$ observations from the same source,

$$p(x_1,\cdots,x_N,\lambda;\mu) = p(\lambda)\prod_{i=1}^{N}Norm(x_i|\lambda;\mu)$$
$$\propto p(\lambda)\lambda^{\frac{N}{2}}\exp\left(-\lambda\sum_{i=1}^{N}\frac{(x_i-\mu)^2}{2}\right)$$

## Conjugate prior of normal distribution for fixed $\mu$

Consider $\mu$ fixed and $\lambda$ as the model parameter

$$p(x|\lambda; \mu) \propto p(x, \lambda; \mu) = p(\lambda) Norm(x|\lambda; \mu)$$
$$\propto p(\lambda) \sqrt{\lambda} \exp\left(-\frac{\lambda(x - \mu)^2}{2}\right)$$

More generally, when we have $N$ observations from the same source,

$$p(x_1, \cdots, x_N, \lambda; \mu) = p(\lambda) \prod_{i=1}^{N} Norm(x_i|\lambda; \mu)$$
$$\propto p(\lambda) \lambda^{\frac{N}{2}} \exp\left(-\lambda \sum_{i=1}^{N} \frac{(x_i - \mu)^2}{2}\right)$$

From inspection, the conjugate prior should have a form $\lambda^a \exp(-b\lambda)$

## Gamma distribution

The distribution with the desired form described in previous slide turns out to be the Gamma distribution. Its pdf, mean, and variance (please verify the mean and variance) are given by

$$Gamma(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} exp(-b\lambda)$$

$$E[\lambda] = \frac{a}{b}$$

$$Var[\lambda] = \frac{a}{b^2},$$

where $a, b > 0$ and $\lambda \geq 0$

# Gamma distribution

The distribution with the desired form described in previous slide turns out to be the Gamma distribution. Its pdf, mean, and variance (please verify the mean and variance) are given by

$$Gamma(\lambda|a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} exp(-b\lambda)$$

$$E[\lambda] = \frac{a}{b}$$

$$Var[\lambda] = \frac{a}{b^2},$$

where $a, b > 0$ and $\lambda \geq 0$

N.B. when $a = 1$, Gamma reduces to the exponential distribution. When $a$ is integer, it reduces to Erlang distribution

# Posterior distribution of normal variable for fixed $\mu$

Posterior probability given Normal likelihood (fixed mean) and Gamma prior

$$p(\lambda|x, a, b; \mu) = Const1 \cdot Gamma(\lambda|a, b) Norm(x|\lambda; \mu)$$

## Posterior distribution of normal variable for fixed $\mu$

Posterior probability given Normal likelihood (fixed mean) and Gamma prior

$$
\begin{aligned}
p(\lambda|x, a, b; \mu) =& Const1 \cdot Gamma(\lambda|a, b)Norm(x|\lambda; \mu) \\
=& Const2 \cdot \lambda^{a-1} \exp(-b\lambda)\sqrt{\lambda} \exp\left(-\lambda\frac{(x-\mu)^2}{2}\right) \\
=& Gamma\left(\lambda; \tilde{a}, \tilde{b}\right),
\end{aligned}
$$

where $\tilde{a} \leftarrow a + \frac{1}{2}$ and $\tilde{b} \leftarrow b + \frac{(x-\mu)^2}{2}$

## Conjugate prior summary

| Distribution | Likelihood $p(\mathbf{x}|\theta)$ | Prior $p(\theta)$ | Distribution |
|---|---|---|---|
| Bernoulli | $(1-\theta)^{(1-x)}\theta^x$ | $\propto (1-\theta)^{(a-1)}\theta^{(b-1)}$ | Beta |
| Binomial | $\propto (1-\theta)^{(N-x)}\theta^x$ | $\propto (1-\theta)^{(a-1)}\theta^{(b-1)}$ | Beta |
| Multinomial | $\propto \theta_1^{x_1}\theta_2^{x_2}\theta_3^{x_3}$ | $\propto \theta_1^{\alpha_1-1}\theta_2^{\alpha_2-1}\theta_3^{\alpha_3-1}$ | Dirichlet |
| Normal (fixed $\sigma^2$) | $\propto \exp\left(-\frac{(x-\theta)^2}{2\sigma^2}\right)$ | $\propto \exp\left(-\frac{(\theta-\mu_0)^2}{2\sigma_0^2}\right)$ | Normal |
| Normal (fixed $\mu$) | $\propto \sqrt{\theta}\exp\left(-\frac{\theta(x-\mu)^2}{2}\right)$ | $\propto \theta^{a-1}exp(-b\theta)$ | Gamma |
| Poisson | $\propto \theta^x \exp(-\theta)$ | $\propto \theta^{a-1}exp(-b\theta)$ | Gamma |

## This time...

- Bayesian Net
- Belief Propagation Algorithm
- LDPC/IRA Codes

# Bayesian Net

- Relationship of variables depicted by a directed graph with no loop
- Given a variable's parents, the variable is conditionally independent of any non-descendants
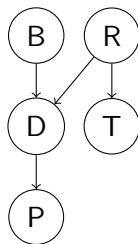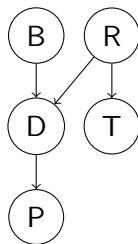- Reduce model complexity
- Facilitate easier inference

## Burlgar and racoon

Burlgar: B; Racoon: R; Dog barked: D; Police called: P; Trash can fell: T

$$p(p, d, b, t, r) = p(p|d, b, t, r)p(d|b, t, r)p(b|t, r)p(t|r)p(r)$$

## Burlgar and racoon

Burlgar: B; Racoon: R; Dog barked: D; Police called: P; Trash can fell: T

$$p(p,d,b,t,r) = p(p|d,b,t,r)p(d|b,t,r)p(b|t,r)p(t|r)p(r)$$
$$= \underbrace{p(p|d,\cancel{b},\cancel{t},\cancel{r})}_{2\ parameters}p(d|b,\cancel{t},r)p(b|\cancel{t},\cancel{r})p(t|r)p(r)$$

## Burlgar and racoon

Burlgar: B; Racoon: R; Dog barked: D; Police called: P; Trash can fell: T

$$p(p, d, b, t, r) = p(p|d, b, t, r)p(d|b, t, r)p(b|t, r)p(t|r)p(r)$$
$$= \underbrace{p(p|d, \not b, \not t, \not r)}_{2 \ parameters}p(d|b, \not t, r)p(b|\not t, \not r)p(t|r)p(r)$$

| P | D | $p(p|d)$ |
|------|------|------|
| p | $\neg d$ | 0.01 |
| p | d | 0.4 |
| $\neg p$ | $\neg d$ | 0.99 |
| $\neg p$ | d | 0.6 |

| T | R | $p(t|r)$ |
|------|------|------|
| t | $\neg r$ | 0.05 |
| t | r | 0.7 |
| $\neg t$ | $\neg r$ | 0.95 |
| $\neg t$ | r | 0.3 |

# Burlgar and racoon

Burlgar: B; Racoon: R; Dog barked: D; Police called: P; Trash can fell: T

$$p(p, d, b, t, r) = p(p|d, b, t, r)p(d|b, t, r)p(b|t, r)p(t|r)p(r)$$
$$= \underbrace{p(p|d, \not{b}, \not{t}, \not{r})}_{2 \ parameters}p(d|b, \not{t}, r)p(b|\not{t}, \not{r})p(t|r)p(r)$$

| P | D | $p(p|d)$ |
|------|------|------|
| p | $\neg d$ | 0.01 |
| p | d | 0.4 |
| $\neg p$ | $\neg d$ | 0.99 |
| $\neg p$ | d | 0.6 |

| T | R | $p(t|r)$ |
|------|------|------|
| t | $\neg r$ | 0.05 |
| t | r | 0.7 |
| $\neg t$ | $\neg r$ | 0.95 |
| $\neg t$ | r | 0.3 |

| D | B | R | $p(d|b, r)$ |
|------|------|------|------|
| d | $\neg b$ | $\neg r$ | 0.1 |
| d | $\neg b$ | r | 0.5 |
| d | b | $\neg r$ | 1 |
| d | b | r | 1 |
| $\neg d$ | $\neg b$ | $\neg r$ | 0.9 |
| $\neg d$ | $\neg b$ | r | 0.5 |
| $\neg d$ | b | $\neg r$ | 0 |
| $\neg d$ | b | r | 0 |

# Comparison of # parameters

- # parameters of complete model: $2^5 - 1 = 31$

# Comparison of # parameters

- # parameters of complete model: $2^5 - 1 = 31$
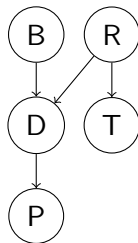- # parameters of Bayesian net:

# Comparison of # parameters

- # parameters of complete model: $2^5 - 1 = 31$
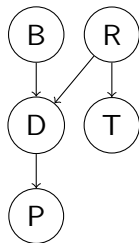- # parameters of Bayesian net:
  - $p(p|d)$: 2

# Comparison of # parameters

- # parameters of complete model: $2^5 - 1 = 31$
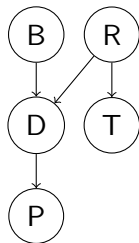- # parameters of Bayesian net:
  - $p(p|d)$: 2
  - $p(d|b, r)$: 4

# Comparison of # parameters

- # parameters of complete model: $2^5 - 1 = 31$
- # parameters of Bayesian net:
  - $p(p|d)$: 2
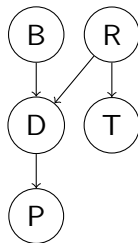  - $p(d|b,r)$: 4
  - $p(b)$: 1

## Comparison of # parameters

- # parameters of complete model: $2^5 - 1 = 31$
- # parameters of Bayesian net:
  - $p(p|d)$: 2
  - $p(d|b, r)$: 4
  - $p(b)$: 1
  - $p(t|r)$: 2

## Comparison of # parameters

- # parameters of complete model: $2^5 - 1 = 31$
- # parameters of Bayesian net:
  - $p(p|d)$: 2
  - $p(d|b, r)$: 4
  - $p(b)$: 1
  - $p(t|r)$: 2
  - $p(r)$: 1
  - Total: $2 + 4 + 1 + 2 + 1 = 10$
- The model size reduces to less than $\frac{1}{3}$!

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

Let $p(r) = 0.2$ and $p(b) = 0.01$

| D | B | R | $p(d|b,r)$ |
|------|----------|----------|------|
| $d$ | $\neg b$ | $\neg r$ | 0.1 |
| $d$ | $\neg b$ | $r$ | 0.5 |
| $d$ | $b$ | $\neg r$ | 1 |
| $d$ | $b$ | $r$ | 1 |
| $\neg d$ | $\neg b$ | $\neg r$ | 0.9 |
| $\neg d$ | $\neg b$ | $r$ | 0.5 |
| $\neg d$ | $b$ | $\neg r$ | 0 |
| $\neg d$ | $b$ | $r$ | 0 |

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

Let $p(r) = 0.2$ and $p(b) = 0.01$

| D | B | R | $p(d|b,r)$ |
|------|----------|----------|------|
| d | ¬b | ¬r | 0.1 |
| d | ¬b | r | 0.5 |
| d | b | ¬r | 1 |
| d | b | r | 1 |
| ¬d | ¬b | ¬r | 0.9 |
| ¬d | ¬b | r | 0.5 |
| ¬d | b | ¬r | 0 |
| ¬d | b | r | 0 |

$\Rightarrow$

| D | B | R | $p(d,b,r)$ |
|------|----------|----------|--------|
| d | ¬b | ¬r | 0.0792 |
| d | ¬b | r | 0.099 |
| d | b | ¬r | 0.008 |
| d | b | r | 0.002 |
| ¬d | ¬b | ¬r | 0.7128 |
| ¬d | ¬b | r | 0.099 |
| ¬d | b | ¬r | 0 |
| ¬d | b | r | 0 |

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

| P | D | $p(p\|d)$ |
|---|---|---|
| $p$ | $\neg d$ | 0.01 |
| $p$ | $d$ | 0.4 |
| $\neg p$ | $\neg d$ | 0.99 |
| $\neg p$ | $d$ | 0.6 |

| P | D | B | R | $p(d, b, r, p)$ |
|---|---|---|---|---|
| $p$ | $d$ | $\neg b$ | $\neg r$ | 0.0792 |
| $p$ | $d$ | $\neg b$ | $r$ | 0.099 |
| $p$ | $d$ | $b$ | $\neg r$ | 0.008 |
| $p$ | $d$ | $b$ | $r$ | 0.002 |
| $p$ | $\neg d$ | $\neg b$ | $\neg r$ | 0.7128 |
| $p$ | $\neg d$ | $\neg b$ | $r$ | 0.099 |
| $p$ | $\neg d$ | $b$ | $\neg r$ | 0 |
| $p$ | $\neg d$ | $b$ | $r$ | 0 |
| | | | $\cdots$ | |

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

| P | D | $p(p \mid d)$ |
|---|---|---|
| p | ¬d | 0.01 |
| p | d | 0.4 |
| ¬p | ¬d | 0.99 |
| ¬p | d | 0.6 |

| P | D | B | R | $p(d, b, r, p)$ |
|---|---|---|---|---|
| p | d | ¬b | ¬r | 0.0792 |
| p | d | ¬b | r | 0.099 |
| p | d | b | ¬r | 0.008 |
| p | d | b | r | 0.002 |
| p | ¬d | ¬b | ¬r | 0.007128 |
| p | ¬d | ¬b | r | 0.00099 |
| p | ¬d | b | ¬r | 0 |
| p | ¬d | b | r | 0 |
| | | | ... | |

# Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

| P | D | $p(p|d)$ |
|---|---|---|
| $p$ | $\neg d$ | 0.01 |
| $p$ | $d$ | 0.4 |
| $\neg p$ | $\neg d$ | 0.99 |
| $\neg p$ | $d$ | 0.6 |

| P | D | B | R | $p(d,b,r,p)$ |
|---|---|---|---|---|
| $p$ | $d$ | $\neg b$ | $\neg r$ | 0.03168 |
| $p$ | $d$ | $\neg b$ | $r$ | 0.0396 |
| $p$ | $d$ | $b$ | $\neg r$ | 0.0032 |
| $p$ | $d$ | $b$ | $r$ | 0.0008 |
| $p$ | $\neg d$ | $\neg b$ | $\neg r$ | 0.007128 |
| $p$ | $\neg d$ | $\neg b$ | $r$ | 0.00099 |
| $p$ | $\neg d$ | $b$ | $\neg r$ | 0 |
| $p$ | $\neg d$ | $b$ | $r$ | 0 |
| | | | $\cdots$ | |

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

| T | R | $p(t|r)$ |
|------|------|------|
| t | $\neg r$ | 0.05 |
| t | r | 0.7 |
| $\neg t$ | $\neg r$ | 0.95 |
| $\neg t$ | r | 0.3 |

| T | P | D | B | R | $p(d, b, r, p, t)$ |
|------|------|------|------|------|------|
| $\neg t$ | p | d | $\neg b$ | $\neg r$ | 0.03168 |
| $\neg t$ | p | d | $\neg b$ | r | 0.0396 |
| $\neg t$ | p | d | b | $\neg r$ | 0.0032 |
| $\neg t$ | p | d | b | r | 0.0008 |
| $\neg t$ | p | $\neg d$ | $\neg b$ | $\neg r$ | 0.007128 |
| $\neg t$ | p | $\neg d$ | $\neg b$ | r | 0.00099 |
| $\neg t$ | p | $\neg d$ | b | $\neg r$ | 0 |
| $\neg t$ | p | $\neg d$ | b | r | 0 |
| | | | $\cdots$ | | |

# Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

| T | R | $p(t|r)$ |
|---|---|---|
| t | $\neg r$ | 0.05 |
| t | r | 0.7 |
| $\neg t$ | $\neg r$ | 0.95 |
| $\neg t$ | r | 0.3 |

| T | P | D | B | R | $p(d,b,r,p,t)$ |
|---|---|---|---|---|---|
| $\neg t$ | p | d | $\neg b$ | $\neg r$ | 0.030096 |
| $\neg t$ | p | d | $\neg b$ | r | 0.0396 |
| $\neg t$ | p | d | b | $\neg r$ | 0.00304 |
| $\neg t$ | p | d | b | r | 0.0008 |
| $\neg t$ | p | $\neg d$ | $\neg b$ | $\neg r$ | 0.0067716 |
| $\neg t$ | p | $\neg d$ | $\neg b$ | r | 0.00099 |
| $\neg t$ | p | $\neg d$ | b | $\neg r$ | 0 |
| $\neg t$ | p | $\neg d$ | b | r | 0 |
| $\cdots$ | | | | | |

# Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

| $T$ | $R$ | $p(t\|r)$ |
|-----|-----|-----------|
| $t$ | $\neg r$ | 0.05 |
| $t$ | $r$ | 0.7 |
| $\neg t$ | $\neg r$ | 0.95 |
| $\neg t$ | $r$ | 0.3 |

| $T$ | $P$ | $D$ | $B$ | $R$ | $p(d, b, r, p, t)$ |
|-----|-----|-----|-----|-----|--------------------|
| $\neg t$ | $p$ | $d$ | $\neg b$ | $\neg r$ | 0.030096 |
| $\neg t$ | $p$ | $d$ | $\neg b$ | $r$ | 0.01188 |
| $\neg t$ | $p$ | $d$ | $b$ | $\neg r$ | 0.00304 |
| $\neg t$ | $p$ | $d$ | $b$ | $r$ | 0.00024 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $\neg r$ | 0.0067716 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $r$ | 0.000297 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $\neg r$ | 0 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $r$ | 0 |
| | | | $\cdots$ | | |

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

Normalize...

| $T$ | $P$ | $D$ | $B$ | $R$ | $p(d, b, r, p)$ |
|------|-----|------|------|------|------------------|
| $\neg t$ | $p$ | $d$ | $\neg b$ | $\neg r$ | 0.030096 |
| $\neg t$ | $p$ | $d$ | $\neg b$ | $r$ | 0.01188 |
| $\neg t$ | $p$ | $d$ | $b$ | $\neg r$ | 0.00304 |
| $\neg t$ | $p$ | $d$ | $b$ | $r$ | 0.00024 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $\neg r$ | 0.0067716 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $r$ | 0.000297 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $\neg r$ | 0 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $r$ | 0 |
| | | | $\cdots$ | | |

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

Normalize...

| $T$ | $P$ | $D$ | $B$ | $R$ | $p(d, b, r, p)$ |
|------|-----|-------|--------|--------|-----------------|
| $\neg t$ | $p$ | $d$ | $\neg b$ | $\neg r$ | 0.57518 |
| $\neg t$ | $p$ | $d$ | $\neg b$ | $r$ | 0.22704 |
| $\neg t$ | $p$ | $d$ | $b$ | $\neg r$ | 0.058099 |
| $\neg t$ | $p$ | $d$ | $b$ | $r$ | 0.0045868 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $\neg r$ | 0.12942 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $r$ | 0.0056761 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $\neg r$ | 0 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $r$ | 0 |
| | | | | $\cdots$ | |

## Burglar and racoon

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

$p(b|\neg t, p)$

$= 0.058099 + 0.0045868$

$\approx 0.0626$

| T | P | D | B | R | $p(d, b, r, p)$ |
|------|---|----------|----------|----------|----------|
| $\neg t$ | $p$ | $d$ | $\neg b$ | $\neg r$ | 0.57518 |
| $\neg t$ | $p$ | $d$ | $\neg b$ | $r$ | 0.22704 |
| $\neg t$ | $p$ | $d$ | $b$ | $\neg r$ | 0.058099 |
| $\neg t$ | $p$ | $d$ | $b$ | $r$ | 0.0045868 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $\neg r$ | 0.12942 |
| $\neg t$ | $p$ | $\neg d$ | $\neg b$ | $r$ | 0.0056761 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $\neg r$ | 0 |
| $\neg t$ | $p$ | $\neg d$ | $b$ | $r$ | 0 |
| $\cdots$ | | | | | |

# Belief Propagation Algorithm

- It is also known to be the sum-product algorithm
- The goal of belief propagation is to efficiently compute the marginal distribution out of the joint distribution of multiple variables. This is essential for inferring the outcome of a particular variable with insufficient information
- The belief propagation algorithm is usually applied to problems modeled by a undirected graph (Markov random field) or a factor graph
- Rather than giving a rigorous proof of the algorithm, we will provide a simple example to illustrate the basic idea

# Factor Graph

- A factor graph is a bipartite graph describing the correlation among several random variables. It generally contains two different types of nodes in the graph: variable nodes and factor nodes
- A variable node that is usually shown as circles corresponds to a random variable
- A factor node that is usually shown as a square connects variable nodes whose corresponding variables are immediately related

## An Example

- A factor graph example is shown below. We have 8 *discrete* random variables, $x_1^4$ and $z_1^4$, depicted by 8 variable nodes

## An Example

- A factor graph example is shown below. We have 8 *discrete* random variables, $x_1^4$ and $z_1^4$, depicted by 8 variable nodes
- Among the variable nodes, random variables $x_1^4$ (indicated by light circles) are unknown and variables $z_1^4$ (indicated by dark circles) are observed with known outcomes $\tilde{z}_1^4$

## An Example

- A factor graph example is shown below. We have 8 *discrete* random variables, $x_1^4$ and $z_1^4$, depicted by 8 variable nodes
- Among the variable nodes, random variables $x_1^4$ (indicated by light circles) are unknown and variables $z_1^4$ (indicated by dark circles) are observed with known outcomes $\tilde{z}_1^4$
- The relationships among variables are captured entirely by the figure. For example, given $x_1^4$, $z_1$, $z_2$, $z_3$, and $z_4$ are conditional independent of each other. Moreover, $(x_3, x_4)$ are conditional independent of $x_1$ given $x_2$

- The joint probability $p(x^4, z^4)$ of all variables can be decomposed into factor functions with subsets of all variables as arguments in the following

$$p(x^4, z^4) = p(x^4)p(z_1|x_1)p(z_2|x_2)p(z_3|x_3)p(z_4|x_4)$$

- Note that each factor function corresponds to a factor node in the factor graph.
- The arguments of the factor function correspond to the variable nodes that the factor node connects to.

- The joint probability $p(x^4, z^4)$ of all variables can be decomposed into factor functions with subsets of all variables as arguments in the following

$$
\begin{aligned}
p(x^4, z^4) &= p(x^4)p(z_1|x_1)p(z_2|x_2)p(z_3|x_3)p(z_4|x_4) \\
&= \underbrace{p(x_1, x_2)}_{f_b(x_1, x_2)} \underbrace{p(x_3, x_4|x_2)}_{f_d(x_2, x_3, x_4)} \underbrace{p(z_3|x_3)}_{f_e(x_3, z_3)} \underbrace{p(z_1|x_1)}_{f_a(x_1, z_1)} \underbrace{p(z_4|x_4)}_{f_f(x_4, z_4)} \underbrace{p(z_2|x_2)}_{f_c(x_2, z_2)}
\end{aligned}
$$

- Note that each factor function corresponds to a factor node in the factor graph.

- The arguments of the factor function correspond to the variable nodes that the factor node connects to.

- The joint probability $p(x^4, z^4)$ of all variables can be decomposed into factor functions with subsets of all variables as arguments in the following

$$
\begin{aligned}
p(x^4, z^4) &= p(x^4)p(z_1|x_1)p(z_2|x_2)p(z_3|x_3)p(z_4|x_4) \\
&= \underbrace{p(x_1, x_2)}_{f_b(x_1, x_2)} \underbrace{p(x_3, x_4|x_2)}_{f_d(x_2, x_3, x_4)} \underbrace{p(z_3|x_3)}_{f_e(x_3, z_3)} \underbrace{p(z_1|x_1)}_{f_a(x_1, z_1)} \underbrace{p(z_4|x_4)}_{f_f(x_4, z_4)} \underbrace{p(z_2|x_2)}_{f_c(x_2, z_2)} \\
&= f_b(x_1, x_2)f_d(x_2, x_3, x_4)f_e(x_3, z_3)f_a(x_1, z_1)f_f(x_4, z_4)f_c(x_2, z_2)
\end{aligned}
$$

- Note that each factor function corresponds to a factor node in the factor graph.

- The arguments of the factor function correspond to the variable nodes that the factor node connects to.

One common problem in probability inference is to estimate the value of a variable given incomplete information. For example, we may want to estimate $x_1$ given $z^4$ as $\tilde{z}^4$. The optimum estimate $\hat{x}_1$ will satisfy

$$\hat{x}_1 = \arg\max_{x_1} p(x_1|\tilde{z}^4) = \arg\max_{x_1} \frac{p(x_1, \tilde{z}^4)}{p(\tilde{z}^4)} = \arg\max_{x_1} p(x_1, \tilde{z}^4).$$

This requires us to compute the marginal distribution $p(x_1, \tilde{z}^4)$ out of the joint probability $p(x^4, \tilde{z}^4)$. Note that

$$p(x_1, \tilde{z}^4) = \sum_{x_2^4} p(x^4, \tilde{z}^4)$$

One common problem in probability inference is to estimate the value of a variable given incomplete information. For example, we may want to estimate $x_1$ given $z^4$ as $\tilde{z}^4$. The optimum estimate $\hat{x}_1$ will satisfy

$$\hat{x}_1 = \arg \max_{x_1} p(x_1 | \tilde{z}^4) = \arg \max_{x_1} \frac{p(x_1, \tilde{z}^4)}{p(\tilde{z}^4)} = \arg \max_{x_1} p(x_1, \tilde{z}^4).$$

This requires us to compute the marginal distribution $p(x_1, \tilde{z}^4)$ out of the joint probability $p(x^4, \tilde{z}^4)$. Note that

$$p(x_1, \tilde{z}^4) = \sum_{x_2^4} p(x^4, \tilde{z}^4)$$

$$= \sum_{x_2^4} f_a(x_1, \tilde{z}_1) f_b(x_1, x_2) f_c(x_2, \tilde{z}_2) f_d(x_2, x_3, x_4) f_e(x_3, \tilde{z}_3) f_f(x_4, \tilde{z}_4)$$

One common problem in probability inference is to estimate the value of a variable given incomplete information. For example, we may want to estimate $x_1$ given $z^4$ as $\tilde{z}^4$. The optimum estimate $\hat{x}_1$ will satisfy

$$\hat{x}_1 = \arg\max_{x_1} p(x_1|\tilde{z}^4) = \arg\max_{x_1} \frac{p(x_1, \tilde{z}^4)}{p(\tilde{z}^4)} = \arg\max_{x_1} p(x_1, \tilde{z}^4).$$

This requires us to compute the marginal distribution $p(x_1, \tilde{z}^4)$ out of the joint probability $p(x^4, \tilde{z}^4)$. Note that

$$p(x_1, \tilde{z}^4) = \sum_{x_2^4} p(x^4, \tilde{z}^4)$$

$$= \sum_{x_2^4} f_a(x_1, \tilde{z}_1) f_b(x_1, x_2) f_c(x_2, \tilde{z}_2) f_d(x_2, x_3, x_4) f_e(x_3, \tilde{z}_3) f_f(x_4, \tilde{z}_4)$$

$$= \underbrace{f_a(x_1, \tilde{z}_1)}_{m_{a1}} \underbrace{\sum_{x_2} f_b(x_1, x_2) \underbrace{f_c(x_2, \tilde{z}_2)}_{m_{c2}} \underbrace{\sum_{x_3, x_4} f_d(x_2, x_3, x_4) \underbrace{f_e(x_3, \tilde{z}_3)}_{m_{3d}} \underbrace{f_f(x_4, \tilde{z}_4)}_{m_{4d}}}_{m_{d2}}}_{\substack{m_{2b} \\ m_{b1}}}$$

We can see from the last equation that the joint probability can be computed by combining a sequence of messages passing from a variable node $i$ to a factor node $a$ ($m_{ia}$) and vice versa ($m_{ai}$). More precisely, we can write

$$m_{a1}(x_1) \leftarrow f_a(x_1, \tilde{z}_1) = \sum_{z_1} f_a(x_1, z_1) \underbrace{p(z_1)}_{m_{1a}},$$

$$m_{c2}(x_2) \leftarrow f_c(x_2, \tilde{z}_2) = \sum_{z_2} f_c(x_2, z_2) \underbrace{p(z_2)}_{m_{2c}},$$

$$m_{e3}(x_3) \leftarrow f_e(x_3, \tilde{z}_3) = \sum_{z_3} f_e(x_3, z_3) \underbrace{p(z_3)}_{m_{3e}},$$

$$m_{f4}(x_4) \leftarrow f_f(x_4, \tilde{z}_4) = \sum_{z_4} f_f(x_4, z_4) \underbrace{p(z_4)}_{m_{4f}},$$

where $p(z_i) = \begin{cases} 1, & z_i = \tilde{z}_i \\ 0, & \text{otherwise} \end{cases}$

$$p(x_1, \tilde{z}^4) = \underbrace{f_a(x_1, \tilde{z}_1)}_{m_{a1}} \underbrace{\sum_{x_2} f_b(x_1, x_2) \underbrace{f_c(x_2, \tilde{z}_2)}_{m_{c2}} \underbrace{\sum_{x_3, x_4} f_d(x_2, x_3, x_4) \underbrace{f_e(x_3, \tilde{z}_3)}_{m_{3d}} \underbrace{f_f(x_4, \tilde{z}_4)}_{m_{4d}}}_{m_{d2}}}_{\substack{m_{2b} \\ m_{b1}}} \quad (1)$$

$$m_{3d}(x_3) \leftarrow m_{e3}(x_3) = f_e(x_3, \tilde{z}_3),$$
$$m_{4d}(x_4) \leftarrow m_{f4}(x_4) = f_f(x_4, \tilde{z}_4),$$

$$p(x_1, \tilde{z}^4) = \underbrace{f_a(x_1, \tilde{z}_1)}_{m_{a1}} \underbrace{\sum_{x_2} f_b(x_1, x_2) \underbrace{f_c(x_2, \tilde{z}_2)}_{m_{c2}} \underbrace{\sum_{x_3, x_4} f_d(x_2, x_3, x_4) \underbrace{f_e(x_3, \tilde{z}_3)}_{m_{3d}} \underbrace{f_f(x_4, \tilde{z}_4)}_{m_{4d}}}_{m_{d2}}}_{\substack{m_{2b} \\ m_{b1}}} \quad (1)$$

$$m_{3d}(x_3) \leftarrow m_{e3}(x_3) = f_e(x_3, \tilde{z}_3),$$
$$m_{4d}(x_4) \leftarrow m_{f4}(x_4) = f_f(x_4, \tilde{z}_4),$$
$$m_{d2}(x_2) \leftarrow \sum_{x_3, x_4} f_d(x_2, x_3, x_4) m_{3d}(x_3) m_{4d}(x_4),$$

$$p(x_1, \tilde{z}^4) = \underbrace{f_a(x_1, \tilde{z}_1)}_{m_{a1}} \underbrace{\sum_{x_2} f_b(x_1, x_2) \underbrace{f_c(x_2, \tilde{z}_2)}_{m_{c2}} \underbrace{\sum_{x_3, x_4} f_d(x_2, x_3, x_4) \underbrace{f_e(x_3, \tilde{z}_3)}_{m_{3d}} \underbrace{f_f(x_4, \tilde{z}_4)}_{m_{4d}}}_{\substack{m_{d2} \\ m_{2b}}}}_{m_{b1}} \quad (1)$$

$$m_{3d}(x_3) \leftarrow m_{e3}(x_3) = f_e(x_3, \tilde{z}_3),$$
$$m_{4d}(x_4) \leftarrow m_{f4}(x_4) = f_f(x_4, \tilde{z}_4),$$
$$m_{d2}(x_2) \leftarrow \sum_{x_3, x_4} f_d(x_2, x_3, x_4) m_{3d}(x_3) m_{4d}(x_4),$$
$$m_{2b}(x_2) \leftarrow m_{c2}(x_2) m_{d2}(x_2),$$

$$p(x_1, \tilde{z}^4) = \underbrace{f_a(x_1, \tilde{z}_1)}_{m_{a1}} \sum_{x_2} f_b(x_1, x_2) \underbrace{f_c(x_2, \tilde{z}_2)}_{m_{c2}} \sum_{x_3, x_4} f_d(x_2, x_3, x_4) \underbrace{f_e(x_3, \tilde{z}_3)}_{m_{3d}} \underbrace{f_f(x_4, \tilde{z}_4)}_{m_{4d}} \quad (1)$$

$$\underbrace{\phantom{f_d(x_2, x_3, x_4) f_e(x_3, \tilde{z}_3) f_f(x_4, \tilde{z}_4)}}_{m_{d2}}$$

$$\underbrace{\phantom{f_c(x_2, \tilde{z}_2) f_d(x_2, x_3, x_4) f_e(x_3, \tilde{z}_3) f_f(x_4, \tilde{z}_4)}}_{\color{red}m_{2b}}$$

$$\underbrace{\phantom{f_b(x_1, x_2) f_c(x_2, \tilde{z}_2) f_d(x_2, x_3, x_4) f_e(x_3, \tilde{z}_3) f_f(x_4, \tilde{z}_4)}}_{m_{b1}}$$

$$m_{3d}(x_3) \leftarrow m_{e3}(x_3) = f_e(x_3, \tilde{z}_3),$$
$$m_{4d}(x_4) \leftarrow m_{f4}(x_4) = f_f(x_4, \tilde{z}_4),$$
$$m_{d2}(x_2) \leftarrow \sum_{x_3, x_4} f_d(x_2, x_3, x_4) m_{3d}(x_3) m_{4d}(x_4),$$
$$m_{2b}(x_2) \leftarrow m_{c2}(x_2) m_{d2}(x_2),$$
$$m_{b1}(x_1) \leftarrow \sum_{x_2} f_b(x_1, x_2) m_{2b}(x_2),$$

$$p(x_1, \tilde{z}^4) = \underbrace{f_a(x_1, \tilde{z}_1)}_{m_{a1}} \sum_{x_2} f_b(x_1, x_2) \underbrace{f_c(x_2, \tilde{z}_2)}_{m_{c2}} \sum_{x_3, x_4} f_d(x_2, x_3, x_4) \underbrace{f_e(x_3, \tilde{z}_3)}_{m_{3d}} \underbrace{f_f(x_4, \tilde{z}_4)}_{m_{4d}} \quad (1)$$

with braces indicating $m_{d2}$, $m_{2b}$, and $m_{b1}$.

$$m_{3d}(x_3) \leftarrow m_{e3}(x_3) = f_e(x_3, \tilde{z}_3),$$

$$m_{4d}(x_4) \leftarrow m_{f4}(x_4) = f_f(x_4, \tilde{z}_4),$$

$$m_{d2}(x_2) \leftarrow \sum_{x_3, x_4} f_d(x_2, x_3, x_4) m_{3d}(x_3) m_{4d}(x_4),$$

$$m_{2b}(x_2) \leftarrow m_{c2}(x_2) m_{d2}(x_2),$$

$$m_{b1}(x_1) \leftarrow \sum_{x_2} f_b(x_1, x_2) m_{2b}(x_2),$$

$$p(x_1, \tilde{z}^4) \leftarrow m_{a1}(x_1) m_{b1}(x_1),$$

$$p(x_1, \tilde{z}^4) = \underbrace{f_a(x_1, \tilde{z}_1)}_{m_{a1}} \sum_{x_2} f_b(x_1, x_2) \underbrace{f_c(x_2, \tilde{z}_2)}_{m_{c2}} \underbrace{\sum_{x_3, x_4} f_d(x_2, x_3, x_4) \underbrace{f_e(x_3, \tilde{z}_3)}_{m_{3d}} \underbrace{f_f(x_4, \tilde{z}_4)}_{m_{4d}}}_{m_{d2}} \quad (1)$$

$$\underbrace{\phantom{\sum}}_{m_{2b}}$$
$$\underbrace{\phantom{\sum}}_{m_{b1}}$$

# Belief propagation algorithm

- **Initialization**: For any variable node $i$, if the prior probability of $x_i$ is known and equal to $p(x_i)$, for $a \in N(i)$,

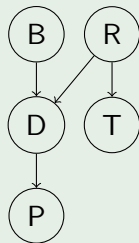- **Message passing**:

- **Belief update**:

- **Stopping criteria**: repeat message update and/or belief update until the algorithm stops when maximum number of iterations is reached or some other conditions are satisfied.

# Belief propagation algorithm

- **Initialization**: For any variable node $i$, if the prior probability of $x_i$ is known and equal to $p(x_i)$, for $a \in N(i)$,

$$m_{ia}(x_i) \leftarrow p(x_i)$$

- **Message passing**:

- **Belief update**:

- **Stopping criteria**: repeat message update and/or belief update until the algorithm stops when maximum number of iterations is reached or some other conditions are satisfied.

# Belief propagation algorithm

- **Initialization**: For any variable node $i$, if the prior probability of $x_i$ is known and equal to $p(x_i)$, for $a \in N(i)$,

$$m_{ia}(x_i) \leftarrow p(x_i)$$

- **Message passing**:

$$m_{ia}(x_i) \leftarrow \prod_{b \in N(i) \setminus a} m_{bi}(x_i),$$

$$m_{ai}(x_i) \leftarrow \sum_{\mathbf{x}_a} f_a(\mathbf{x}_a) \prod_{j \in N(a) \setminus i} m_{ja}(x_j) \qquad (\text{``sum-product''})$$

- **Belief update**:

- **Stopping criteria**: repeat message update and/or belief update until the algorithm stops when maximum number of iterations is reached or some other conditions are satisfied.

# Belief propagation algorithm

- **Initialization**: For any variable node $i$, if the prior probability of $x_i$ is known and equal to $p(x_i)$, for $a \in N(i)$,

$$m_{ia}(x_i) \leftarrow p(x_i)$$

- **Message passing**:

$$m_{ia}(x_i) \leftarrow \prod_{b \in N(i) \setminus a} m_{bi}(x_i),$$

$$m_{ai}(x_i) \leftarrow \sum_{\mathbf{x}_a} f_a(\mathbf{x}_a) \prod_{j \in N(a) \setminus i} m_{ja}(x_j) \qquad (\text{"sum-product"})$$

- **Belief update**:

$$\beta_i(x_i) \leftarrow \prod_{a \in N(i)} m_{ai}(x_i)$$

- **Stopping criteria**: repeat message update and/or belief update until the algorithm stops when maximum number of iterations is reached or some other conditions are satisfied.

# Remark

- We have not assumed the precise phyical meanings of the factor functions themselves. The only assumption we made is that the joint probability can be decomposed into the factor functions and apparently this decomposition is not unique

- The belief propagation algorithm as shown above is exact only because the corresponding graph is a tree and has no loop. If loop exists, the algorithm is not exact and generally the final belief may not even converge

- While the result is no longer exact, applying BP algorithm for general graphs (sometimes refer to as loopy BP) works well in many applications such as LDPC decoding

# Burglar and racoon revisit

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?

## Burglar and racoon revisit

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?



Moralization...

# Burglar and racoon revisit

Question: What is the probability of a burglar visit if police was called but trash can stayed untouched?



Convert to factor graph..

# Using belief propagation...

$$\begin{cases} f_P(p) & = 1 \\ f_P(\neg p) & = 0 \end{cases}$$

$$\begin{cases} f_T(t) & = 0 \\ f_T(\neg t) & = 1 \end{cases}$$

$$f_{B,D,R}(b,d,r) = p(b,d,r)$$
$$f_{T,R}(t,r) = p(t|r)$$
$$f_{D,P}(d,p) = p(p|d)$$

# Some History of LDPC Codes

- Before 1990's, the strategy for channel code has always been looking for codes that can be decoded optimally. This leads to a wide range of so-called algebraic codes. It turns out the "optimally-decodable" codes are usually poor codes

- Until early 1990's, researchers had basically agreed that the Shannon capacity was restricted to theoretical interest and could hardly be reached in practice

- The introduction of turbo codes gave a huge shock to the research community. The community were so dubious about the amazing performance of turbo codes that they did not accept the finding initially until independent researchers had verified the results

- The low-density parity-check (LDPC) codes were later rediscovered and both LDPC codes and turbo codes are based on the same philosophy differs from codes in the past. Instead of designing and using codes that can be decoded "optimally", let us just pick some *random* codes and perform decoding "sub-optimally"

# LDPC Codes

- As its name suggests, LDPC codes refer to codes that with sparse (low-density) parity check matrices. In other words, there are only few ones in a parity check matrix and the rest are all zeros

# LDPC Codes

- As its name suggests, LDPC codes refer to codes that with sparse (low-density) parity check matrices. In other words, there are only few ones in a parity check matrix and the rest are all zeros

- We learn from the proof of Channel Coding Theorem that random code is asymptotically optimum. This suggests that if we just generate a code randomly with a very long code length. It is likely that we will get a very good code.

# LDPC Codes

- As its name suggests, LDPC codes refer to codes that with sparse (low-density) parity check matrices. In other words, there are only few ones in a parity check matrix and the rest are all zeros

- We learn from the proof of Channel Coding Theorem that random code is asymptotically optimum. This suggests that if we just generate a code randomly with a very long code length. It is likely that we will get a very good code.

- The problem is: how do we perform decoding? Due to the lack of structure of a random code, tricks that enable fast decoding for structured algebraic codes that were widely used before 1990's are unrealizable here

- Solution: Belief propagation!

## Tanner Graph

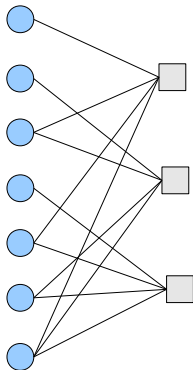- An LDPC code can be represented using a Tanner graph as shown on the right

## Tanner Graph

- An LDPC code can be represented using a Tanner graph as shown on the right
- Each circle $x_i$ represents a code bit sent to the decoder
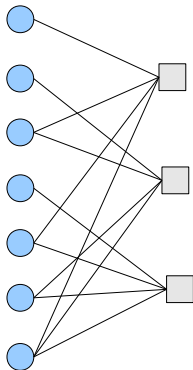
## Tanner Graph

- An LDPC code can be represented using a Tanner graph as shown on the right
- Each circle $x_i$ represents a code bit sent to the decoder
- Each square represents a check bit with value equal to the sum of code bit connecting to it
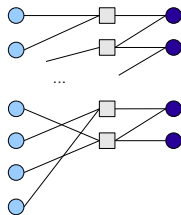
## Tanner Graph

- An LDPC code can be represented using a Tanner graph as shown on the right
- Each circle $x_i$ represents a code bit sent to the decoder
- Each square represents a check bit with value equal to the sum of code bit connecting to it
- The vector $x_1, x_2, \cdots, x_N$ is a codeword only if all checks are zero

## Tanner Graph

- An LDPC code can be represented using a Tanner graph as shown on the right
- Each circle $x_i$ represents a code bit sent to the decoder
- Each square represents a check bit with value equal to the sum of code bit connecting to it
- The vector $x_1, x_2, \cdots, x_N$ is a codeword only if all checks are zero
- By default, the mapping between a codeword to the actual message is non-trivial for an LDPC code

## Tanner Graph

- An LDPC code can be represented using a Tanner graph as shown on the right
- Each circle $x_i$ represents a code bit sent to the decoder
- Each square represents a check bit with value equal to the sum of code bit connecting to it
- The vector $x_1, x_2, \cdots, x_N$ is a codeword only if all checks are zero
- By default, the mapping between a codeword to the actual message is non-trivial for an LDPC code
- It would be great if the actual message is included in the codeword. That is, some of the bits in the codeword spell out the actual message $\Rightarrow$ IRA codes
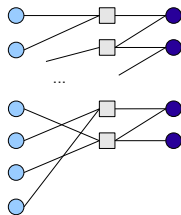
# IRA Codes

- Irregular repeated accumulate (IRA) code a type of systematic LDPC code, i.e., each codeword can be partitioned into message bits and syndrome bits
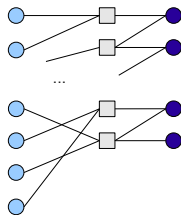
# IRA Codes

- Irregular repeated accumulate (IRA) code a type of systematic LDPC code, i.e., each codeword can be partitioned into message bits and syndrome bits
- As shown on the right, light blue circles correspond to the input message bits and the dark blue circle correspond to the syndrome bits

# IRA Codes

- Irregular repeated accumulate (IRA) code a type of systematic LDPC code, i.e., each codeword can be partitioned into message bits and syndrome bits

- As shown on the right, light blue circles correspond to the input message bits and the dark blue circle correspond to the syndrome bits

- To ensure the top check bit is satisfied, the top syndrome bit will be set to be the sum of message bits connecting to the check
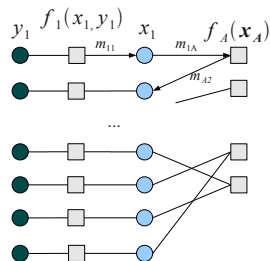
# IRA Codes

- Irregular repeated accumulate (IRA) code a type of systematic LDPC code, i.e., each codeword can be partitioned into message bits and syndrome bits

- As shown on the right, light blue circles correspond to the input message bits and the dark blue circle correspond to the syndrome bits



- To ensure the top check bit is satisfied, the top syndrome bit will be set to be the sum of message bits connecting to the check

- The computed syndrome bit will then pass to the next check and again we can ensure the next check bit is satisfied by setting that second syndrome bit as the sum of message bits conecting to the check + *last syndrome bit*. All (dark blue) syndrome bits can be assigned in similar token
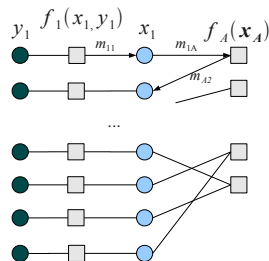
# LDPC Decoding

- $x_1, \cdots, x_N$ (light blue): transmitted bits
- $y_1, \cdots, y_N$ (dark grey): received bits

# LDPC Decoding

- $x_1, \cdots, x_N$ (light blue): transmitted bits
- $y_1, \cdots, y_N$ (dark grey): received bits
- $p(x^N, y^N) = \prod_i \underbrace{p(y_i|x_i)}_{f_i(x_i,y_i)} \underbrace{p(x^N)}_{\prod_A f_A(\mathbf{x}_A)}$
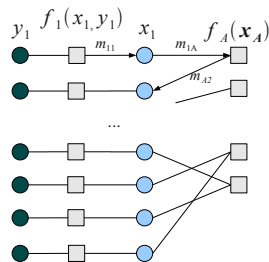
# LDPC Decoding

- $x_1, \cdots, x_N$ (light blue): transmitted bits
- $y_1, \cdots, y_N$ (dark grey): received bits
- $p(x^N, y^N) = \prod_i \underbrace{p(y_i|x_i)}_{f_i(x_i,y_i)} \underbrace{p(x^N)}_{\prod_A f_A(\mathbf{x}_A)}$
- $f_i(x_i, y_i) = p(y_i|x_i)$ and

$$f_A(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} \text{ contains even number of 1,} \\ 1, & \mathbf{x} \text{ contains odd number of 1.} \end{cases}$$

## Variable Node Update

- Since the unknown variables are binary, it is more convenient to represent the messages using likelihood or log-likelihood ratios. Define

$$l_{ai} \triangleq \frac{m_{ai}(0)}{m_{ai}(1)}, \qquad L_{ai} \triangleq \log l_{ai} \tag{2}$$

and

$$l_{ia} \triangleq \frac{m_{ia}(0)}{m_{ia}(1)}, \qquad L_{ia} \triangleq \log l_{ia} \tag{3}$$

for any variable node $i$ and factor node $a$.

- Then,

$$L_{ia} \leftarrow \sum_{b \in N(i) \setminus i} L_{ai}. \tag{4}$$

## Check Node Update

- Assuming that we have three variable nodes 1,2, and 3 connecting to the check node $a$, then the check to variable node updates become

$$m_{a1}(1) \leftarrow m_{2a}(1)m_{3a}(0) + m_{2a}(0)m_{3a}(1) \tag{5}$$
$$m_{a1}(0) \leftarrow m_{2a}(0)m_{3a}(0) + m_{2a}(1)m_{3a}(1) \tag{6}$$

- Substitute in the likelihood ratios and log-likelihood ratios, we have

$$l_{a1} \triangleq \frac{m_{a1}(0)}{m_{a1}(1)} \leftarrow \frac{1 + l_{2a}l_{3a}}{l_{2a} + l_{3a}} \tag{7}$$

and

$$e^{L_{a1}} = l_{a1} \leftarrow \frac{1 + e^{L_{2a}}e^{L_{3a}}}{e^{L_{2a}} + e^{L_{3a}}}. \tag{8}$$

- Note that

$$\tanh\left(\frac{L_{a1}}{2}\right) = \frac{e^{\frac{L_{a1}}{2}} - e^{-\frac{L_{a1}}{2}}}{e^{\frac{L_{a1}}{2}} + e^{-\frac{L_{a1}}{2}}} = \frac{e^{L_{a1}} - 1}{e^{L_{a1}} + 1} \qquad (9)$$

$$\leftarrow \frac{1 + e^{L_{2a}}e^{L_{3a}} - e^{L_{2a}} - e^{L_{3a}}}{1 + e^{L_{2a}}e^{L_{3a}} + e^{L_{2a}} + e^{L_{3a}}} \qquad (10)$$

$$= \frac{(e^{L_{2a}} - 1)(e^{L_{3a}} - 1)}{(e^{L_{2a}} + 1)(e^{L_{3a}} + 1)} \qquad (11)$$

$$= \tanh\left(\frac{L_{2a}}{2}\right)\tanh\left(\frac{L_{3a}}{2}\right). \qquad (12)$$

- When we have more than 3 variable nodes connecting to the check node $a$, it is easy to show using induction that

$$\tanh\left(\frac{L_{ai}}{2}\right) \leftarrow \prod_{j \in N(a)\setminus i} \tanh\left(\frac{L_{ja}}{2}\right). \qquad (13)$$

## More inequalities

### Lemma (Anup Rao, CSE 533, Lecture 2, Lemma 3)

If $k \leq n/2$, then $\sum_{i=0}^{k} \binom{n}{i} \leq 2^{nH(k/n)}$

### Proof.

Consider length-$n$ binary sequence $X_1, X_2, \cdots, X_n$ uniformly sampled from a set of binary sequences with at most $k$ 1's. Since there are $\sum_{i=0}^{k} \binom{n}{i}$ so many sequences, $H(X_1, X_2, \cdots, X_n) = \log \sum_{i=0}^{k} \binom{n}{i}$. On the other hand, $H(X_1, X_2, \cdots, X_n) \leq \sum_{i=1}^{n} H(X_i) = nH(k/n)$. Raise both sides with the power of two and we get the proof                     $\square$

## Example

Say we have $2^n$ people watching a subset of $2n$ movies. Each of them have at least watch 90% of all movies. At least two people actually watch the same set

### Proof.

Let's count how many different subsets a person can watch, which is

$$\sum_{i=0.9(2n)}^{2n} \binom{2n}{i} = \sum_{i=0}^{0.1(2n)} \binom{2n}{i} \leq 2^{2nH(0.1)} < 2^n$$

since $H(0.1) = 0.469 < 0.5$.

As we have $2^n$ people, by pigeon hole principle, there must be at least a pair who watched the same set