

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: $\hat{x} = \sum_{\theta} p(\theta|o) \sum_x xp(x|\theta)$
- For zero-mean \mathbf{X} , $\Sigma_{\mathbf{X}} = E[\mathbf{X}\mathbf{X}^T]$ and say we have $P^T \Sigma_{\mathbf{X}} P = D$. 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) $\Sigma_{\mathbf{X}\mathbf{X}}^{-1} = \Lambda_{\mathbf{X}\mathbf{X}} - \Lambda_{\mathbf{X}\mathbf{Y}} \Lambda_{\mathbf{Y}\mathbf{Y}}^{-1} \Lambda_{\mathbf{Y}\mathbf{X}}$
- (b) $\det(a\Sigma) = \det(a\Sigma_{\mathbf{Y}\mathbf{Y}}) \det(a\Lambda_{\mathbf{X}\mathbf{X}}^{-1})$ for any constant a

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})])$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
 $= E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T])$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
 $= E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^n \lambda_i$
 - Similarly, if we “reconstruct” \mathbf{X} as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = \text{tr}(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T])$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
 $= E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^n \lambda_i$
 - Similarly, if we “reconstruct” \mathbf{X} as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = \text{tr}(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) = \text{tr}(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]P^T)$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
 $= E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^n \lambda_i$
 - Similarly, if we “reconstruct” \mathbf{X} as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = \text{tr}(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) = \text{tr}(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]P^T) = \text{tr}(P^TPE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T])$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
 $= E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^n \lambda_i$
 - Similarly, if we “reconstruct” \mathbf{X} as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = \text{tr}(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) = \text{tr}(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]P^T) = \text{tr}(P^TPE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^n \lambda_i$

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

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 $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$
 - Generate an approximate $\hat{\mathbf{Y}}$ of \mathbf{Y} by setting all components except first k as 0
 - The mean square error (mse) of¹ $\hat{\mathbf{Y}} = E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]$
 $= \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}})]) = E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})^T(\mathbf{Y} - \hat{\mathbf{Y}}))]$
 $= E[\text{tr}((\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T)] = \text{tr}(E[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^n \lambda_i$
 - Similarly, if we “reconstruct” \mathbf{X} as $\hat{\mathbf{X}} = P\hat{\mathbf{Y}}$. The mse of $\hat{\mathbf{X}} = E[(\mathbf{X} - \hat{\mathbf{X}})^T(\mathbf{X} - \hat{\mathbf{X}})] = \text{tr}(E[(\mathbf{X} - \hat{\mathbf{X}})(\mathbf{X} - \hat{\mathbf{X}})^T]) = \text{tr}(PE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]P^T) = \text{tr}(P^TPE[(\mathbf{Y} - \hat{\mathbf{Y}})(\mathbf{Y} - \hat{\mathbf{Y}})^T]) = \sum_{i=k+1}^n \lambda_i$
 - Note that the eigenvectors of Σ (columns of P) are known as the principal components

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

Practical PCA

In practice, we typically are given a dataset with samples of \mathbf{X} instead of the distribution or covariance matrix of \mathbf{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

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

³Note 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 \mathbf{X} instead of the distribution or covariance matrix of \mathbf{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 is² $\mathcal{X} \leftarrow \mathcal{X} - \text{ones}(m, 1)\text{mean}(\mathcal{X})$

²I used the matlab notations for $\text{ones}(\cdot)$ and $\text{mean}(\cdot)$ here

³Note 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 \mathbf{X} instead of the distribution or covariance matrix of \mathbf{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 is² $\mathcal{X} \leftarrow \mathcal{X} - \text{ones}(m, 1)\text{mean}(\mathcal{X})$
- 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 approximate³

²I used the matlab notations for $\text{ones}(\cdot)$ and $\text{mean}(\cdot)$ here

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

Practical PCA

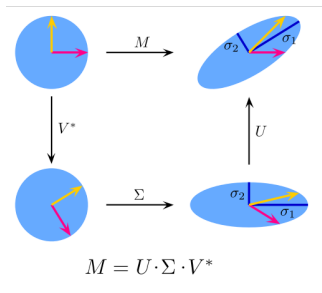
In practice, we typically are given a dataset with samples of \mathbf{X} instead of the distribution or covariance matrix of \mathbf{X} . Denote the data as \mathcal{X} with each row is a data point and a total of m data points. Thus \mathcal{X} is an m by n matrix

- Data are rarely zero-mean to begin with, but we can easily preprocess it by subtracting the mean. That is² $\mathcal{X} \leftarrow \mathcal{X} - \text{ones}(m, 1)\text{mean}(\mathcal{X})$
- 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 approximate³
 - A more common approach is to decompose \mathcal{X} with singular value decomposition (SVD) instead

²I used the matlab notations for $\text{ones}(\cdot)$ and $\text{mean}(\cdot)$ here

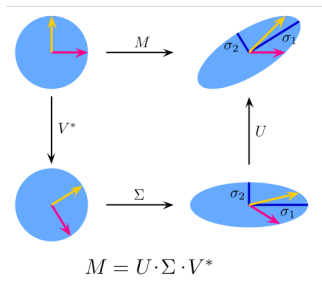
³Note that $\hat{\Sigma}$ won't be full rank and positive definite as one would hope 

Singular value decomposition (SVD)



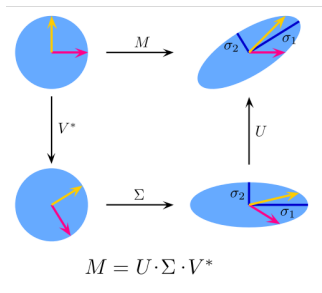
- Every matrix M can be decomposed as $M = UDV^\dagger$, where D is diagonal and U, V are unitary. The diagonal terms in Σ are known to be the **singular values**

Singular value decomposition (SVD)



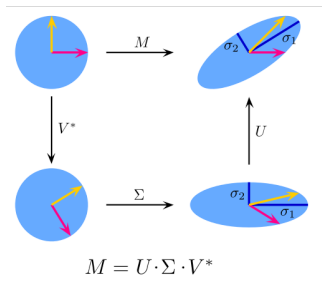
- Every matrix M can be decomposed as $M = UDV^\dagger$, where D is diagonal and U, V are unitary. The diagonal terms in Σ are known to be the **singular values**
- For real matrix M , we can write $M = UDV^T$ instead. U, V are now “real unitary” or orthogonal

Singular value decomposition (SVD)



- Every matrix M can be decomposed as $M = UDV^\dagger$, where D is diagonal and U, V are unitary. The diagonal terms in Σ are known to be the **singular values**
- For real matrix M , we can write $M = UDV^T$ instead. U, V are now “real unitary” or orthogonal
 - Note that $M^T M = VD^T U^T U D V^T = VD^2 V^T$. Therefore, V are really eigenvectors of $M^T M$ with eigenvalues equal to the square of the singular values

Singular value decomposition (SVD)



- Every matrix M can be decomposed as $M = UDV^\dagger$, where D is diagonal and U, V are unitary. The diagonal terms in Σ are known to be the **singular values**
- For real matrix M , we can write $M = UDV^T$ instead. U, V are now “real unitary” or orthogonal
 - Note that $M^T M = VD^T U^T U D V^T = VD^2 V^T$. 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

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: $\hat{x} = \sum_{\theta} p(\theta|o) \sum_x xp(x|\theta)$
- For zero-mean \mathbf{X} , $\Sigma_{\mathbf{X}} = E[\mathbf{X}\mathbf{X}^T]$ and say we have $P^T \Sigma_{\mathbf{X}} P = D$. 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) $\Sigma_{\mathbf{X}\mathbf{X}}^{-1} = \Lambda_{\mathbf{X}\mathbf{X}} - \Lambda_{\mathbf{X}\mathbf{Y}} \Lambda_{\mathbf{Y}\mathbf{Y}}^{-1} \Lambda_{\mathbf{Y}\mathbf{X}}$
- (b) $\det(a\Sigma) = \det(a\Sigma_{\mathbf{Y}\mathbf{Y}}) \det(a\Lambda_{\mathbf{X}\mathbf{X}}^{-1})$ for any constant a

Conditioning of normal distribution

- Consider the same $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{Z}}, \boldsymbol{\Sigma}_{\mathbf{Z}})$ 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 $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{Z}}, \boldsymbol{\Sigma}_{\mathbf{Z}})$ 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 $\mathbf{Z} \sim \mathcal{N}(\boldsymbol{\mu}_Z, \Sigma_Z)$ 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 $p(\mathbf{y}) = \mathcal{N}(\mathbf{y}; \boldsymbol{\mu}_Y, \Sigma_{YY})$. Therefore,

$$\begin{aligned}
 p(\mathbf{x}|\mathbf{y}) &\propto \exp \left(-\frac{1}{2} \left[\begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \end{pmatrix}^T \Sigma^{-1} \begin{pmatrix} \tilde{\mathbf{x}} \\ \tilde{\mathbf{y}} \end{pmatrix} - \tilde{\mathbf{y}}^T \Sigma_{YY}^{-1} \tilde{\mathbf{y}} \right] \right) \\
 &\propto \exp \left(-\frac{1}{2} [\tilde{\mathbf{x}}^T \Lambda_{XX} \tilde{\mathbf{x}} + \tilde{\mathbf{x}}^T \Lambda_{XY} \tilde{\mathbf{y}} + \tilde{\mathbf{y}}^T \Lambda_{YX} \tilde{\mathbf{x}}] \right),
 \end{aligned}$$

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

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &\propto \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})^T \Lambda_{\mathbf{XX}}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})\right) \\ &= \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))^T \Lambda_{\mathbf{XX}}\right. \\ &\quad \left. (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))\right) \end{aligned}$$

Conditioning of normal distribution

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

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &\propto \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})^T \Lambda_{\mathbf{XX}}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})\right) \\ &= \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))^T \Lambda_{\mathbf{XX}}\right. \\ &\quad \left. (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))\right) \end{aligned}$$

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

Conditioning of normal distribution

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

$$\begin{aligned} p(\mathbf{x}|\mathbf{y}) &\propto \exp\left(-\frac{1}{2}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})^T \Lambda_{\mathbf{XX}}(\tilde{\mathbf{x}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}\tilde{\mathbf{y}})\right) \\ &= \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))^T \Lambda_{\mathbf{XX}}\right. \\ &\quad \left. (\mathbf{x} - \boldsymbol{\mu}_{\mathbf{X}} + \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}))\right) \end{aligned}$$

- Therefore $\mathbf{X}|\mathbf{y}$ is Gaussian distributed with mean $\boldsymbol{\mu}_{\mathbf{X}} - \Lambda_{\mathbf{XX}}^{-1}\Lambda_{\mathbf{XY}}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}})$ 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

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

Interpretation of conditioning

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

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

Interpretation of conditioning

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{X}} + \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}^{-1}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}), \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{X}} - \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}^{-1}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{X}})$$

- 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 $\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}$, the variance of \mathbf{Y} for the 1-D case.
 - The observation is less reliable with the increase of $\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}$. The adjustment is finally scaled by $\boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}$, which translates the variation of \mathbf{Y} to the variation of \mathbf{X}

Interpretation of conditioning

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{X}} + \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}^{-1}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}), \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{X}} - \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}^{-1}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{X}})$$

- 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 $\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}$, the variance of \mathbf{Y} for the 1-D case.
 - The observation is less reliable with the increase of $\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}$. The adjustment is finally scaled by $\boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}$, 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

Interpretation of conditioning

$$\mathbf{X}|\mathbf{y} \sim \mathcal{N}(\boldsymbol{\mu}_{\mathbf{X}} + \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}^{-1}(\mathbf{y} - \boldsymbol{\mu}_{\mathbf{Y}}), \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{X}} - \boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}^{-1}\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{X}})$$

- 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 $\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}$, the variance of \mathbf{Y} for the 1-D case.
 - The observation is less reliable with the increase of $\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}$. The adjustment is finally scaled by $\boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}$, 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 $\boldsymbol{\Sigma}_{\mathbf{Y}\mathbf{Y}}$ is smaller and $\boldsymbol{\Sigma}_{\mathbf{X}\mathbf{Y}}$ is larger (\mathbf{X} and \mathbf{Y} are more correlated)

$X \perp\!\!\!\perp Y|Z$ 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 coefficient between X and Z . Similarly, ρ_{YZ} and ρ_{XY} are the correlation coefficients between Y and Z , and X and Y , respectively.

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

Proof.

- Without loss of generality, we can assume the variables with mean 0 and variance 1. Thus, $\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \Sigma)$, where $\Sigma = \begin{pmatrix} 1 & \rho_{XY} & \rho_{XZ} \\ \rho_{XY} & 1 & \rho_{YZ} \\ \rho_{XZ} & \rho_{YZ} & 1 \end{pmatrix}$

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

Proof.

- Without loss of generality, we can assume the variables with mean 0 and variance 1. Thus, $\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \Sigma)$, where $\Sigma = \begin{pmatrix} 1 & \rho_{XY} & \rho_{XZ} \\ \rho_{XY} & 1 & \rho_{YZ} \\ \rho_{XZ} & \rho_{YZ} & 1 \end{pmatrix}$
- Then from the conditioning result, we have

$$\begin{aligned} \Sigma \begin{pmatrix} X \\ Y \end{pmatrix} | Z &= \begin{pmatrix} 1 & \rho_{XY} \\ \rho_{XY} & 1 \end{pmatrix} - (\rho_{XZ} \quad \rho_{YZ}) \sigma_{YY}^{-1} \begin{pmatrix} \rho_{XZ} \\ \rho_{YZ} \end{pmatrix} \\ &= \begin{pmatrix} 1 - \rho_{XZ}^2 & \rho_{XY} - \rho_{XZ}\rho_{YZ} \\ \rho_{XY} - \rho_{XZ}\rho_{YZ} & 1 - \rho_{YZ}^2 \end{pmatrix} \end{aligned}$$

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

Proof.

- Without loss of generality, we can assume the variables with mean 0 and variance 1. Thus, $\begin{pmatrix} X \\ Y \\ Z \end{pmatrix} \sim \mathcal{N}(\mathbf{0}, \Sigma)$, where $\Sigma = \begin{pmatrix} 1 & \rho_{XY} & \rho_{XZ} \\ \rho_{XY} & 1 & \rho_{YZ} \\ \rho_{XZ} & \rho_{YZ} & 1 \end{pmatrix}$

- Then from the conditioning result, we have

$$\begin{aligned} \Sigma \begin{pmatrix} X \\ Y \end{pmatrix} | Z &= \begin{pmatrix} 1 & \rho_{XY} \\ \rho_{XY} & 1 \end{pmatrix} - (\rho_{XZ} \quad \rho_{YZ}) \sigma_{YY}^{-1} \begin{pmatrix} \rho_{XZ} \\ \rho_{YZ} \end{pmatrix} \\ &= \begin{pmatrix} 1 - \rho_{XZ}^2 & \rho_{XY} - \rho_{XZ}\rho_{YZ} \\ \rho_{XY} - \rho_{XZ}\rho_{YZ} & 1 - \rho_{YZ}^2 \end{pmatrix} \end{aligned}$$

- Therefore, X and Y are uncorrelated given Z when $\sigma_{XY|Z} = \rho_{XY} - \rho_{XZ}\rho_{YZ} = 0$ or $\rho_{XY} = \rho_{XZ}\rho_{YZ}$. Since for Gaussian variables, uncorrelatedness implies independence. This concludes the proof. □

Product of normal distributions

- Assume that we try to recover some vector parameter \mathbf{x} , which is subject to multivariate Gaussian noise

Product of normal distributions

- Assume that we try 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

Product of normal distributions

- Assume that we try 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

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

Product of normal distributions

- Assume that we try 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

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

- 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.

$$\begin{aligned} & \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1}) \mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) \\ & \propto \exp \left(-\frac{1}{2} [(\mathbf{x} - \mathbf{y}_1)^T \Lambda_{\mathbf{Y}_1} (\mathbf{x} - \mathbf{y}_1) + (\mathbf{x} - \mathbf{y}_2)^T \Lambda_{\mathbf{Y}_2} (\mathbf{x} - \mathbf{y}_2)] \right) \end{aligned}$$

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.

$$\begin{aligned} & \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1}) \mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) \\ & \propto \exp \left(-\frac{1}{2} [(\mathbf{x} - \mathbf{y}_1)^T \Lambda_{\mathbf{Y}_1} (\mathbf{x} - \mathbf{y}_1) + (\mathbf{x} - \mathbf{y}_2)^T \Lambda_{\mathbf{Y}_2} (\mathbf{x} - \mathbf{y}_2)] \right) \\ & \propto \exp \left(-\frac{1}{2} [\mathbf{x}^T (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2}) \mathbf{x} - (\mathbf{y}_2^T \Lambda_{\mathbf{Y}_2} + \mathbf{y}_1^T \Lambda_{\mathbf{Y}_1}) \mathbf{x} - \mathbf{x}^T (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1)] \right) \end{aligned}$$

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.

$$\begin{aligned}
 & \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1}) \mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) \\
 & \propto \exp \left(-\frac{1}{2} [(\mathbf{x} - \mathbf{y}_1)^T \Lambda_{\mathbf{Y}_1} (\mathbf{x} - \mathbf{y}_1) + (\mathbf{x} - \mathbf{y}_2)^T \Lambda_{\mathbf{Y}_2} (\mathbf{x} - \mathbf{y}_2)] \right) \\
 & \propto \exp \left(-\frac{1}{2} [\mathbf{x}^T (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2}) \mathbf{x} - (\mathbf{y}_2^T \Lambda_{\mathbf{Y}_2} + \mathbf{y}_1^T \Lambda_{\mathbf{Y}_1}) \mathbf{x} - \mathbf{x}^T (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1)] \right) \\
 & \propto e^{-\frac{1}{2} [(\mathbf{x} - (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1} (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1))^T (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2}) (\mathbf{x} - (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1} (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1))] }
 \end{aligned}$$

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.

$$\begin{aligned}
 & \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1}) \mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) \\
 & \propto \exp \left(-\frac{1}{2} [(\mathbf{x} - \mathbf{y}_1)^T \Lambda_{\mathbf{Y}_1} (\mathbf{x} - \mathbf{y}_1) + (\mathbf{x} - \mathbf{y}_2)^T \Lambda_{\mathbf{Y}_2} (\mathbf{x} - \mathbf{y}_2)] \right) \\
 & \propto \exp \left(-\frac{1}{2} [\mathbf{x}^T (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2}) \mathbf{x} - (\mathbf{y}_2^T \Lambda_{\mathbf{Y}_2} + \mathbf{y}_1^T \Lambda_{\mathbf{Y}_1}) \mathbf{x} - \mathbf{x}^T (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1)] \right) \\
 & \propto e^{-\frac{1}{2} [(\mathbf{x} - (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1} (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1))^T (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2}) (\mathbf{x} - (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1} (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1))]} \\
 & \propto \mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1} (\Lambda_{\mathbf{Y}_2} \mathbf{y}_2 + \Lambda_{\mathbf{Y}_1} \mathbf{y}_1), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})
 \end{aligned}$$

Therefore,

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

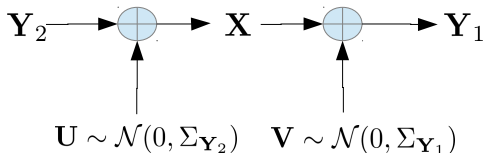
for some scaling factor $K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})$ independent of \mathbf{x} .

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

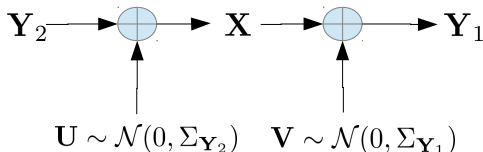
Product of normal distributions

- 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



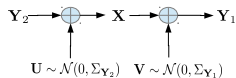
- Since $\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2}) = \mathcal{N}(\mathbf{x}; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2})$ and $\mathbf{Y}_1 \perp\!\!\!\perp \mathbf{Y}_2 | \mathbf{X}$, we have

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

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,

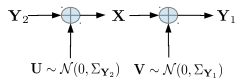


$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2) d\mathbf{x} = p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{Y_2} + \Sigma_{Y_1})$$

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,



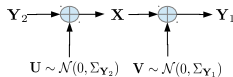
$$\int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} = p(\mathbf{y}_1|\mathbf{y}_2) = \mathcal{N}(\mathbf{y}_1; \mathbf{y}_2, \Sigma_{\mathbf{Y}_2} + \Sigma_{\mathbf{Y}_1})$$

- On the other hand,

$$\begin{aligned} \int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} &= \int \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})d\mathbf{x} \\ &= \int K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})d\mathbf{x} \\ &= K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}). \end{aligned}$$

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,

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

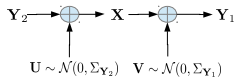
- On the other hand,

$$\begin{aligned} \int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} &= \int \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})d\mathbf{x} \\ &= \int K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})d\mathbf{x} \\ &= K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}). \end{aligned}$$

- 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,

$$\begin{aligned} \int p(\mathbf{y}_1, \mathbf{x}|\mathbf{y}_2)d\mathbf{x} &= \int \mathcal{N}(\mathbf{y}_1; \mathbf{x}, \Sigma_{\mathbf{Y}_1})\mathcal{N}(\mathbf{y}_2; \mathbf{x}, \Sigma_{\mathbf{Y}_2})d\mathbf{x} \\ &= \int K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2})\mathcal{N}(\mathbf{x}; (\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1), (\Lambda_{\mathbf{Y}_2} + \Lambda_{\mathbf{Y}_1})^{-1})d\mathbf{x} \\ &= K(\mathbf{y}_1, \mathbf{y}_2, \Sigma_{\mathbf{Y}_1}, \Sigma_{\mathbf{Y}_2}). \end{aligned}$$

- 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

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

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 scalar

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 scalar

- 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 essentially 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 scalar

- The mean considering both observations, $(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}y_2 + \Lambda_{\mathbf{Y}_1}y_1)$, is essentially a weighted average of observations y_2 and 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 y_1 and 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 scalar

- The mean considering both observations, $(\Lambda_{\mathbf{Y}_1} + \Lambda_{\mathbf{Y}_2})^{-1}(\Lambda_{\mathbf{Y}_2}\mathbf{y}_2 + \Lambda_{\mathbf{Y}_1}\mathbf{y}_1)$, is essentially 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

Division of normal distributions

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

$$\begin{aligned} & \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \mathcal{N}(\mathbf{x}; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}) \\ &= \mathcal{N}(\boldsymbol{\mu}_1; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), \boldsymbol{\Lambda}_2^{-1} + (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \end{aligned}$$

Division of normal distributions

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

$$\begin{aligned} & \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \mathcal{N}(\mathbf{x}; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}) \\ &= \mathcal{N}(\boldsymbol{\mu}_1; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), \boldsymbol{\Lambda}_2^{-1} + (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \end{aligned}$$

- Therefore,

$$\begin{aligned} \frac{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} &= \frac{\mathcal{N}(\mathbf{x}; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_1, (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2); \boldsymbol{\Lambda}_2^{-1} + (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})} \\ &= \frac{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_1; \boldsymbol{\mu}, \boldsymbol{\Lambda}_2^{-1} + (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})}, \end{aligned}$$

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

Division of normal distributions

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

$$\begin{aligned} & \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2) \mathcal{N}(\mathbf{x}; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}) \\ &= \mathcal{N}(\boldsymbol{\mu}_1; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), \boldsymbol{\Lambda}_2^{-1} + (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}) \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1) \end{aligned}$$

- Therefore,

$$\begin{aligned} \frac{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)}{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)} &= \frac{\mathcal{N}(\mathbf{x}; (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2), (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_1, (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1}(\boldsymbol{\Lambda}_1 \boldsymbol{\mu}_1 - \boldsymbol{\Lambda}_2 \boldsymbol{\mu}_2); \boldsymbol{\Lambda}_2^{-1} + (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})} \\ &= \frac{\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})}{\mathcal{N}(\boldsymbol{\mu}_1; \boldsymbol{\mu}, \boldsymbol{\Lambda}_2^{-1} + (\boldsymbol{\Lambda}_1 - \boldsymbol{\Lambda}_2)^{-1})}, \end{aligned}$$

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

- Note that the final pdf will be Gaussian-like if $\boldsymbol{\Lambda}_1 \succeq \boldsymbol{\Lambda}_2$. 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, S behaves like $\mathcal{N}(0, 1)$

Mixture of Gaussians

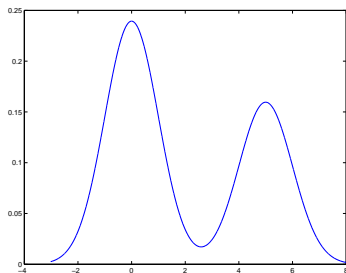
Consider an electrical system that outputs signal of different statistics when it is on and off

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

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

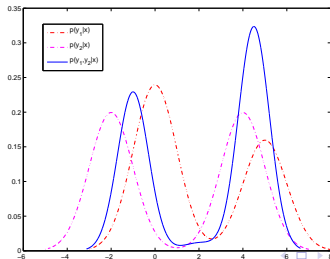
- We can see that the weight for the component at mean 1.5 is very small. And the component at mean 2 has a rather small weight also.

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

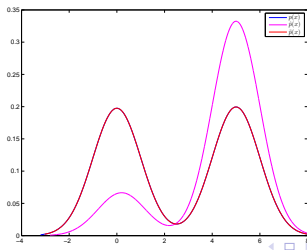
Consider

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

- 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

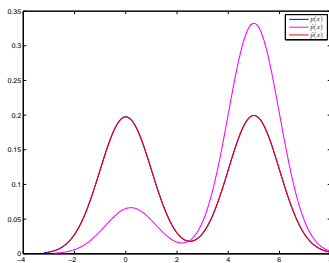
- The problem is that while the first five components are all relatively small compared to the last one, they are all quite similar and their combined contribution is comparable to the latter

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 we gauge the similarity between two components.

Which Components to Merge?

It is reasonable to pick similar components to merge. The question is how do we 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 we 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 we 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,

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

Which Components to Merge?

It is reasonable to pick similar components to merge. The question is how do we 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,

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

- 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

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

Similarity measure

- Let's define

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

- 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)

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

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 $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), \dots, \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ with weights w_1, w_2, \dots, w_n . What should the combined component be like?

How to Merge Components?

Say we have n components $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), \dots, \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ with weights w_1, w_2, \dots, 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 $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), \dots, \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ with weights w_1, w_2, \dots, 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 $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), \dots, \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ with weights w_1, w_2, \dots, 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 $\sum_{i=1}^n \hat{w}_i \boldsymbol{\Sigma}_i$.

How to Merge Components?

Say we have n components $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), \dots, \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ with weights w_1, w_2, \dots, 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 $\sum_{i=1}^n \hat{w}_i \boldsymbol{\Sigma}_i$.
 - However, it is an underestimate

How to Merge Components?

Say we have n components $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), \dots, \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ with weights w_1, w_2, \dots, 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 $\sum_{i=1}^n \hat{w}_i \boldsymbol{\Sigma}_i$.
 - 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 $\mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1), \mathcal{N}(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2), \dots, \mathcal{N}(\boldsymbol{\mu}_n, \boldsymbol{\Sigma}_n)$ with weights w_1, w_2, \dots, 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 $\sum_{i=1}^n \hat{w}_i \boldsymbol{\Sigma}_i$.
 - 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 \mathbf{X} as the variable sampled from the mixture. That is, $\mathbf{X} \sim \mathcal{N}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ with probability \hat{w}_i . Then, we have (please verify)

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

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

