

Neural Networks

Samuel Cheng

School of ECE
University of Oklahoma

Spring, 2019

Table of Contents

- 1 Review
- 2 SVM
- 3 Introduction to neural networks
- 4 Back-propagation
- 5 Activation functions
- 6 Initialization
- 7 Regularization

- Need your presentation preference by this Thursday. Please give me two package names with order of preference. The final decision will be computed by minimizing the following cost function :)
 - $\sum_{\text{student}} \text{student cost} + \sum_{\text{package}} \text{package cost}$
 - $\text{student cost} = \begin{cases} 0, & \text{first priority} \\ 3, & \text{second priority} \end{cases}$
 - $\text{package cost} = \alpha \cdot 2^{(\text{num presentations covered})}$

- Need your presentation preference by this Thursday. Please give me two package names with order of preference. The final decision will be computed by minimizing the following cost function :)
 - $\sum_{\text{student}} \text{student cost} + \sum_{\text{package}} \text{package cost}$
 - $\text{student cost} = \begin{cases} 0, & \text{first priority} \\ 3, & \text{second priority} \end{cases}$
 - $\text{package cost} = \alpha \cdot 2^{(\text{num presentations covered})}$
- Most popular package (in terms of first priority pick) will be presented first. If there is a tie, I will break it with popularity based all choices regardless of priority. If there is a tie, I will break it by random

- Need your presentation preference by this Thursday. Please give me two package names with order of preference. The final decision will be computed by minimizing the following cost function :)
 - $\sum_{\text{student}} \text{student cost} + \sum_{\text{package}} \text{package cost}$
 - $\text{student cost} = \begin{cases} 0, & \text{first priority} \\ 3, & \text{second priority} \end{cases}$
 - $\text{package cost} = \alpha \cdot 2^{(\text{num presentations covered})}$
- Most popular package (in terms of first priority pick) will be presented first. If there is a tie, I will break it with popularity based all choices regardless of priority. If there is a tie, I will break it by random
- Students presenting the same packages will be ordered randomly

- Package choice due this Thursday

Review

In the last couple classes, we discussed

- Basic concepts of regression and classification
- Examples of regularization such as ridge (l_2) regression and lasso (l_1)
- Linear classifiers including logistic regression and softmax classifier

Review

In the last couple classes, we discussed

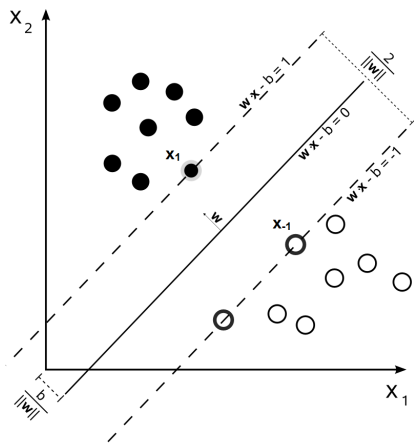
- Basic concepts of regression and classification
- Examples of regularization such as ridge (l_2) regression and lasso (l_1)
- Linear classifiers including logistic regression and softmax classifier
 - We introduced loss functions and the concept of training a classifier through minimizing the loss function

Review

In the last couple classes, we discussed

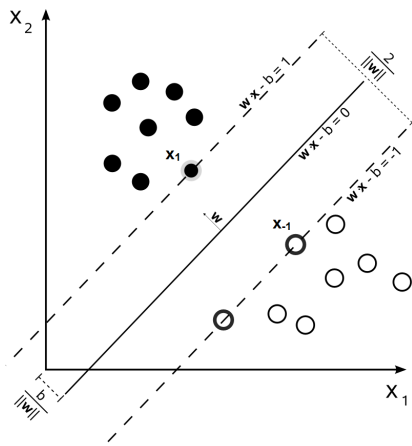
- Basic concepts of regression and classification
- Examples of regularization such as ridge (l_2) regression and lasso (l_1)
- Linear classifiers including logistic regression and softmax classifier
 - We introduced loss functions and the concept of training a classifier through minimizing the loss function
 - We described stochastic gradient descent and momentum trick for classification

SVM



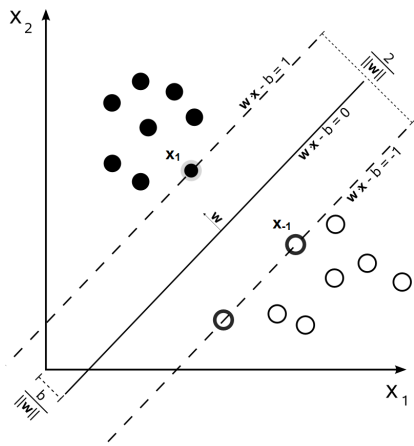
- Denote $\hat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$, $\hat{\mathbf{w}} \cdot \mathbf{x}_1$
 $(\hat{\mathbf{w}} \cdot \mathbf{x}_1)$ is the distance of the
 boundary line of \mathbf{x}_1 (\mathbf{x}_1)
 from the origin

SVM



- Denote $\hat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$, $\hat{\mathbf{w}} \cdot \mathbf{x}_1$
 $(\hat{\mathbf{w}} \cdot \mathbf{x}_{-1})$ is the distance of the
 boundary line of \mathbf{x}_1 (\mathbf{x}_{-1})
 from the origin
- Thus, the distance between
 the two boundary lines is
 $\hat{\mathbf{w}} \cdot (\mathbf{x}_1 - \mathbf{x}_{-1}) = \frac{2}{\|\mathbf{w}\|}$

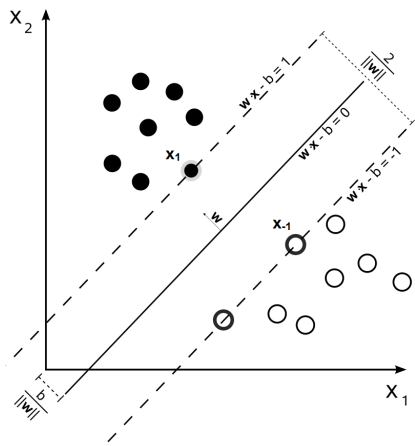
SVM



- Denote $\hat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$, $\hat{\mathbf{w}} \cdot \mathbf{x}_1$
 $(\hat{\mathbf{w}} \cdot \mathbf{x}_{-1})$ is the distance of the boundary line of \mathbf{x}_1 (\mathbf{x}_{-1}) from the origin
- Thus, the distance between the two boundary lines is $\hat{\mathbf{w}} \cdot (\mathbf{x}_1 - \mathbf{x}_{-1}) = \frac{2}{\|\mathbf{w}\|}$
- SVM: for all \mathbf{x}_i

$$\max \frac{2}{\|\mathbf{w}\|} \quad \text{s.t.} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1$$

SVM



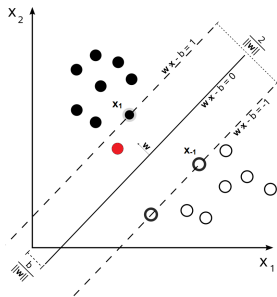
- Denote $\hat{\mathbf{w}} = \frac{\mathbf{w}}{\|\mathbf{w}\|}$, $\hat{\mathbf{w}} \cdot \mathbf{x}_1$ ($\hat{\mathbf{w}} \cdot \mathbf{x}_{-1}$) is the distance of the boundary line of \mathbf{x}_1 (\mathbf{x}_{-1}) from the origin
- Thus, the distance between the two boundary lines is $\hat{\mathbf{w}} \cdot (\mathbf{x}_1 - \mathbf{x}_{-1}) = \frac{2}{\|\mathbf{w}\|}$
- SVM: for all \mathbf{x}_i

$$\max \frac{2}{\|\mathbf{w}\|} \quad \text{s.t.} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1$$

Equivalently,

$$\min \|\mathbf{w}\| \quad \text{s.t.} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i - b) \geq 1$$

Soft-margin SVM and hinge loss



- Hard-margin SVM

$$\min \|\mathbf{w}\| \quad s.t. \quad y_i(\mathbf{w} \cdot \mathbf{x}_i - b) - 1 \geq 0$$

- Soft-margin SVM (allow constrain to be violate)

- Define “hinge” loss function
- Want to minimize hinge loss

$$\sum_i h(1 - y_i(\mathbf{w} \cdot \mathbf{x}_i - b))$$

- Soft-margin SVM

$$\min \lambda \|\mathbf{w}\|^2 + \sum_i h(1 - y_i(\mathbf{w} \cdot \mathbf{x}_i - b))$$

Multi-class SVM

- We can easily extend soft-margin SVM to multi-class case. Let $s_l(\mathbf{x}) = \mathbf{w}_l^T \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$ be the score for class l .

Multi-class SVM

- We can easily extend soft-margin SVM to multi-class case. Let $s_l(\mathbf{x}) = \mathbf{w}_l^T \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$ be the score for class l . We can define the hinge loss for sample \mathbf{x} as

$$\sum_{l \neq j} h(s_l(\mathbf{x}) - s_j(\mathbf{x}) + \Delta) = \sum_{l \neq j} \max(0, s_l(\mathbf{x}) - s_j(\mathbf{x}) + \Delta),$$

where j is the true label of \mathbf{x} and Δ contributes a margin ensuring that the true label score has to be at least Δ more than the rest to be penalty free

Multi-class SVM

- We can easily extend soft-margin SVM to multi-class case. Let $s_l(\mathbf{x}) = \mathbf{w}_l^T \begin{bmatrix} 1 \\ \mathbf{x} \end{bmatrix}$ be the score for class l . We can define the hinge loss for sample \mathbf{x} as

$$\sum_{l \neq j} h(s_l(\mathbf{x}) - s_j(\mathbf{x}) + \Delta) = \sum_{l \neq j} \max(0, s_l(\mathbf{x}) - s_j(\mathbf{x}) + \Delta),$$

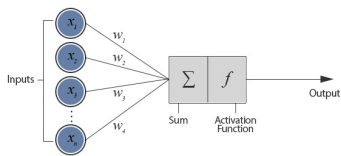
where j is the true label of \mathbf{x} and Δ contributes a margin ensuring that the true label score has to be at least Δ more than the rest to be penalty free

- Multi-class SVM:

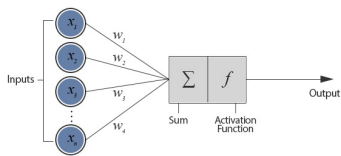
$$\min \lambda \|\mathbf{w}\|^2 + \sum_i \sum_{l \neq j(\mathbf{x}_i)} h(s_l(\mathbf{x}_i) - s_{j(\mathbf{x}_i)}(\mathbf{x}_i) + \Delta)$$

Perceptron

- Perceptron is an artificial neuron with step function as activation function

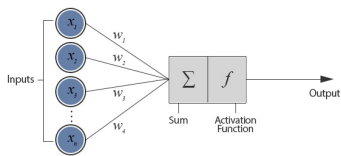


Perceptron



- Perceptron is an artificial neuron with step function as activation function
- It is impossible to extend perceptron to multilayer. Multilayer perceptron (MLP) is a misnomer. Step activation function is never used multilayer neural networks (not trainable)

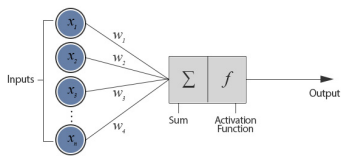
Perceptron



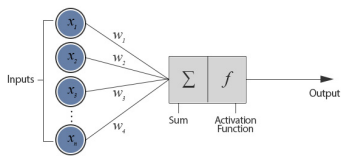
- Perceptron is an artificial neuron with step function as activation function
- It is impossible to extend perceptron to multilayer. Multilayer perceptron (MLP) is a misnomer. Step activation function is never used multilayer neural networks (not trainable)
- According to Hinton, perceptrons are still used in systems with large number (millions) of features. Other than that, it has relatively limited use since most problems are not linearly separable

Perceptron

- In most cases, perceptron would be useful if only one manages to handcode inputs into separable features

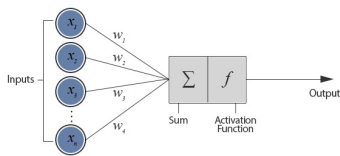


Perceptron



- In most cases, perceptron would be useful if only one manages to handcode inputs into separable features
- That was the main area of research in many machine learning applications—finding efficient ways to generate good features

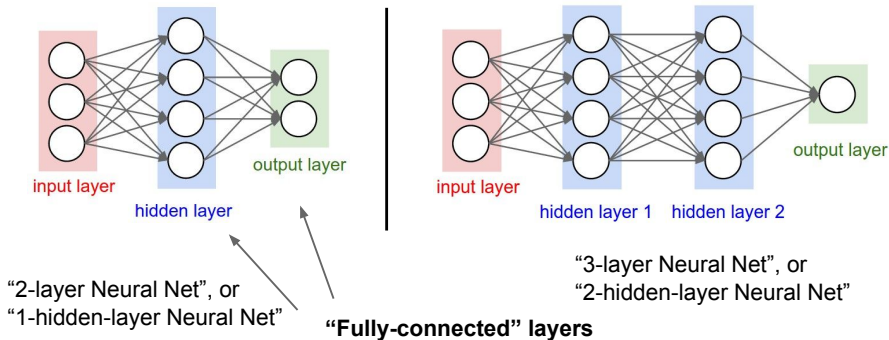
Perceptron



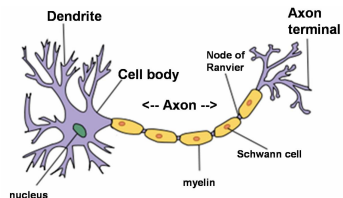
- In most cases, perceptron would be useful if only one manages to handcode inputs into separable features
- That was the main area of research in many machine learning applications—finding efficient ways to generate good features
- One attractive characteristic of deep learning (neural networks) is that we not only can train the classifier but also can learn the appropriate features automatically

Nomenclature of basic network architectures

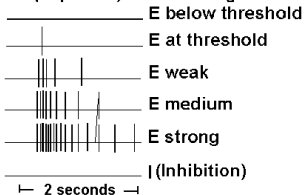
Neural Networks: Architectures



Caveat: don't go too far for the brain analogy



Axon's All-Or-Nothing Action Potentials (impulses) to Increasing Excitation



Biological neurons:

- Many different types
- Dendrite can perform complex non-linear operations
- Synapses are not a single weight but a complex non-linear dynamical system
- Rate code model may not be adequate

Also see London 2005 (Slide credit: CS231n)

Back-propagation and computational graph

- As described in last lecture, training in supervised learning system often boils down to minimizing of loss function w.r.t. some parameters

Back-propagation and computational graph

- As described in last lecture, training in supervised learning system often boils down to minimizing of loss function w.r.t. some parameters
- For neural networks, it is thus necessary to find $\frac{\partial L(\mathbf{w}; \mathbf{x})}{\partial w}$ for a weight in each layer

Back-propagation and computational graph

- As described in last lecture, training in supervised learning system often boils down to minimizing of loss function w.r.t. some parameters
- For neural networks, it is thus necessary to find $\frac{\partial L(\mathbf{w}; \mathbf{x})}{\partial w}$ for a weight in each layer
- Back-propagation (BP) is an efficient way to find such derivation. Actually it is in fact just another way of spelling out the chain rule $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$ in calculus

Back-propagation and computational graph

- As described in last lecture, training in supervised learning system often boils down to minimizing of loss function w.r.t. some parameters
- For neural networks, it is thus necessary to find $\frac{\partial L(\mathbf{w}; \mathbf{x})}{\partial w}$ for a weight in each layer
- Back-propagation (BP) is an efficient way to find such derivation. Actually it is in fact just another way of spelling out the chain rule $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$ in calculus
- It is often easier to explain BP in terms of **computational graph**

Back-propagation and computational graph

- As described in last lecture, training in supervised learning system often boils down to minimizing of loss function w.r.t. some parameters
- For neural networks, it is thus necessary to find $\frac{\partial L(\mathbf{w}; \mathbf{x})}{\partial w}$ for a weight in each layer
- Back-propagation (BP) is an efficient way to find such derivation. Actually it is in fact just another way of spelling out the chain rule $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$ in calculus
- It is often easier to explain BP in terms of **computational graph**
 - Computational graph can be interpreted as generalization of a neural networks
 - Neuron no longer will be restricted to summation and activation function but can be any computation as well (e.g., max)

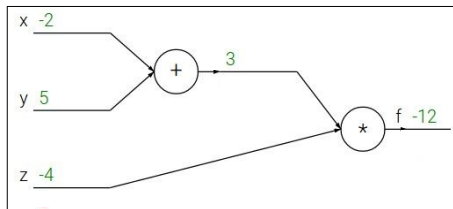
Back-propagation and computational graph

- As described in last lecture, training in supervised learning system often boils down to minimizing of loss function w.r.t. some parameters
- For neural networks, it is thus necessary to find $\frac{\partial L(\mathbf{w}; \mathbf{x})}{\partial w}$ for a weight in each layer
- Back-propagation (BP) is an efficient way to find such derivation. Actually it is in fact just another way of spelling out the chain rule $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} \frac{\partial y}{\partial x}$ in calculus
- It is often easier to explain BP in terms of **computational graph**
 - Computational graph can be interpreted as generalization of a neural networks
 - Neuron no longer will be restricted to summation and activation function but can be any computation as well (e.g., max)
- Let me try to explain through an example

A simple BP example

$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$



A simple BP example

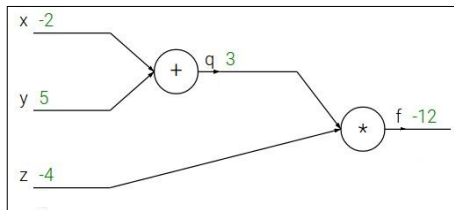
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



A simple BP example

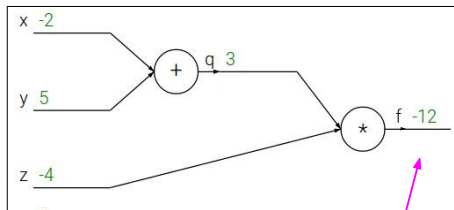
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial f}$$

A simple BP example

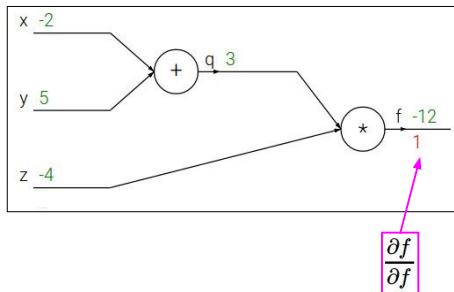
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



A simple BP example

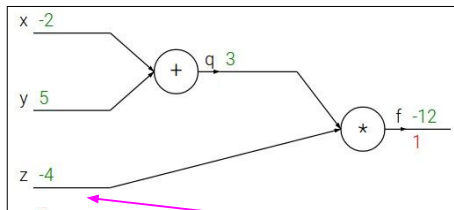
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

A simple BP example

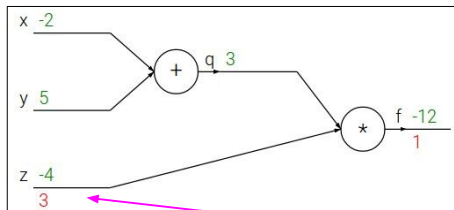
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

A simple BP example

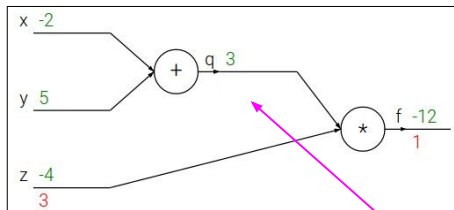
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$

A simple BP example

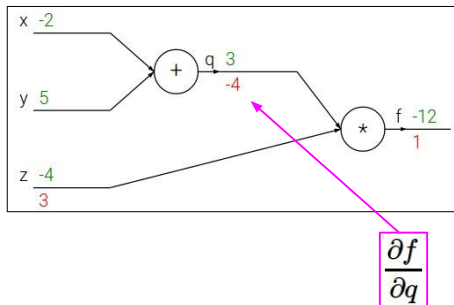
$$f(x, y, z) = (x + y)z$$

$$\text{e.g. } x = -2, y = 5, z = -4$$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

$$\text{Want: } \frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$$



A simple BP example

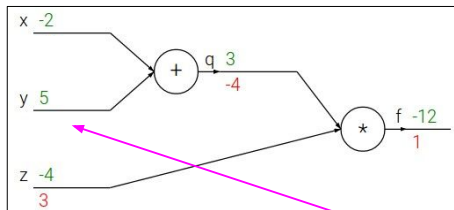
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial y}$$

A simple BP example

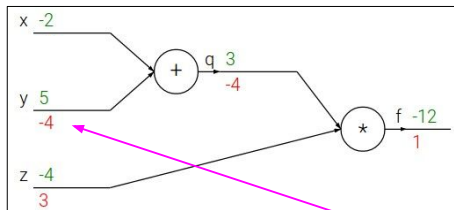
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial y}$$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

A simple BP example

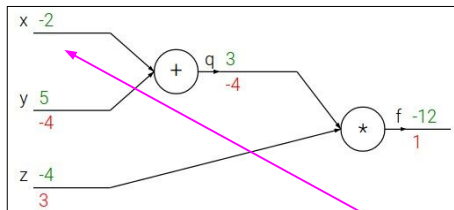
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x}$$

A simple BP example

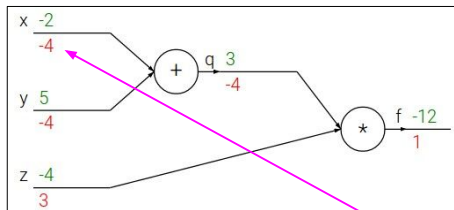
$$f(x, y, z) = (x + y)z$$

e.g. $x = -2, y = 5, z = -4$

$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \quad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



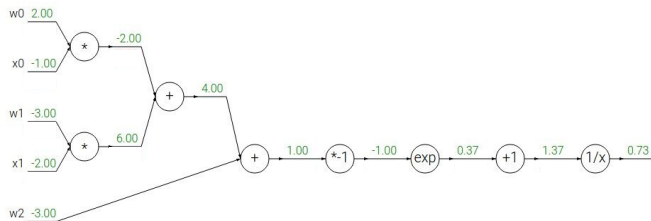
Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

$$\frac{\partial f}{\partial x}$$

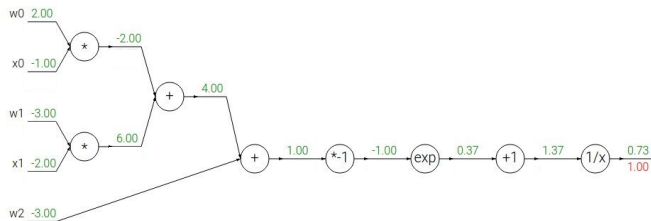
Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



Yet another BP example

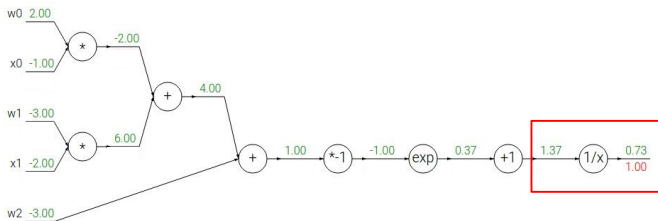
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{lcl}
 f(x) = e^x & \rightarrow & \frac{df}{dx} = e^x \\
 f_a(x) = ax & \rightarrow & \frac{df}{dx} = a
 \end{array}
 \quad \Bigg| \quad
 \begin{array}{lcl}
 f(x) = \frac{1}{x} & \rightarrow & \frac{df}{dx} = -1/x^2 \\
 f_c(x) = c + x & \rightarrow & \frac{df}{dx} = 1
 \end{array}$$

Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

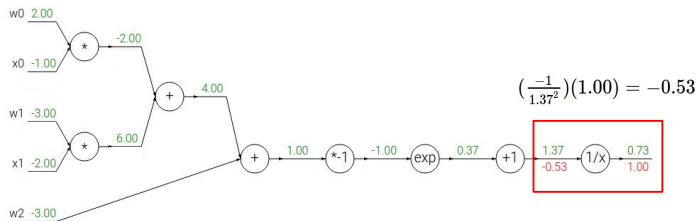
$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

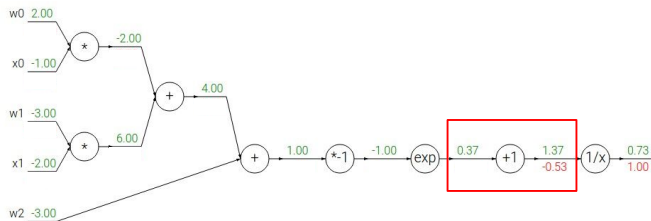
$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

Yet another BP example

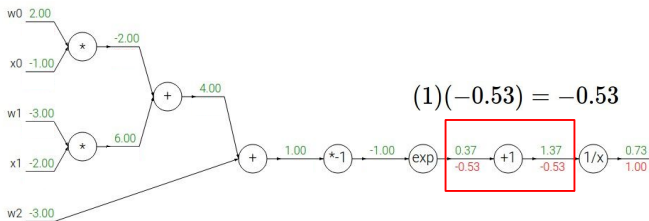
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$f(x) = e^x$	\rightarrow	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	\rightarrow	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	\rightarrow	$\frac{df}{dx} = a$		$f_c(x) = c + x$	\rightarrow	$\frac{df}{dx} = 1$

Yet another BP example

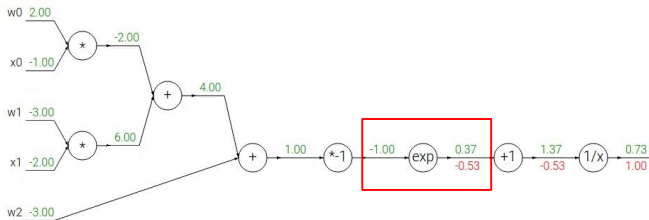
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



$f(x) = e^x$	\rightarrow	$\frac{df}{dx} = e^x$		$f(x) = \frac{1}{x}$	\rightarrow	$\frac{df}{dx} = -1/x^2$
$f_a(x) = ax$	\rightarrow	$\frac{df}{dx} = a$		$f_c(x) = c + x$	\rightarrow	$\frac{df}{dx} = 1$

Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

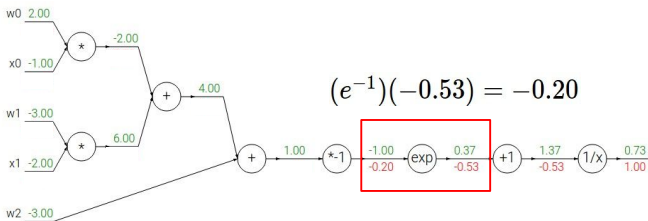
$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



$$(e^{-1})(-0.53) = -0.20$$

$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

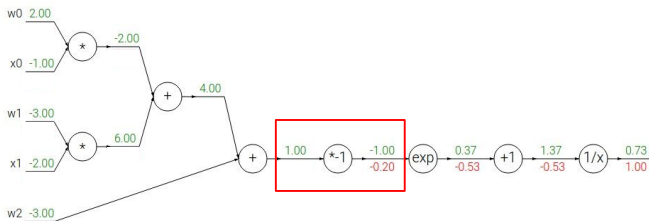
$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

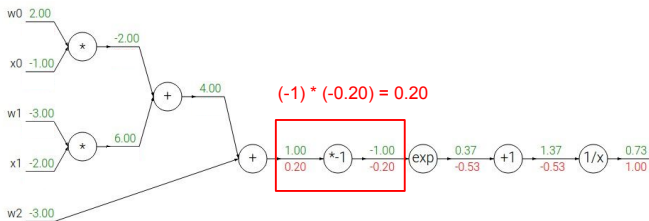
$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



$$f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x$$

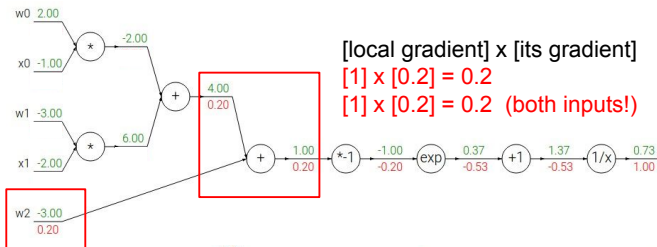
$$f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1$$

Yet another BP example

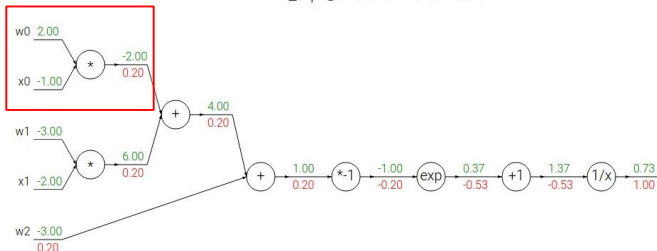
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



$$\begin{array}{l}
 f(x) = e^x \quad \rightarrow \quad \frac{df}{dx} = e^x \\
 f_a(x) = ax \quad \rightarrow \quad \frac{df}{dx} = a
 \end{array}
 \quad \Bigg| \quad
 \begin{array}{l}
 f(x) = \frac{1}{x} \quad \rightarrow \quad \frac{df}{dx} = -1/x^2 \\
 f_c(x) = c + x \quad \rightarrow \quad \frac{df}{dx} = 1
 \end{array}$$

Yet another BP example

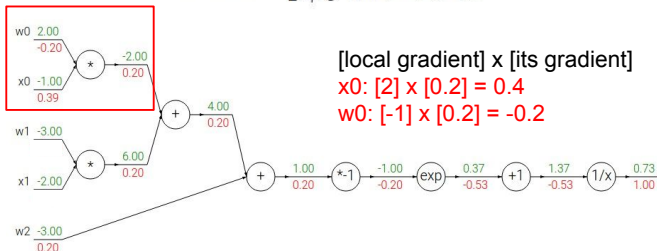
Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$



$$\begin{array}{lcl}
 f(x) = e^x & \rightarrow & \frac{df}{dx} = e^x \\
 f_a(x) = ax & \rightarrow & \frac{df}{dx} = a
 \end{array}
 \quad \Bigg| \quad
 \begin{array}{lcl}
 f(x) = \frac{1}{x} & \rightarrow & \frac{df}{dx} = -1/x^2 \\
 f_c(x) = c + x & \rightarrow & \frac{df}{dx} = 1
 \end{array}$$

Yet another BP example

Another example: $f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$



$$\begin{array}{lcl}
 f(x) = e^x & \rightarrow & \frac{df}{dx} = e^x \\
 f_a(x) = ax & \rightarrow & \frac{df}{dx} = a
 \end{array}
 \quad \Bigg| \quad
 \begin{array}{lcl}
 f(x) = \frac{1}{x} & \rightarrow & \frac{df}{dx} = -1/x^2 \\
 f_c(x) = c + x & \rightarrow & \frac{df}{dx} = 1
 \end{array}$$

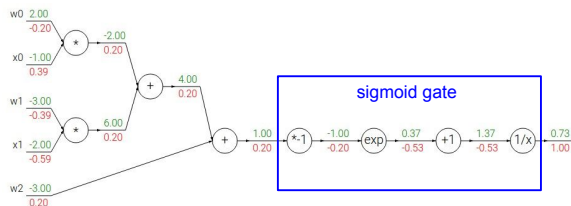
Breaking down at different granularities

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2)}}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$



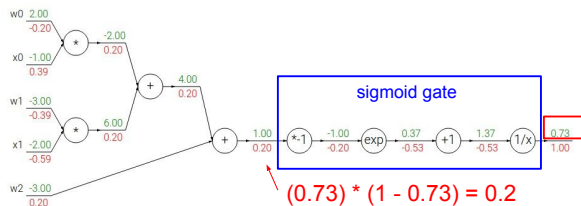
Breaking down at different granularities

$$f(w, x) = \frac{1}{1 + e^{-(w_0x_0 + w_1x_1 + w_2x_2)}}$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

sigmoid function

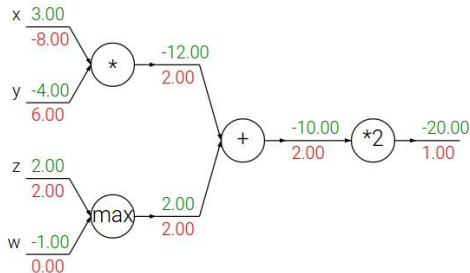
$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}} \right) \left(\frac{1}{1 + e^{-x}} \right) = (1 - \sigma(x))\sigma(x)$$



Think, pair, share

Patterns in backward flow

add gate: gradient distributor

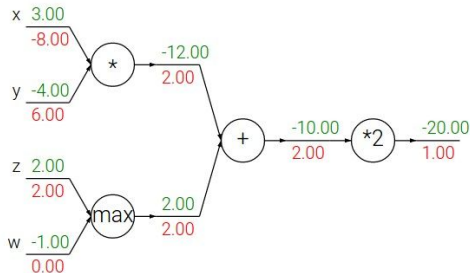


Think, pair, share

Patterns in backward flow

add gate: gradient distributor

Q: What is a **max** gate?

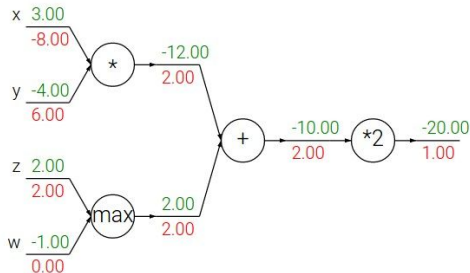


Think, pair, share

Patterns in backward flow

add gate: gradient distributor

max gate: gradient router



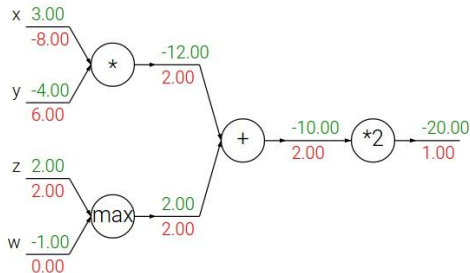
Think, pair, share

Patterns in backward flow

add gate: gradient distributor

max gate: gradient router

Q: What is a **mul** gate?



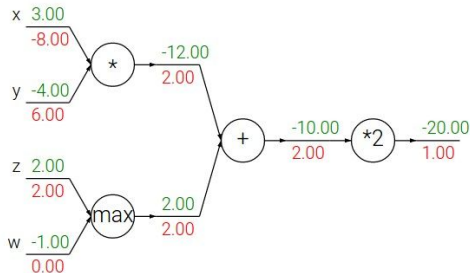
Think, pair, share

Patterns in backward flow

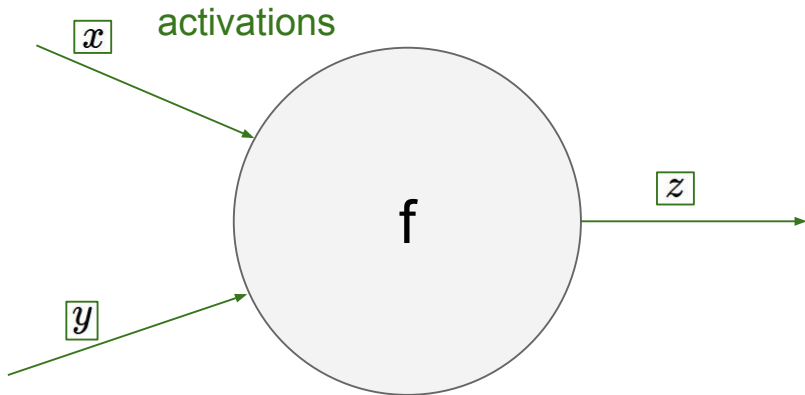
add gate: gradient distributor

max gate: gradient router

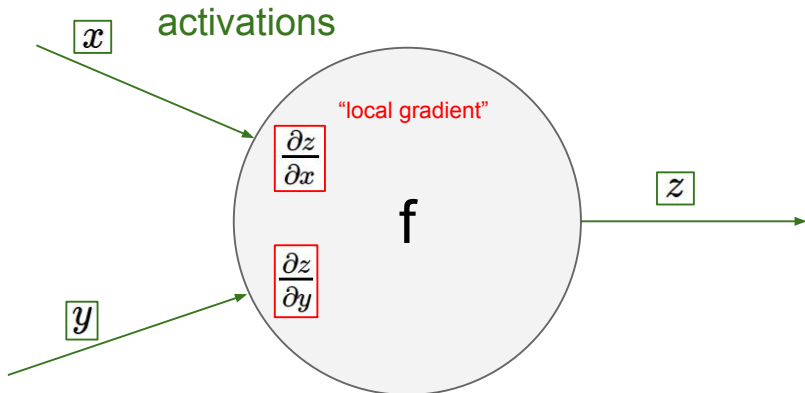
mul gate: gradient switcher



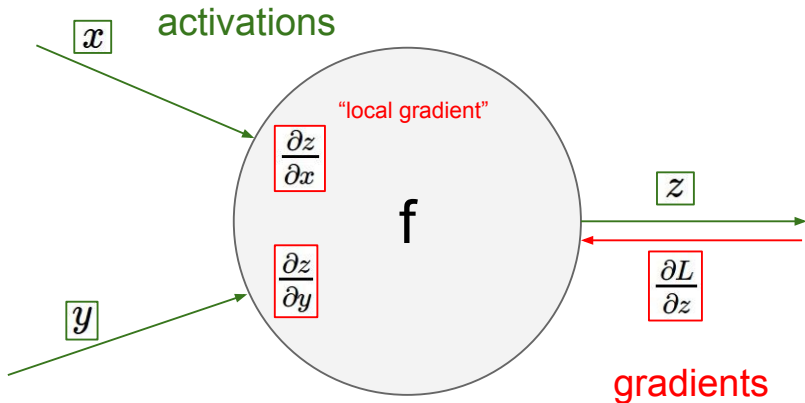
BP at one node



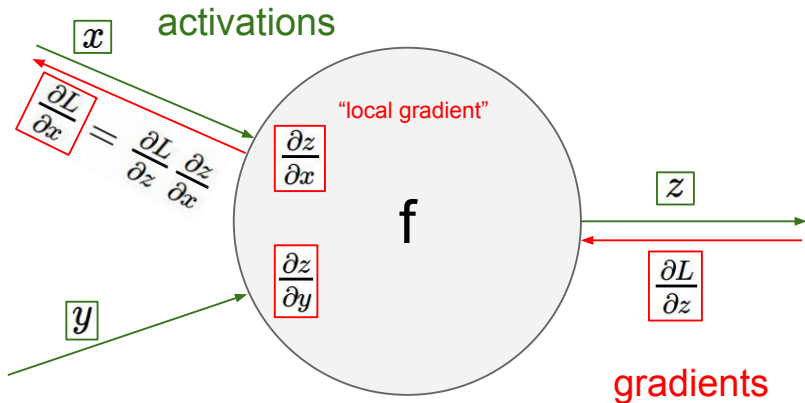
BP at one node



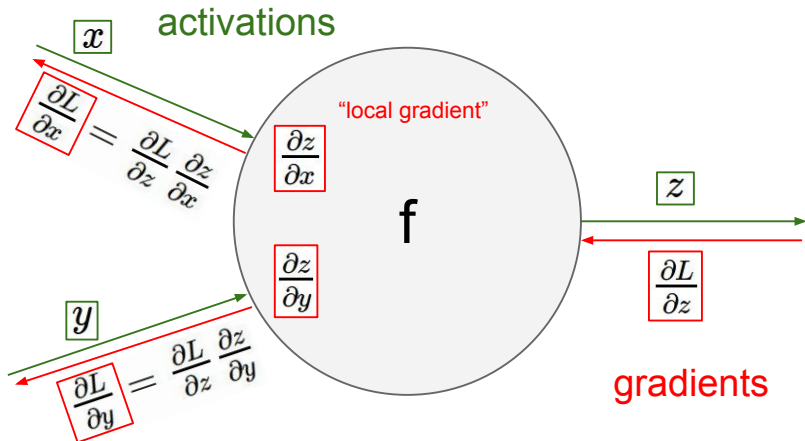
BP at one node



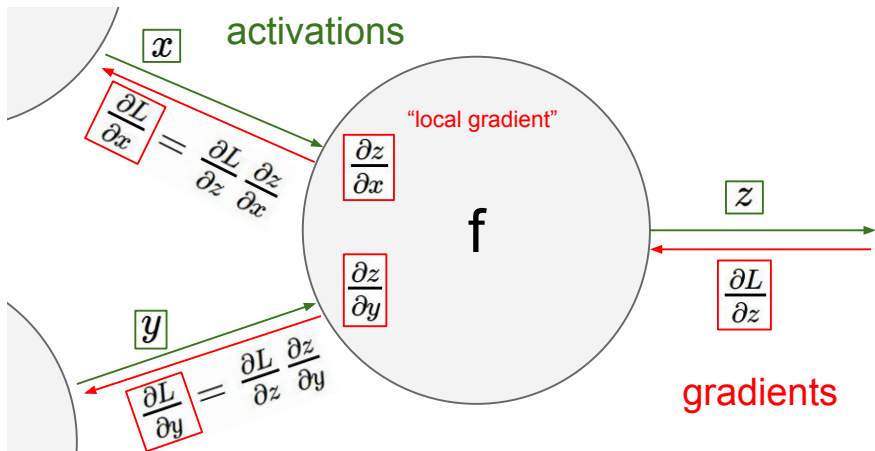
BP at one node



BP at one node



BP at one node

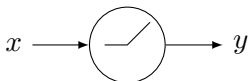


More examples: RELU

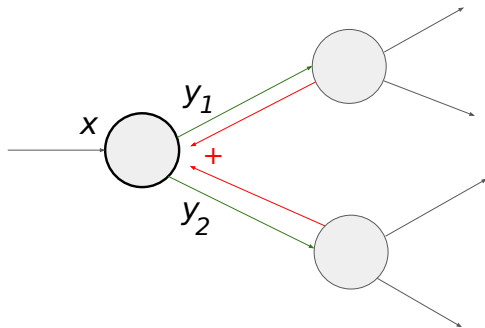
- Consider a “half-linear” function with negative side chopped off. That is,

$$f(x) = \begin{cases} x & x \geq 0 \\ 0 & \text{otherwise} \end{cases}$$

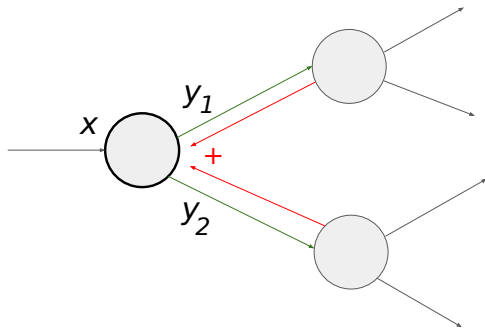
- This is known to be the rectified linear unit (RELU)
- How should the gradient be propagated back?



Merging gradients



Merging gradients

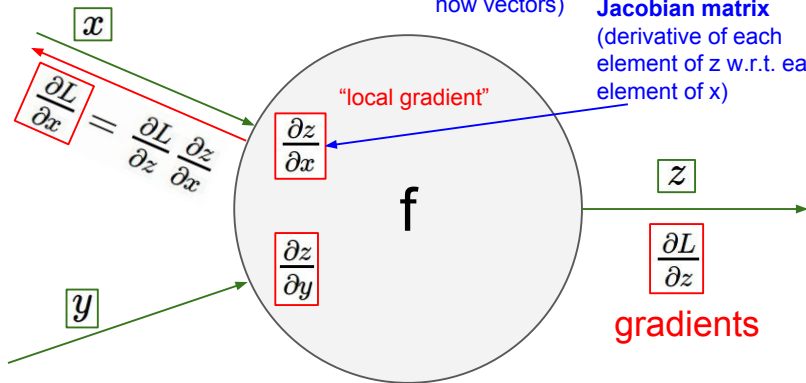


$$\frac{\partial L(y_1(x), y_2(x))}{\partial x} = \frac{\partial L}{\partial y_1} \frac{\partial y_1}{\partial x} + \frac{\partial L}{\partial y_2} \frac{\partial y_2}{\partial x}$$

Handling vector variables

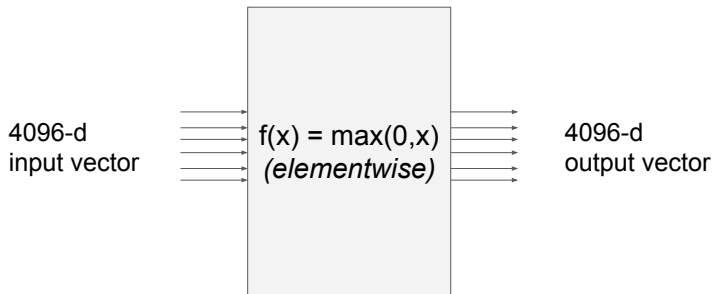
Gradients for vectorized code

(x,y,z are now vectors)

This is now the **Jacobian matrix** (derivative of each element of z w.r.t. each element of x)

Handling vector variables

Vectorized operations



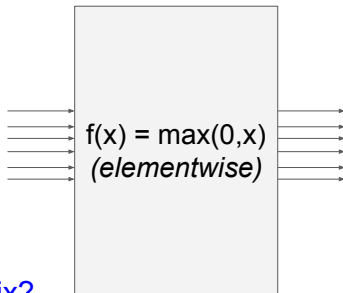
Handling vector variables

Vectorized operations

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}} \frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector



4096-d
output vector

Q: what is the
size of the
Jacobian matrix?

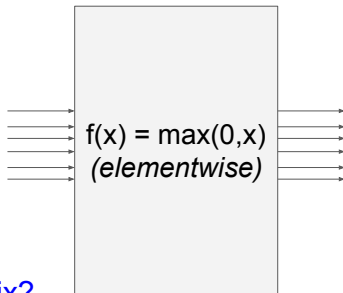
Handling vector variables

Vectorized operations

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}} \frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector



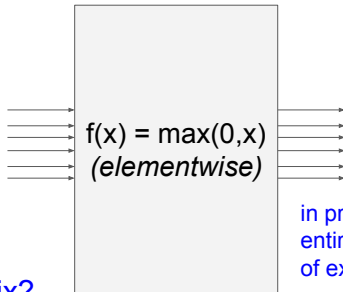
4096-d
output vector

Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Handling vector variables

Vectorized operations

4096-d
input vector



4096-d
output vector

Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

in practice we process an
entire minibatch (e.g. 100)
of examples at one time:

i.e. Jacobian would technically be a
[409,600 x 409,600] matrix :)

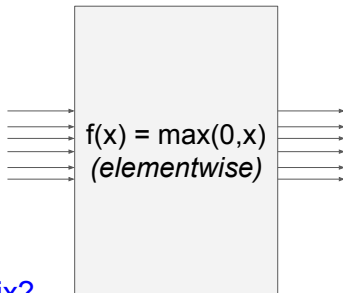
Handling vector variables

Vectorized operations

$$\frac{\partial L}{\partial x} = \boxed{\frac{\partial f}{\partial x}} \frac{\partial L}{\partial f}$$

Jacobian matrix

4096-d
input vector



4096-d
output vector

Q: what is the
size of the
Jacobian matrix?
[4096 x 4096!]

Q2: what does it
look like?

Handling vector variables

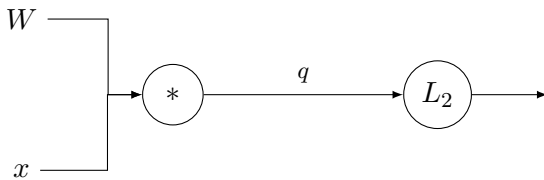
A vectorized example: $f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$

Handling vector variables

A vectorized example: $f(\underset{\mathbb{R}^n}{x}, \underset{\mathbb{R}^{n \times n}}{W}) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$

Handling vector variables

A vectorized example: $f(\underset{\mathbb{R}^n}{x}, \underset{\mathbb{R}^{n \times n}}{W}) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$

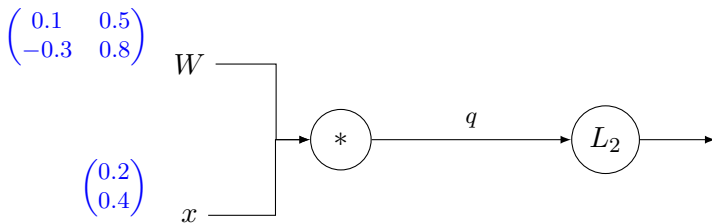


$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

Handling vector variables

A vectorized example: $f(\underset{\mathbb{R}^n}{x}, \underset{\mathbb{R}^{n \times n}}{W}) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$



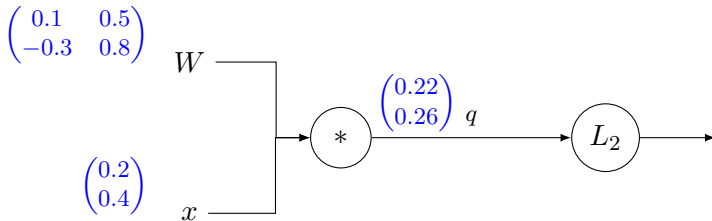
$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

Handling vector variables

A vectorized example: $f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$

$\underbrace{\quad}_{\mathbb{R}^n}$ $\underbrace{\quad}_{\mathbb{R}^{n \times n}}$



$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

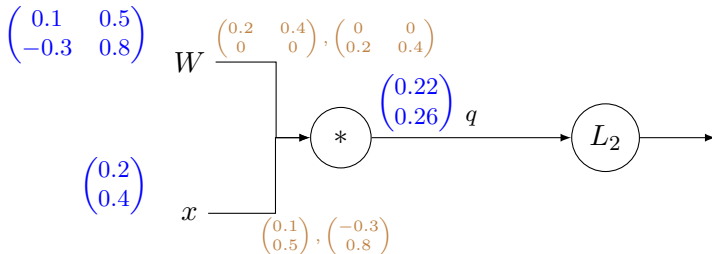
$$\frac{\partial q_k}{\partial W_{i,j}} = \delta_{i,k} x_j$$

$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

Handling vector variables

A vectorized example: $f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$

$\underbrace{\quad}_{\mathbb{R}^n}$ $\underbrace{\quad}_{\mathbb{R}^{n \times n}}$



$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

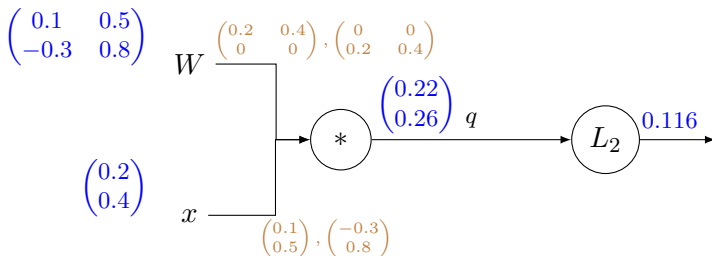
$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial q_k}{\partial W_{i,j}} = \delta_{i,k} x_j$$

$$\frac{\partial q_k}{\partial x_i} = W_{k,i}$$

Handling vector variables

A vectorized example: $f(\underbrace{x}_{\mathbb{R}^n}, \underbrace{W}_{\mathbb{R}^{n \times n}}) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$



$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

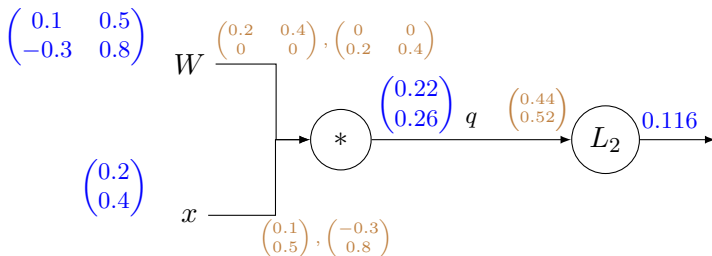
$$\frac{\partial f}{\partial q_i} = 2q_i$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

Handling vector variables

A vectorized example: $f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$

$\underbrace{\quad}_{\mathbb{R}^n} \quad \underbrace{\quad}_{\mathbb{R}^{n \times n}}$



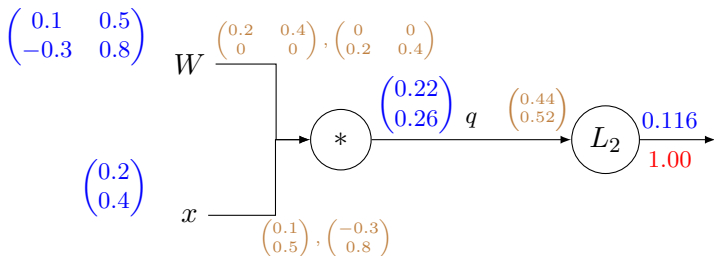
$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial f}{\partial q_i} = 2q_i$$

Handling vector variables

A vectorized example: $f(\underbrace{x}_{\mathbb{R}^n}, \underbrace{W}_{\mathbb{R}^{n \times n}}) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$



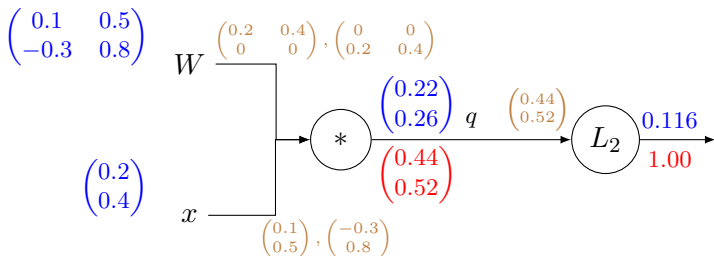
$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial f}{\partial q_i} = 1.00 \cdot \frac{\partial f}{\partial q_i}$$

Handling vector variables

A vectorized example: $f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$
 $\mathbb{R}^n \quad \mathbb{R}^{n \times n}$



$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

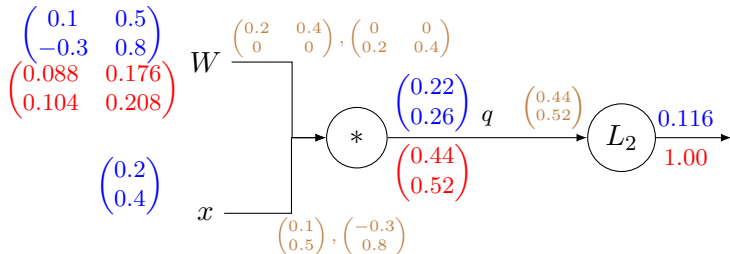
$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial f}{\partial q_i} = 1.00 \cdot \frac{\partial f}{\partial q_i}$$

Handling vector variables

A vectorized example: $f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$

$\underbrace{\quad}_{\mathbb{R}^n}$ $\underbrace{\quad}_{\mathbb{R}^{n \times n}}$



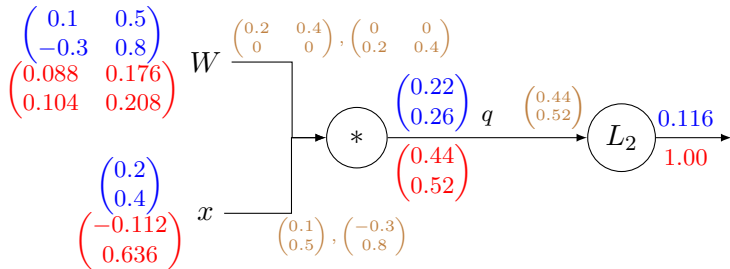
$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial f}{\partial W_{i,j}} = \frac{\partial f}{\partial q_1} \frac{\partial q_1}{\partial W_{i,j}} + \frac{\partial f}{\partial q_2} \frac{\partial q_2}{\partial W_{i,j}}$$

Handling vector variables

A vectorized example: $f(x, W) = \|Wx\|^2 = \sum_{i=1}^n (Wx)_i^2$
 $\mathbb{R}^n \quad \mathbb{R}^{n \times n}$



$$q = Wx = \begin{pmatrix} W_{1,1}x_1 + \cdots + W_{1,n}x_n \\ \vdots \\ W_{n,1}x_1 + \cdots + W_{n,n}x_n \end{pmatrix}$$

$$f(q) = \|q\|^2 = q_1^2 + \cdots + q_n^2$$

$$\frac{\partial f}{\partial x_i} = \frac{\partial f}{\partial q_1} \frac{\partial q_1}{\partial x_i} + \frac{\partial f}{\partial q_2} \frac{\partial q_2}{\partial x_i}$$

More examples: IoU (reference)

- Intersection over union is commonly used to quantify segmentation quality for image segmentation
- For pixel v , X_v is the estimated mask and $Y_v \in \{0, 1\}$ is the ground truth

More examples: IoU (reference)

- Intersection over union is commonly used to quantify segmentation quality for image segmentation
- For pixel v , X_v is the estimated mask and $Y_v \in \{0, 1\}$ is the ground truth
- $IoU(X) = \frac{I(X)}{U(X)}$, where $I(X) \approx \sum_v X_v Y_v$ and $U(X) \approx \sum_v (X_v + Y_v - X_v Y_v)$
- $\frac{\partial IoU(X)}{\partial X_v}$

More examples: IoU (reference)

- Intersection over union is commonly used to quantify segmentation quality for image segmentation
- For pixel v , X_v is the estimated mask and $Y_v \in \{0, 1\}$ is the ground truth
- $IoU(X) = \frac{I(X)}{U(X)}$, where $I(X) \approx \sum_v X_v Y_v$ and $U(X) \approx \sum_v (X_v + Y_v - X_v Y_v)$
- $\frac{\partial IoU(X)}{\partial X_v} = \frac{U(X) \frac{\partial I(X)}{\partial X_v} - I(X) \frac{\partial U(X)}{\partial X_v}}{U^2(X)}$

More examples: IoU (reference)

- Intersection over union is commonly used to quantify segmentation quality for image segmentation
- For pixel v , X_v is the estimated mask and $Y_v \in \{0, 1\}$ is the ground truth
- $IoU(X) = \frac{I(X)}{U(X)}$, where $I(X) \approx \sum_v X_v Y_v$ and $U(X) \approx \sum_v (X_v + Y_v - X_v Y_v)$
- $\frac{\partial IoU(X)}{\partial X_v} = \frac{U(X) \frac{\partial I(X)}{\partial X_v} - I(X) \frac{\partial U(X)}{\partial X_v}}{U^2(X)} = \frac{U(X) Y_v - I(X) (1 - Y_v)}{U(X)^2}$

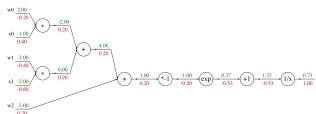
More examples: IoU (reference)

- Intersection over union is commonly used to quantify segmentation quality for image segmentation
 - For pixel v , X_v is the estimated mask and $Y_v \in \{0, 1\}$ is the ground truth
 - $IoU(X) = \frac{I(X)}{U(X)}$, where $I(X) \approx \sum_v X_v Y_v$ and $U(X) \approx \sum_v (X_v + Y_v - X_v Y_v)$
 - $\frac{\partial IoU(X)}{\partial X_v} = \frac{U(X) \frac{\partial I(X)}{\partial X_v} - I(X) \frac{\partial U(X)}{\partial X_v}}{U^2(X)} = \frac{U(X) Y_v - I(X) (1 - Y_v)}{U(X)^2}$
- $$\Rightarrow \frac{\partial IoU(X)}{\partial X_v} = \begin{cases} \frac{1}{U(X)} & Y_v = 1 \\ -\frac{I(X)}{U(X)^2} & Y_v = 0 \end{cases}$$

Implementation

Modularized implementation: forward / backward API

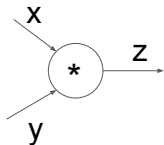
Graph (or Net) object (*rough psuedo code*)



```
class ComputationalGraph(object):
    #...
    def forward(inputs):
        # 1. [pass inputs to input gates...]
        # 2. forward the computational graph:
        for gate in self.graph.nodes_topologically_sorted():
            gate.forward()
        return loss # the final gate in the graph outputs the loss
    def backward():
        for gate in reversed(self.graph.nodes_topologically_sorted()):
            gate.backward() # little piece of backprop (chain rule applied)
        return inputs_gradients
```

Implementation

Modularized implementation: forward / backward API



(x,y,z are scalars)

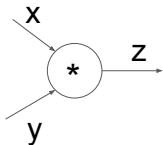
```
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        return z
    def backward(dz):
        # dx = ... #todo
        # dy = ... #todo
        return [dx, dy]
```

$$\frac{\partial L}{\partial z}$$

$$\frac{\partial L}{\partial x}$$

Implementation

Modularized implementation: forward / backward API



(x,y,z are scalars)

```
class MultiplyGate(object):
    def forward(x,y):
        z = x*y
        self.x = x # must keep these around!
        self.y = y
        return z
    def backward(dz):
        dx = self.y * dz # [dz/dx * dL/dz]
        dy = self.x * dz # [dz/dy * dL/dz]
        return [dx, dy]
```

Remark of BP

- During the **forward** pass, each computing unit will evaluate the output and also the corresponding local derivatives of the output w.r.t. the inputs

Remark of BP

- During the **forward** pass, each computing unit will evaluate the output and also the corresponding local derivatives of the output w.r.t. the inputs
- During the **backward** pass, the local derivatives and the evaluated outputs will be “consumed” to compute the overall derivatives

Remark of BP

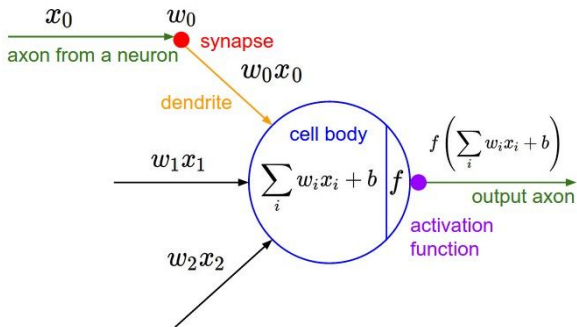
- During the **forward** pass, each computing unit will evaluate the output and also the corresponding local derivatives of the output w.r.t. the inputs
- During the **backward** pass, the local derivatives and the evaluated outputs will be “consumed” to compute the overall derivatives
 - For a large network, there can be a large spike of memory consumption during the forward pass

Remark of BP

- During the **forward** pass, each computing unit will evaluate the output and also the corresponding local derivatives of the output w.r.t. the inputs
- During the **backward** pass, the local derivatives and the evaluated outputs will be “consumed” to compute the overall derivatives
 - For a large network, there can be a large spike of memory consumption during the forward pass
- Note that BP only computes the gradients. It does not perform the optimization. Sometimes you may hear people said that they trained their networks with BP. What they said was not literally right. We will discuss more on optimizer later today

Activation functions

Activation Functions

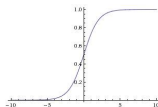


Activation functions

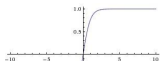
Activation Functions

Sigmoid

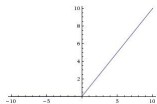
$$\sigma(x) = 1/(1 + e^{-x})$$



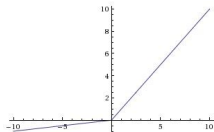
tanh tanh(x)



ReLU max(0,x)



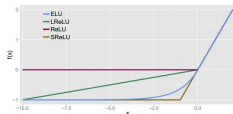
Leaky ReLU max(0.1x, x)



Maxout $\max(w_1^T x + b_1, w_2^T x + b_2)$

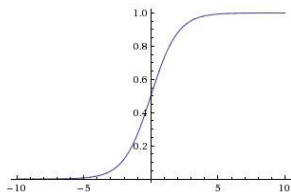
ELU

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$



Activation functions

Activation Functions



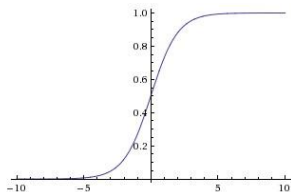
Sigmoid

$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

Activation functions

Activation Functions



Sigmoid

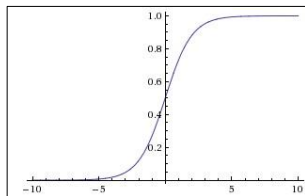
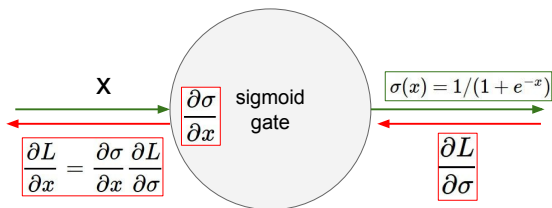
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients

Activation functions



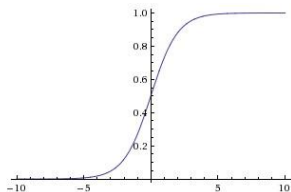
What happens when $x = -10$?

What happens when $x = 0$?

What happens when $x = 10$?

Activation functions

Activation Functions



Sigmoid

$$\sigma(x) = 1/(1 + e^{-x})$$

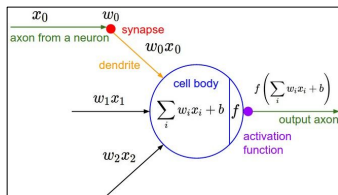
- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered

Activation functions

Consider what happens when the input to a neuron (x) is always positive:



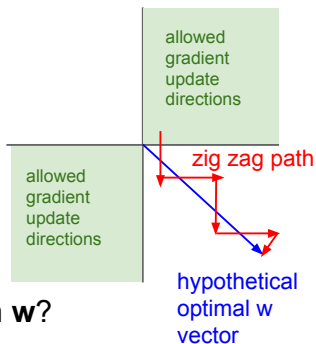
$$f\left(\sum_i w_i x_i + b\right)$$

What can we say about the gradients on \mathbf{w} ?

Activation functions

Consider what happens when the input to a neuron is always positive...

$$f\left(\sum_i w_i x_i + b\right)$$

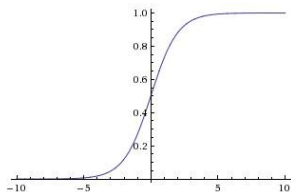


What can we say about the gradients on \mathbf{w} ?

Always all positive or all negative :(
(this is also why you want zero-mean data!)

Activation functions

Activation Functions



Sigmoid

$$\sigma(x) = 1/(1 + e^{-x})$$

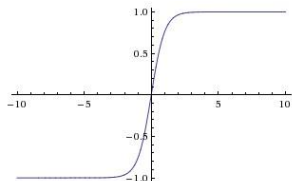
- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3. $\exp()$ is a bit compute expensive

Activation functions

Activation Functions



$\tanh(x)$

- Squashes numbers to range $[-1,1]$
- zero centered (nice)
- still kills gradients when saturated :(

[LeCun et al., 1991]

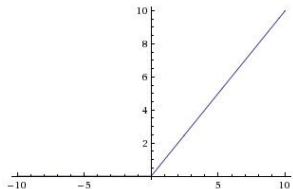
Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 37

20 Jan 2016

Activation functions

Activation Functions



- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

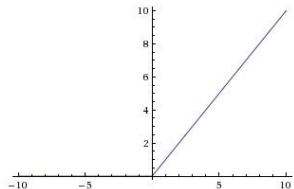
ReLU

(Rectified Linear Unit)

[Krizhevsky et al., 2012]

Activation functions

Activation Functions



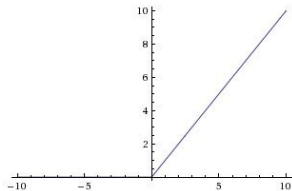
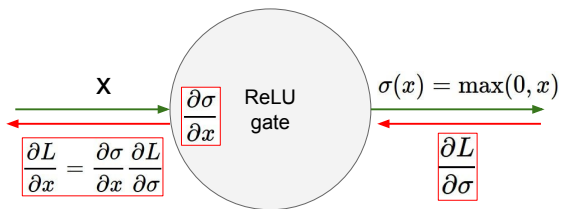
ReLU

(Rectified Linear Unit)

- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Not zero-centered output
- An annoyance:

hint: what is the gradient when $x < 0$?

Activation functions

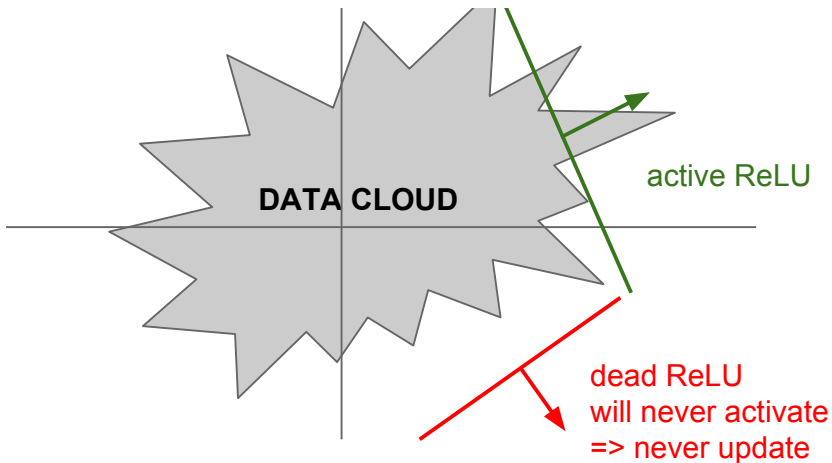


What happens when $x = -10$?

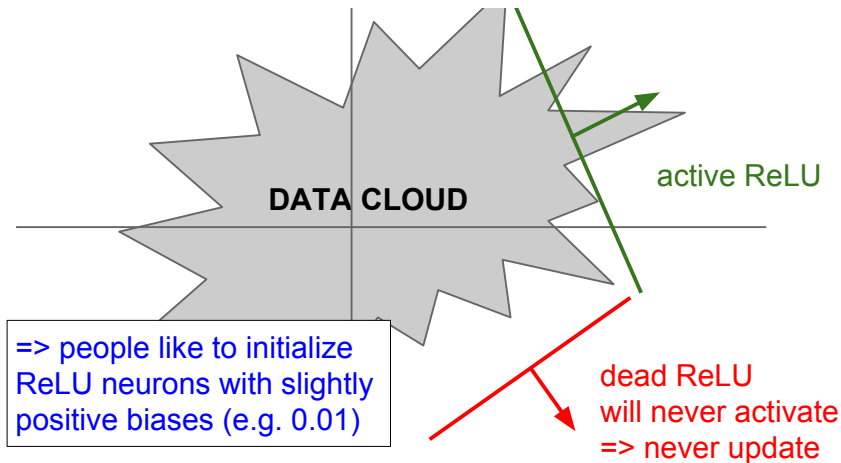
What happens when $x = 0$?

What happens when $x = 10$?

Activation functions



Activation functions

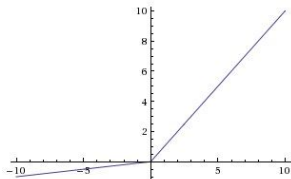


Activation functions

Activation Functions

[Mass et al., 2013]

[He et al., 2015]



- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”**.

Leaky ReLU

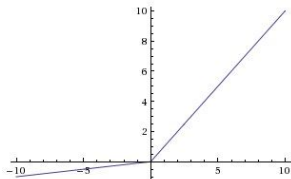
$$f(x) = \max(0.01x, x)$$

Activation functions

Activation Functions

[Mass et al., 2013]

[He et al., 2015]



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

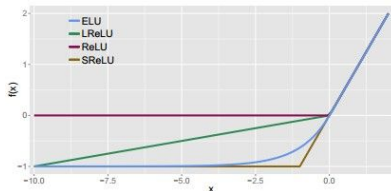
backprop into α
(parameter)

Activation functions

Activation Functions

[Clevert et al., 2015]

Exponential Linear Units (ELU)



$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

- All benefits of ReLU
- Does not die
- Closer to zero mean outputs
- Computation requires $\exp()$

Activation functions

Maxout "Neurons" [Goodfellow et al., 2013]

- Try to generalize ReLU and leaky ReLU

$$\max(\mathbf{w}_1^T \mathbf{x} + b_1, \mathbf{w}_2^T \mathbf{x} + b_2)$$

Activation functions

Maxout "Neurons" [Goodfellow et al., 2013]

- Try to generalize ReLU and leaky ReLU

$$\max(\mathbf{w}_1^T \mathbf{x} + b_1, \mathbf{w}_2^T \mathbf{x} + b_2)$$

Pros

- Linear regime
- Does not saturate
- Does not die

Activation functions

Maxout "Neurons" [Goodfellow et al., 2013]

- Try to generalize ReLU and leaky ReLU

$$\max(\mathbf{w}_1^T \mathbf{x} + b_1, \mathbf{w}_2^T \mathbf{x} + b_2)$$

Pros

- Linear regime
- Does not saturate
- Does not die

Cons

- Double amount of parameters

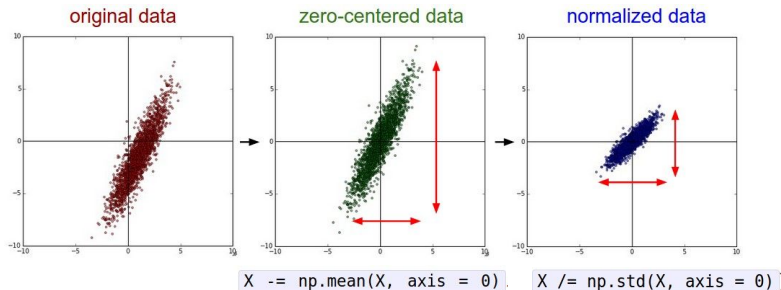
Activation functions

TLDR: In practice:

- Use **ReLU**. Be careful with your learning rates
- Try out **Leaky ReLU / Maxout / ELU**
- Try out **tanh** but don't expect much
- **Don't use sigmoid**

Input preprocessing

Step 1: Preprocess the data

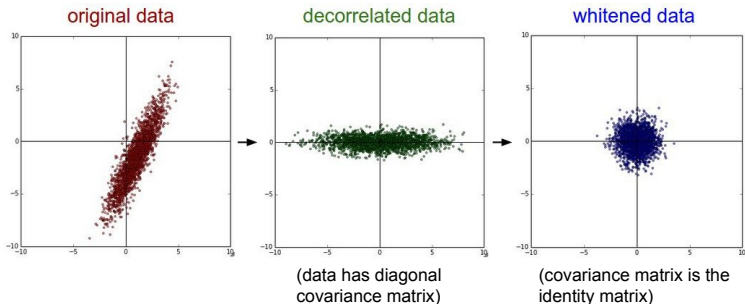


(Assume X [NxD] is data matrix,
each example in a row)

Input preprocessing

Step 1: Preprocess the data

In practice, you may also see **PCA** and **Whitening** of the data



Input preprocessing

TLDR: In practice for Images: center only

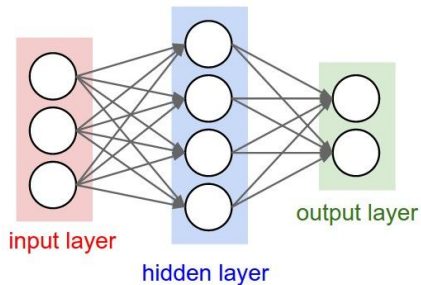
e.g. consider CIFAR-10 example with [32,32,3] images

- Subtract the mean image (e.g. AlexNet)
(mean image = [32,32,3] array)
- Subtract per-channel mean (e.g. VGGNet)
(mean along each channel = 3 numbers)

Not common to normalize
variance, to do PCA or
whitening

Weight initialization

- Q: what happens when $W=0$ init is used?



Weight initialization

- First idea: **Small random numbers**
(gaussian with zero mean and $1e-2$ standard deviation)

```
W = 0.01* np.random.randn(D,H)
```

Weight initialization

- First idea: **Small random numbers**
(gaussian with zero mean and $1e-2$ standard deviation)

```
W = 0.01* np.random.randn(D,H)
```

Works ~okay for small networks, but can lead to non-homogeneous distributions of activations across the layers of a network.

Weight initialization

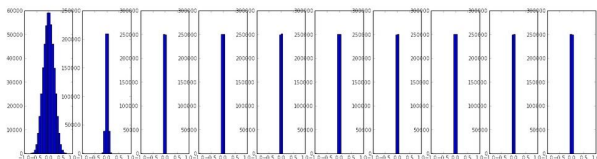
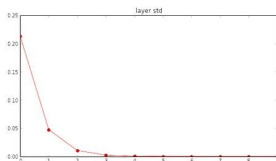
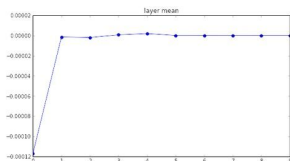
Let's look at some activation statistics

- 10 layers
- 500 neurons per layer
- $\tanh(\cdot)$ for activation
- $W = 0.01 * \text{np.random.randn}(\text{fan_in}, \text{fan_out})$ as described in the last slide

Weight initialization

```

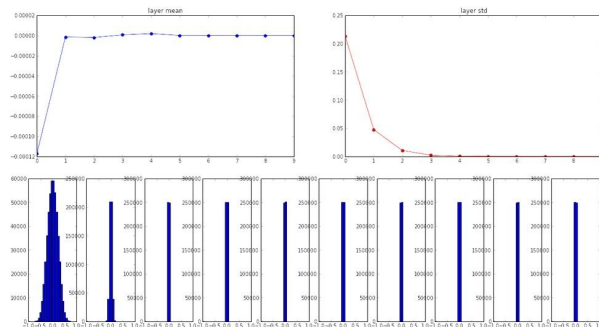
input layer had mean 0.000927 and std 0.998388
hidden layer 1 had mean -0.000117 and std 0.213081
hidden layer 2 had mean -0.000001 and std 0.047551
hidden layer 3 had mean -0.000002 and std 0.010630
hidden layer 4 had mean 0.000001 and std 0.002378
hidden layer 5 had mean 0.000002 and std 0.000532
hidden layer 6 had mean -0.000000 and std 0.000119
hidden layer 7 had mean 0.000000 and std 0.000026
hidden layer 8 had mean -0.000000 and std 0.000006
hidden layer 9 had mean 0.000000 and std 0.000001
hidden layer 10 had mean -0.000000 and std 0.000000
  
```



Weight initialization

```

input layer had mean 0.000927 and std 0.998388
hidden layer 1 had mean -0.000117 and std 0.213081
hidden layer 2 had mean -0.000001 and std 0.047551
hidden layer 3 had mean -0.000002 and std 0.010630
hidden layer 4 had mean 0.000001 and std 0.002378
hidden layer 5 had mean 0.000002 and std 0.000532
hidden layer 6 had mean -0.000000 and std 0.000119
hidden layer 7 had mean 0.000000 and std 0.000026
hidden layer 8 had mean -0.000000 and std 0.000006
hidden layer 9 had mean 0.000000 and std 0.000001
hidden layer 10 had mean -0.000000 and std 0.000000
  
```



All activations become zero!

Q: think about the backward pass. What do the gradients look like?

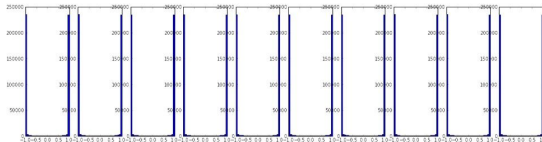
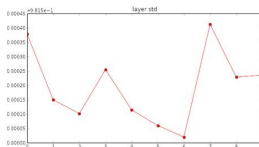
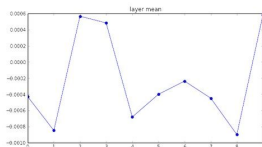
Hint: think about backward pass for a $W \cdot X$ gate.

Weight initialization

```
W = np.random.randn(fan_in, fan_out) * 1.0 # layer initialization
```

```
input layer had mean 0.001800 and std 1.001311
hidden layer 1 had mean -0.000430 and std 0.981879
hidden layer 2 had mean -0.000849 and std 0.981649
hidden layer 3 had mean 0.000566 and std 0.981601
hidden layer 4 had mean -0.000483 and std 0.981755
hidden layer 5 had mean -0.000682 and std 0.981614
hidden layer 6 had mean -0.000401 and std 0.981560
hidden layer 7 had mean -0.000237 and std 0.981520
hidden layer 8 had mean -0.000448 and std 0.981913
hidden layer 9 had mean -0.000899 and std 0.981728
hidden layer 10 had mean 0.000584 and std 0.981736
```

*1.0 instead of *0.01



Almost all neurons completely saturated, either -1 and 1. Gradients will be all zero.

Variance calibration for linear layer

Assume linear activation and zero-mean weights and inputs. And number of inputs is n . Then,

$$\text{Var}(y) = \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i)$$

Variance calibration for linear layer

Assume linear activation and zero-mean weights and inputs. And number of inputs is n . Then,

$$\begin{aligned}\text{Var}(y) &= \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i) \\ &= \sum_i^n [E(w_i)]^2 \text{Var}(x_i) + E[(x_i)]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i)\end{aligned}$$

$$\text{Var}(XY) = E[X]^2 \text{Var}(Y) + E[Y]^2 \text{Var}(X) + \text{Var}(X) \text{Var}(Y)$$

$$\text{Var}(XY) = E[(XY)^2] - E[XY]^2$$

$$\begin{aligned} \text{Var}(XY) = \\ E[X]^2 \text{Var}(Y) + E[Y]^2 \text{Var}(X) + \text{Var}(X) \text{Var}(Y) \end{aligned}$$

$$\begin{aligned} \text{Var}(XY) &= E[(XY)^2] - E[XY]^2 \\ &= E[X^2]E[Y^2] - E[X]^2E[Y]^2 \end{aligned}$$

$$\begin{aligned} \text{Var}(XY) = \\ E[X]^2 \text{Var}(Y) + E[Y]^2 \text{Var}(X) + \text{Var}(X) \text{Var}(Y) \end{aligned}$$

$$\begin{aligned} \text{Var}(XY) &= E[(XY)^2] - E[XY]^2 \\ &= E[X^2]E[Y^2] - E[X]^2E[Y]^2 \end{aligned}$$

$$\begin{aligned} &\text{Var}(X)\text{Var}(Y) \\ &= (E[X^2] - E[X]^2)(E[Y^2] - E[Y]^2) \end{aligned}$$

$$\begin{aligned} \text{Var}(XY) = \\ E[X]^2 \text{Var}(Y) + E[Y]^2 \text{Var}(X) + \text{Var}(X) \text{Var}(Y) \end{aligned}$$

$$\begin{aligned} \text{Var}(XY) &= E[(XY)^2] - E[XY]^2 \\ &= E[X^2]E[Y^2] - E[X]^2E[Y]^2 \end{aligned}$$

$$\begin{aligned} &\text{Var}(X)\text{Var}(Y) \\ &= (E[X^2] - E[X]^2)(E[Y^2] - E[Y]^2) \\ &= E[X^2]E[Y^2] - E[X]^2E[Y^2] - E[X^2]E[Y]^2 + E[X]^2E[Y]^2 \end{aligned}$$

$$\text{Var}(XY) = E[X]^2 \text{Var}(Y) + E[Y]^2 \text{Var}(X) + \text{Var}(X) \text{Var}(Y)$$

$$\begin{aligned} \text{Var}(XY) &= E[(XY)^2] - E[XY]^2 \\ &= E[X^2]E[Y^2] - E[X]^2E[Y]^2 \end{aligned}$$

$$\begin{aligned} &\text{Var}(X)\text{Var}(Y) \\ &= (E[X^2] - E[X]^2)(E[Y^2] - E[Y]^2) \\ &= E[X^2]E[Y^2] - E[X]^2E[Y^2] - E[X^2]E[Y]^2 + E[X]^2E[Y]^2 \\ &= E[X^2]E[Y^2] - E[X]^2(E[Y^2] - E[Y]^2) \\ &\quad E[Y]^2(E[X^2] - E[X]^2) - E[X]^2E[Y]^2 \end{aligned}$$

$$\text{Var}(XY) = E[X]^2 \text{Var}(Y) + E[Y]^2 \text{Var}(X) + \text{Var}(X) \text{Var}(Y)$$

$$\begin{aligned} \text{Var}(XY) &= E[(XY)^2] - E[XY]^2 \\ &= E[X^2]E[Y^2] - E[X]^2E[Y]^2 \end{aligned}$$

$$\begin{aligned} &\text{Var}(X)\text{Var}(Y) \\ &= (E[X^2] - E[X]^2)(E[Y^2] - E[Y]^2) \\ &= E[X^2]E[Y^2] - E[X]^2E[Y^2] - E[X^2]E[Y]^2 + E[X]^2E[Y]^2 \\ &= E[X^2]E[Y^2] - E[X]^2(E[Y^2] - E[Y]^2) \\ &\quad E[Y]^2(E[X^2] - E[X]^2) - E[X]^2E[Y]^2 \\ &= \text{Var}(XY) - E[X]^2\text{Var}(Y) - E[Y]^2\text{Var}(X) \end{aligned}$$

Variance calibration for linear layer

Assume linear activation and zero-mean weights and inputs. And number of inputs is n . Then,

$$\begin{aligned}\text{Var}(y) &= \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i) \\ &= \sum_i^n E[w_i]^2 \text{Var}(x_i) + E[x_i]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i)\end{aligned}$$

Variance calibration for linear layer

Assume linear activation and zero-mean weights and inputs. And number of inputs is n . Then,

$$\begin{aligned}\text{Var}(y) &= \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i) \\ &= \sum_i^n E[w_i]^2 \text{Var}(x_i) + E[x_i]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i) \\ &= \sum_i^n \text{Var}(x_i) \text{Var}(w_i) \\ &= (n \text{Var}(w)) \text{Var}(x)\end{aligned}$$

Variance calibration for linear layer

Assume linear activation and zero-mean weights and inputs. And number of inputs is n . Then,

$$\begin{aligned}\text{Var}(y) &= \text{Var}\left(\sum_i^n w_i x_i\right) = \sum_i^n \text{Var}(w_i x_i) \\ &= \sum_i^n E[w_i]^2 \text{Var}(x_i) + E[x_i]^2 \text{Var}(w_i) + \text{Var}(x_i) \text{Var}(w_i) \\ &= \sum_i^n \text{Var}(x_i) \text{Var}(w_i) \\ &= (n \text{Var}(w)) \text{Var}(x)\end{aligned}$$

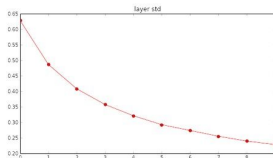
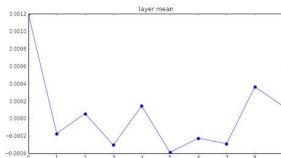
Thus, output will have same variance as input if $n \text{Var}(w) = 1$. This is known as Xavier weight initialization

Weight initialization

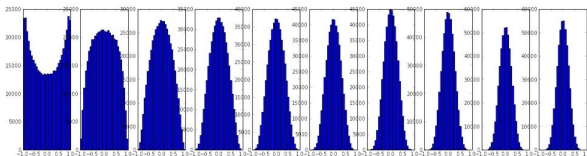
input layer had mean 0.001800 and std 1.001311
 hidden layer 1 had mean 0.001198 and std 0.627953
 hidden layer 2 had mean -0.000175 and std 0.486051
 hidden layer 3 had mean 0.000055 and std 0.407723
 hidden layer 4 had mean -0.000306 and std 0.357108
 hidden layer 5 had mean 0.000142 and std 0.320917
 hidden layer 6 had mean -0.000389 and std 0.292116
 hidden layer 7 had mean -0.000228 and std 0.273387
 hidden layer 8 had mean -0.000291 and std 0.254935
 hidden layer 9 had mean 0.000361 and std 0.239266
 hidden layer 10 had mean 0.000139 and std 0.228008

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in) # layer initialization
```

“Xavier initialization”
 [Glorot et al., 2010]



Reasonable initialization.
 (Mathematical derivation
 assumes linear activations)

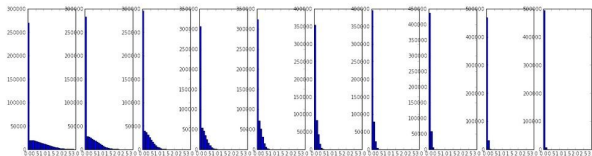
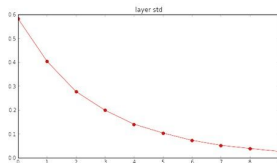
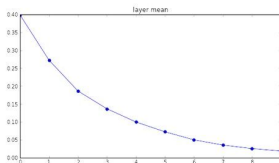


Weight initialization

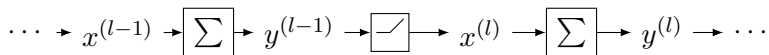
input layer had mean 0.009501 and std 0.999444
 hidden layer 1 had mean 0.398623 and std 0.582273
 hidden layer 2 had mean 0.272352 and std 0.403795
 hidden layer 3 had mean 0.186076 and std 0.276912
 hidden layer 4 had mean 0.136442 and std 0.198685
 hidden layer 5 had mean 0.099568 and std 0.140299
 hidden layer 6 had mean 0.072234 and std 0.103280
 hidden layer 7 had mean 0.049775 and std 0.072748
 hidden layer 8 had mean 0.035138 and std 0.051572
 hidden layer 9 had mean 0.025404 and std 0.038583
 hidden layer 10 had mean 0.018408 and std 0.026076

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in) # layer initialization
```

but when using the ReLU nonlinearity it breaks.



Variance calibration for ReLU

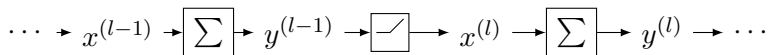


Note that it doesn't work when the activation layer is ReLU. But...¹

$$\text{Var}(y^{(l)}) = \text{Var} \left(\sum_i^n w_i^{(l)} x_i^{(l)} \right)$$

¹Note that $y^{(l)}$ now denotes the sum of input before going through the activation function.

Variance calibration for ReLU

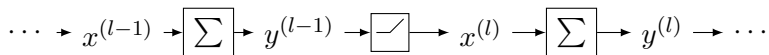


Note that it doesn't work when the activation layer is ReLU. But...¹

$$\text{Var}(y^{(l)}) = \text{Var}\left(\sum_i^n w_i^{(l)} x_i^{(l)}\right) = \sum_i^n \text{Var}(w_i^{(l)} x_i^{(l)}) = n \text{Var}(w^{(l)} x^{(l)})$$

¹Note that $y^{(l)}$ now denotes the sum of input before going through the activation function.

Variance calibration for ReLU

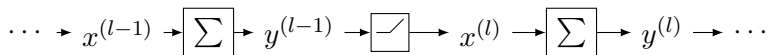


Note that it doesn't work when the activation layer is ReLU. But...¹

$$\begin{aligned} \text{Var}(y^{(l)}) &= \text{Var}\left(\sum_i^n w_i^{(l)} x_i^{(l)}\right) = \sum_i^n \text{Var}(w_i^{(l)} x_i^{(l)}) = n \text{Var}(w^{(l)} x^{(l)}) \\ &= n E[w^{(l)}]^2 \text{Var}(x^{(l)}) + n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \end{aligned}$$

¹Note that $y^{(l)}$ now denotes the sum of input before going through the activation function.

Variance calibration for ReLU

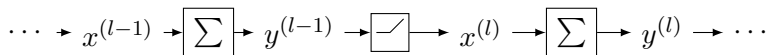


Note that it doesn't work when the activation layer is ReLU. But...¹

$$\begin{aligned} \text{Var}(y^{(l)}) &= \text{Var}\left(\sum_i^n w_i^{(l)} x_i^{(l)}\right) = \sum_i^n \text{Var}(w_i^{(l)} x_i^{(l)}) = n \text{Var}(w^{(l)} x^{(l)}) \\ &= n E[w^{(l)}]^2 \text{Var}(x^{(l)}) + n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \\ &= n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \end{aligned}$$

¹Note that $y^{(l)}$ now denotes the sum of input before going through the activation function.

Variance calibration for ReLU

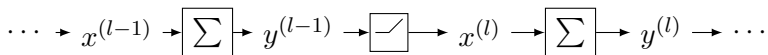


Note that it doesn't work when the activation layer is ReLU. But...¹

$$\begin{aligned} \text{Var}(y^{(l)}) &= \text{Var}\left(\sum_i^n w_i^{(l)} x_i^{(l)}\right) = \sum_i^n \text{Var}(w_i^{(l)} x_i^{(l)}) = n \text{Var}(w^{(l)} x^{(l)}) \\ &= n E[w^{(l)}]^2 \text{Var}(x^{(l)}) + n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \\ &= n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \\ &= n E[(x^{(l)})^2] \text{Var}(w^{(l)}) \end{aligned}$$

¹Note that $y^{(l)}$ now denotes the sum of input before going through the activation function.

Variance calibration for ReLU

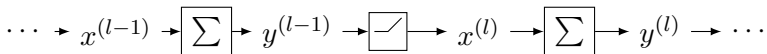


Note that it doesn't work when the activation layer is ReLU. But...¹

$$\begin{aligned} \text{Var}(y^{(l)}) &= \text{Var}\left(\sum_i^n w_i^{(l)} x_i^{(l)}\right) = \sum_i^n \text{Var}(w_i^{(l)} x_i^{(l)}) = n \text{Var}(w^{(l)} x^{(l)}) \\ &= n E[w^{(l)}]^2 \text{Var}(x^{(l)}) + n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \\ &= n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \\ &= n E[(x^{(l)})^2] \text{Var}(w^{(l)}) \\ &= n (\text{Var}(y^{(l-1)})/2) \text{Var}(w^{(l)}) = \left(\frac{n}{2} \text{Var}(w^{(l)})\right) \text{Var}(y^{(l-1)}) \end{aligned}$$

¹Note that $y^{(l)}$ now denotes the sum of input before going through the activation function.

Variance calibration for ReLU



Note that it doesn't work when the activation layer is ReLU. But...¹

$$\begin{aligned} \text{Var}(y^{(l)}) &= \text{Var}\left(\sum_i^n w_i^{(l)} x_i^{(l)}\right) = \sum_i^n \text{Var}(w_i^{(l)} x_i^{(l)}) = n \text{Var}(w^{(l)} x^{(l)}) \\ &= n E[w^{(l)}]^2 \text{Var}(x^{(l)}) + n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \\ &= n E[x^{(l)}]^2 \text{Var}(w^{(l)}) + n \text{Var}(x^{(l)}) \text{Var}(w^{(l)}) \\ &= n E[(x^{(l)})^2] \text{Var}(w^{(l)}) \\ &= n (\text{Var}(y^{(l-1)})/2) \text{Var}(w^{(l)}) = \left(\frac{n}{2} \text{Var}(w^{(l)})\right) \text{Var}(y^{(l-1)}) \end{aligned}$$

Variance of y conserved across a layer if $\frac{n}{2} \text{Var}(w) = 1$

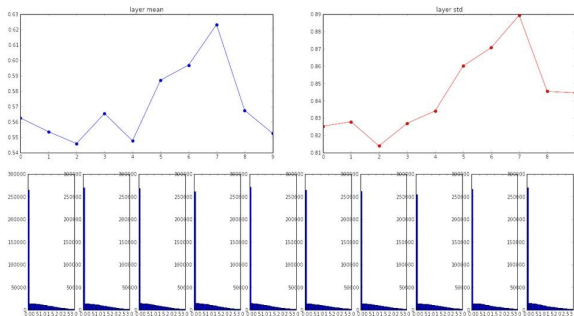
¹Note that $y^{(l)}$ now denotes the sum of input before going through the activation function.

Weight initialization

input layer had mean 0.000501 and std 0.999444
 hidden layer 1 had mean 0.562488 and std 0.825232
 hidden layer 2 had mean 0.553614 and std 0.827835
 hidden layer 3 had mean 0.545867 and std 0.813855
 hidden layer 4 had mean 0.565396 and std 0.826902
 hidden layer 5 had mean 0.547678 and std 0.834092
 hidden layer 6 had mean 0.587103 and std 0.860035
 hidden layer 7 had mean 0.596867 and std 0.870610
 hidden layer 8 had mean 0.623214 and std 0.809340
 hidden layer 9 had mean 0.567490 and std 0.845357
 hidden layer 10 had mean 0.552531 and std 0.844523

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in/2) # layer initialization
```

He et al., 2015
 (note additional /2)

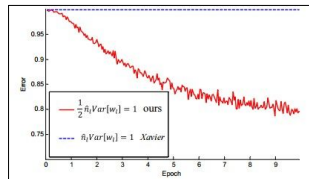
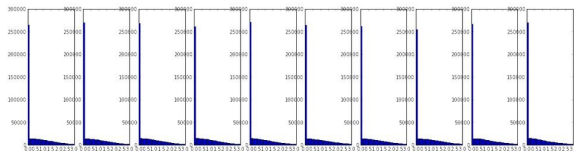
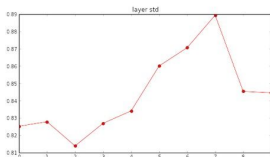
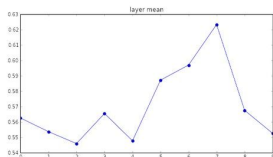


Weight initialization

input layer had mean 0.000501 and std 0.999444
 hidden layer 1 had mean 0.562488 and std 0.825232
 hidden layer 2 had mean 0.553614 and std 0.827835
 hidden layer 3 had mean 0.545867 and std 0.813855
 hidden layer 4 had mean 0.565396 and std 0.826902
 hidden layer 5 had mean 0.547678 and std 0.834092
 hidden layer 6 had mean 0.587103 and std 0.860035
 hidden layer 7 had mean 0.596867 and std 0.870610
 hidden layer 8 had mean 0.623214 and std 0.809340
 hidden layer 9 had mean 0.567490 and std 0.845357
 hidden layer 10 had mean 0.552531 and std 0.844523

```
W = np.random.randn(fan_in, fan_out) / np.sqrt(fan_in/2) # layer initialization
```

He et al., 2015
 (note additional /2)



Weight initialization

Proper initialization is an active area of research...

Understanding the difficulty of training deep feedforward neural networks

by Glorot and Bengio, 2010

Exact solutions to the nonlinear dynamics of learning in deep linear neural networks by

Saxe et al, 2013

Random walk initialization for training very deep feedforward networks by Sussillo and

Abbott, 2014

Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification by He et al., 2015

Data-dependent Initializations of Convolutional Neural Networks by Krähenbühl et al., 2015

All you need is a good init, Mishkin and Matas, 2015

...

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 64

20 Jan 2016

Batch normalization

Batch Normalization

[Ioffe and Szegedy, 2015]

“you want unit gaussian activations? just make them so.”

consider a batch of activations at some layer.
To make each dimension unit gaussian, apply:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

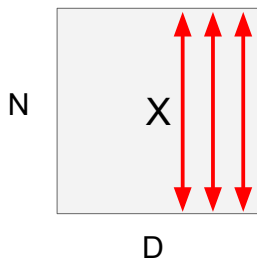
this is a vanilla
differentiable function...

Batch normalization

Batch Normalization

[Ioffe and Szegedy, 2015]

“you want unit gaussian activations?
just make them so.”



1. compute the empirical mean and variance independently for each dimension.

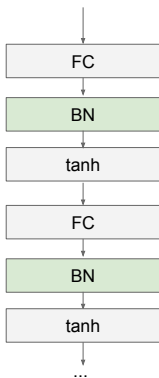
2. Normalize

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Batch normalization

Batch Normalization

[Ioffe and Szegedy, 2015]



Usually inserted after Fully Connected / (or Convolutional, as we'll see soon) layers, and before nonlinearity.

Problem: do we necessarily want a unit gaussian input to a tanh layer?

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

Batch normalization

Batch Normalization

[Ioffe and Szegedy, 2015]

Normalize:

$$\hat{x}^{(k)} = \frac{x^{(k)} - \mathbb{E}[x^{(k)}]}{\sqrt{\text{Var}[x^{(k)}]}}$$

And then allow the network to squash the range if it wants to:

$$y^{(k)} = \gamma^{(k)} \hat{x}^{(k)} + \beta^{(k)}$$

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\text{Var}[x^{(k)}]}$$

$$\beta^{(k)} = \mathbb{E}[x^{(k)}]$$

to recover the identity mapping.

Batch normalization

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \quad // \text{ scale and shift}$$

- Improves gradient flow through the network
- Allows higher learning rates
- Reduces the strong dependence on initialization
- Acts as a form of regularization in a funny way, and slightly reduces the need for dropout, maybe

Batch normalization

Batch Normalization

[Ioffe and Szegedy, 2015]

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;
Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

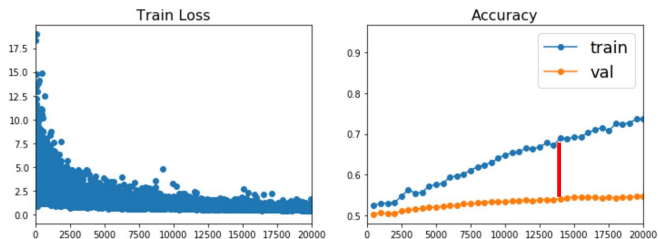
Note: at test time BatchNorm layer functions differently:

The mean/std are not computed based on the batch. Instead, a single fixed empirical mean of activations during training is used.

(e.g. can be estimated during training with running averages)

Reducing testing error

How to improve single-model performance?



Ensemble trick

1. Train multiple independent models
2. At test time average their results

Enjoy 2% extra performance

Ensemble trick

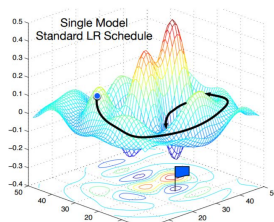
Fun Tips/Tricks:

- can also get a small boost from averaging multiple model checkpoints of a single model.

Ensemble trick

Model Ensembles: Tips and Tricks

Instead of training independent models, use multiple snapshots of a single model during training!



Loshchilov and Hutter, "SGDR: Stochastic gradient descent with restarts", arXiv 2016

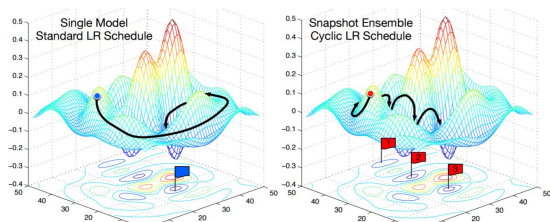
Huang et al, "Snapshot ensembles: train 1, get M for free", ICLR 2017

Figures copyright Yixuan Li and Geoff Pleiss, 2017. Reproduced with permission.

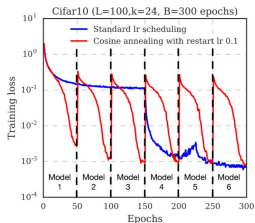
Ensemble trick

Model Ensembles: Tips and Tricks

Instead of training independent models, use multiple snapshots of a single model during training!



Loshchilov and Hutter, "SGDR: Stochastic gradient descent with restarts", arXiv 2016
 Huang et al, "Snapshot ensembles: train 1, get M for free", ICLR 2017
 Figures copyright Yixuan Li and Geoff Pleiss, 2017. Reproduced with permission.



Cyclic learning rate schedules can make this work even better!

Ensemble trick

Model Ensembles: Tips and Tricks

Instead of using actual parameter vector, keep a moving average of the parameter vector and use that at test time (Polyak averaging)

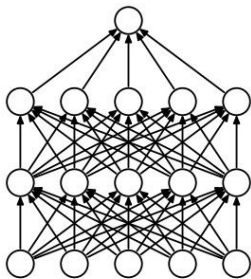
```
while True:
    data_batch = dataset.sample_data_batch()
    loss = network.forward(data_batch)
    dx = network.backward()
    x += - learning_rate * dx
    x_test = 0.995*x_test + 0.005*x # use for test set
```

Polyak and Juditsky, "Acceleration of stochastic approximation by averaging". SIAM Journal on Control and Optimization, 1992.

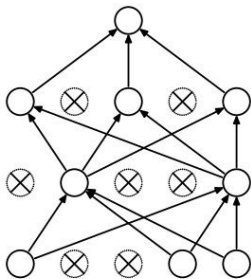
Dropout

Regularization: **Dropout**

“randomly set some neurons to zero in the forward pass”



(a) Standard Neural Net



(b) After applying dropout.

[Srivastava et al., 2014]

Dropout

```

p = 0.5 # probability of keeping a unit active. higher = less dropout

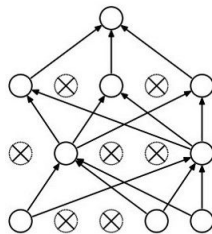
def train_step(X):
    """ X contains the data """

    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

```

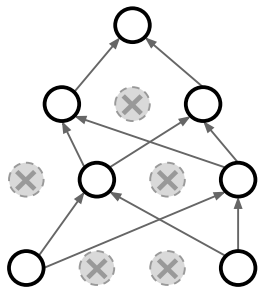
Example forward pass with a 3-layer network using dropout



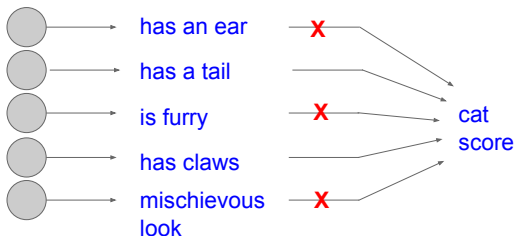
Dropout

Regularization: Dropout

How can this possibly be a good idea?



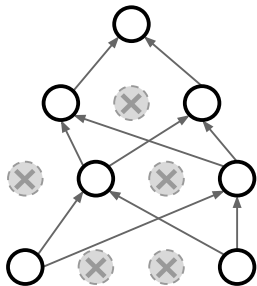
Forces the network to have a redundant representation;
Prevents co-adaptation of features



Dropout

Regularization: Dropout

How can this possibly be a good idea?



Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model

An FC layer with 4096 units has $2^{4096} \sim 10^{1233}$ possible masks!

Only $\sim 10^{82}$ atoms in the universe...

Dropout

Dropout: Test time

Dropout makes our output random!

$$\begin{array}{l}
 \text{Output} \\
 \text{(label)}
 \end{array}
 \quad
 \begin{array}{l}
 \text{Input} \\
 \text{(image)}
 \end{array}
 \quad
 \begin{array}{l}
 \text{Random} \\
 \text{mask}
 \end{array}$$

$$\boxed{y} = f_W(\boxed{x}, \boxed{z})$$

Want to “average out” the randomness at test-time

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

But this integral seems hard ...

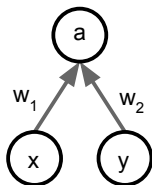
Dropout

Dropout: Test time

Want to approximate
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.



Dropout

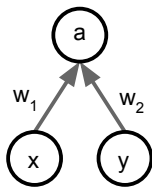
Dropout: Test time

Want to approximate
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.

At test time we have: $E[a] = w_1x + w_2y$



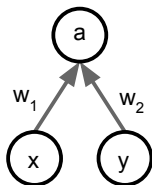
Dropout

Dropout: Test time

Want to approximate
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.



At test time we have: $E[a] = w_1x + w_2y$

During training we have:
$$E[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) = \frac{1}{2}(w_1x + w_2y)$$

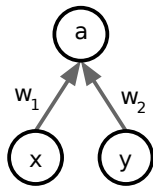
Dropout

Dropout: Test time

Want to approximate
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron.



At test time we have: $E[a] = w_1x + w_2y$

During training we have: $E[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y)$
 $+ \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y)$
 $= \frac{1}{2}(w_1x + w_2y)$

At test time, multiply
by probability p

Dropout

Dropout: Test time

```
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

At test time all neurons are active always

=> We must scale the activations so that for each neuron:

output at test time = expected output at training time

Dropout

```

""" Vanilla Dropout: Not recommended implementation (see notes below) """

p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """

    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients.. (not shown)
    # perform parameter update.. (not shown)

def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
    out = np.dot(W3, H2) + b3

```

Dropout Summary

drop in forward pass

scale at test time

Dropout

More common: “Inverted dropout”

```

p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensemble forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    out = np.dot(W3, H2) + b3

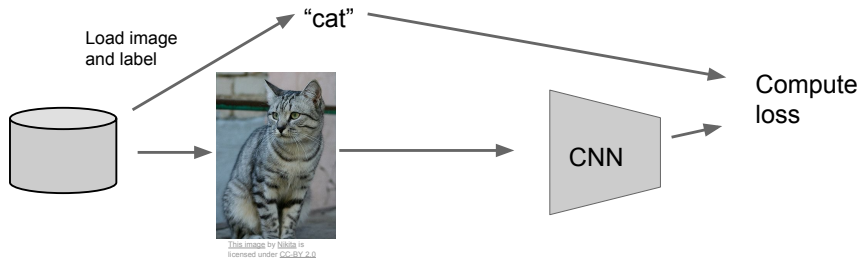
```

test time is unchanged!



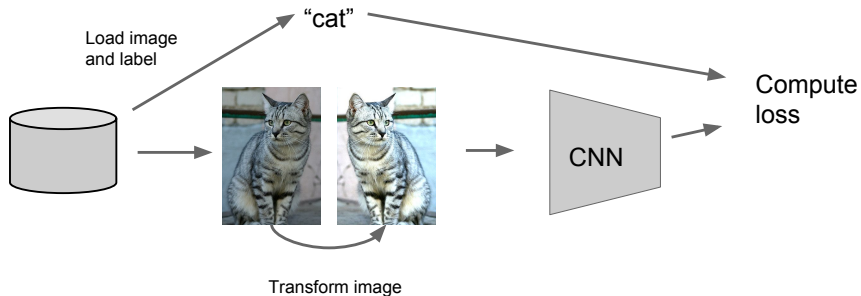
Data augmentation

Regularization: Data Augmentation



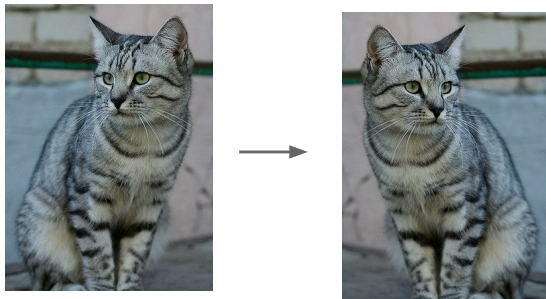
Data augmentation

Regularization: Data Augmentation



Data augmentation

Data Augmentation Horizontal Flips



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 76

April 25, 2017

Data augmentation

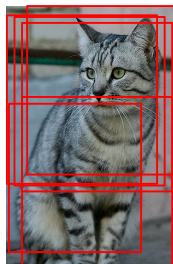
Data Augmentation

Random crops and scales

Training: sample random crops / scales

ResNet:

1. Pick random L in range $[256, 480]$
2. Resize training image, short side = L
3. Sample random 224×224 patch



Data augmentation

Data Augmentation

Random crops and scales

Training: sample random crops / scales

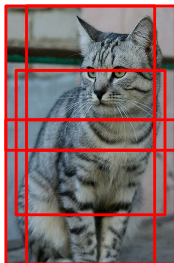
ResNet:

1. Pick random L in range $[256, 480]$
2. Resize training image, short side = L
3. Sample random 224×224 patch

Testing: average a fixed set of crops

ResNet:

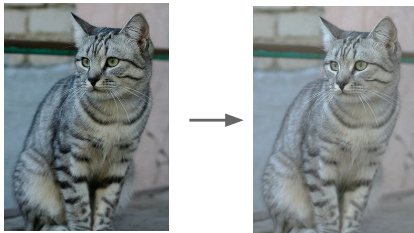
1. Resize image at 5 scales: $\{224, 256, 384, 480, 640\}$
2. For each size, use 10 224×224 crops: 4 corners + center, + flips



Data augmentation

Data Augmentation Color Jitter

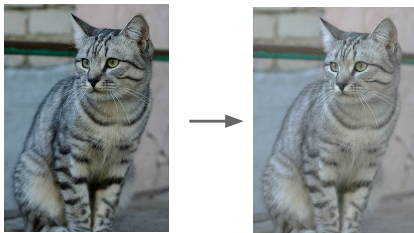
Simple: Randomize
contrast and brightness



Data augmentation

Data Augmentation Color Jitter

Simple: Randomize
contrast and brightness



More Complex:

1. Apply PCA to all [R, G, B] pixels in training set
2. Sample a “color offset” along principal component directions
3. Add offset to all pixels of a training image

(As seen in [Krizhevsky et al. 2012], ResNet, etc)

Data augmentation

Data Augmentation

Get creative for your problem!

Random mix/combinations of :

- translation
- rotation
- stretching
- shearing,
- lens distortions, ... (go crazy)

Other regularization techniques

Regularization: A common pattern

Training: Add some kind of randomness

$$y = f_W(x, z)$$

Testing: Average out randomness (sometimes approximate)

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Other regularization techniques

Regularization: A common pattern

Training: Add random noise

Testing: Marginalize over the noise

Examples:

Dropout

Batch Normalization

Data Augmentation

Other regularization techniques

Regularization: A common pattern

Training: Add random noise

Testing: Marginalize over the noise

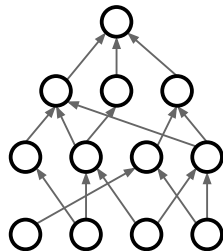
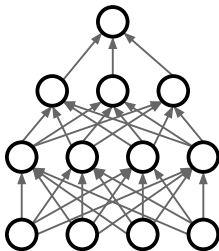
Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect



Wan et al, "Regularization of Neural Networks using DropConnect", ICML 2013

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 83

April 25, 2017

Other regularization techniques

Regularization: A common pattern

Training: Add random noise

Testing: Marginalize over the noise

Examples:

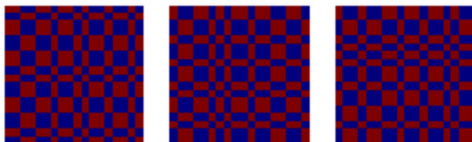
Dropout

Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling



Graham, "Fractional Max Pooling", arXiv 2014

Other regularization techniques

Regularization: A common pattern

Training: Add random noise

Testing: Marginalize over the noise

Examples:

Dropout

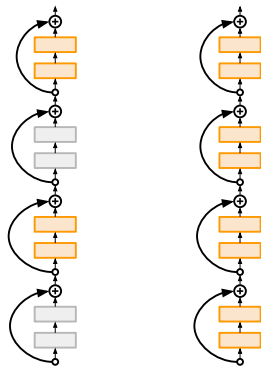
Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth



Huang et al, "Deep Networks with Stochastic Depth", ECCV 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 85

April 25, 2017

Optimizers

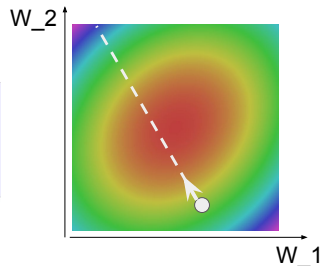
Optimization

```
# Vanilla Gradient Descent
```

```
while True:
```

```
    weights_grad = evaluate_gradient(loss_fun, data, weights)
```

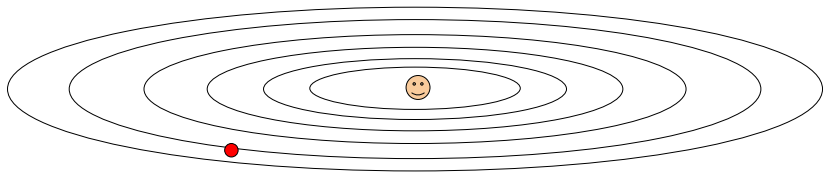
```
    weights += - step_size * weights_grad # perform parameter update
```



Optimizers

Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?
What does gradient descent do?



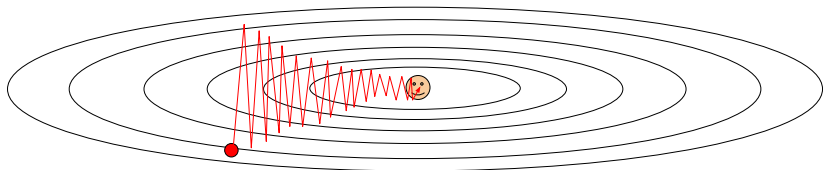
Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

Optimizers

Optimization: Problems with SGD

What if loss changes quickly in one direction and slowly in another?
What does gradient descent do?

Very slow progress along shallow dimension, jitter along steep direction

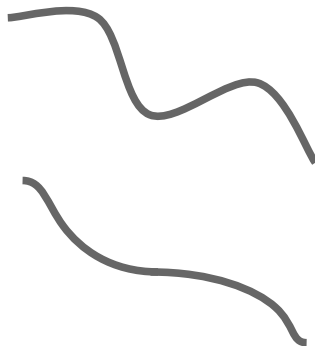


Loss function has high **condition number**: ratio of largest to smallest singular value of the Hessian matrix is large

Optimizers

Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point**?

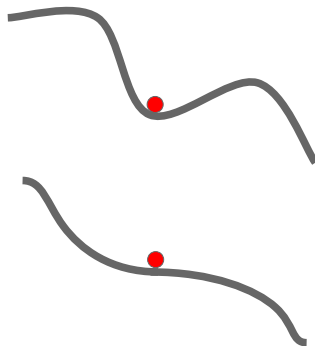


Optimizers

Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point**?

Zero gradient,
gradient descent
gets stuck

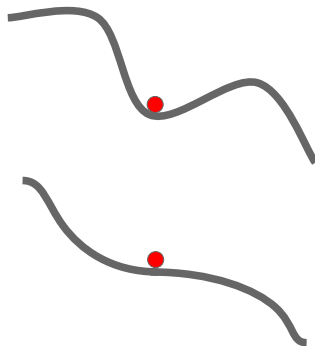


Optimizers

Optimization: Problems with SGD

What if the loss function has a **local minima** or **saddle point**?

Saddle points much more common in high dimension



Dauphin et al, "Identifying and attacking the saddle point problem in high-dimensional non-convex optimization", NIPS 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 19

April 25, 2017

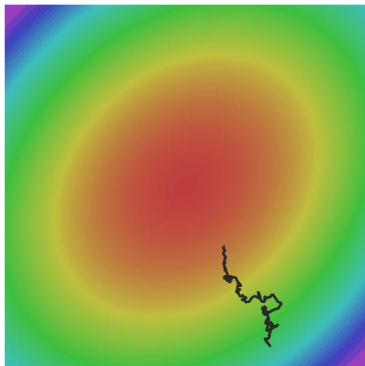
Optimizers

Optimization: Problems with SGD

Our gradients come from minibatches so they can be noisy!

$$L(W) = \frac{1}{N} \sum_{i=1}^N L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^N \nabla_W L_i(x_i, y_i, W)$$



Exponential moving average

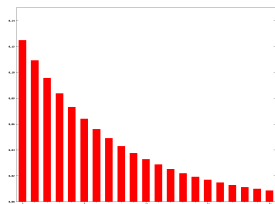
- $$S_t = \begin{cases} Y_1, & t = 1 \\ \alpha \cdot Y_t + (1 - \alpha) \cdot S_{t-1}, & t > 1 \end{cases}$$

Exponential moving average

- $S_t = \begin{cases} Y_1, & t = 1 \\ \alpha \cdot Y_t + (1 - \alpha) \cdot S_{t-1}, & t > 1 \end{cases}$
- $S_t = \alpha [Y_{t-1} + (1 - \alpha)Y_{t-2} + (1 - \alpha)^2 Y_{t-3} + \dots]$

Exponential moving average

- $S_t = \begin{cases} Y_1, & t = 1 \\ \alpha \cdot Y_t + (1 - \alpha) \cdot S_{t-1}, & t > 1 \end{cases}$
- $S_t = \alpha [Y_{t-1} + (1 - \alpha)Y_{t-2} + (1 - \alpha)^2 Y_{t-3} + \dots]$
 $= \frac{Y_{t-1} + (1 - \alpha)Y_{t-2} + (1 - \alpha)^2 Y_{t-3} + \dots}{1 + (1 - \alpha) + (1 - \alpha)^2 + \dots}$



Optimizers

Momentum update

```
# Gradient descent update  
x += - learning_rate * dx
```



```
# Momentum update  
v = mu * v - learning_rate * dx # integrate velocity  
x += v # integrate position
```

- Physical interpretation as ball rolling down the loss function + friction (μ coefficient).
- μ = usually ~ 0.5 , 0.9 , or 0.99 (Sometimes annealed over time, e.g. from $0.5 \rightarrow 0.99$)

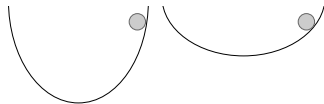
Optimizers

Momentum update

```
# Gradient descent update  
x += - learning_rate * dx
```

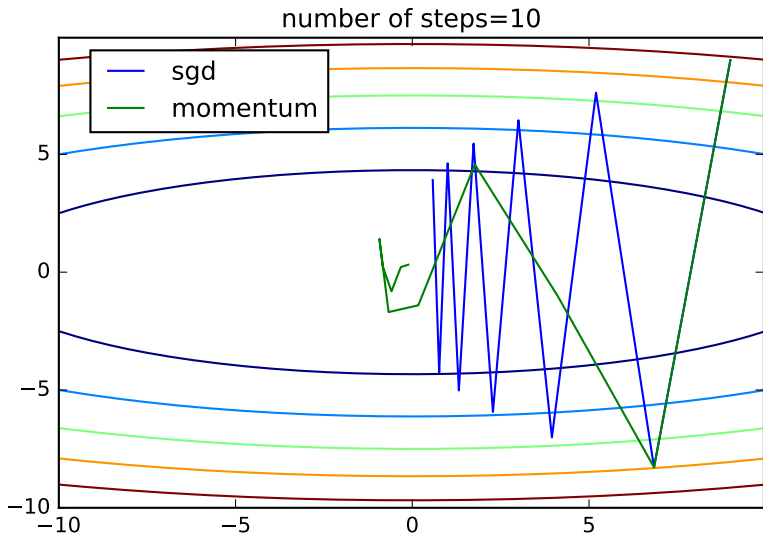


```
# Momentum update  
v = mu * v - learning_rate * dx # integrate velocity  
x += v # integrate position
```



- Allows a velocity to “build up” along shallow directions
- Velocity becomes damped in steep direction due to quickly changing sign

Optimizers

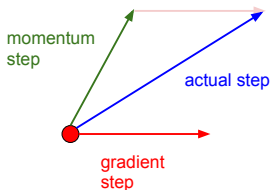


Optimizers

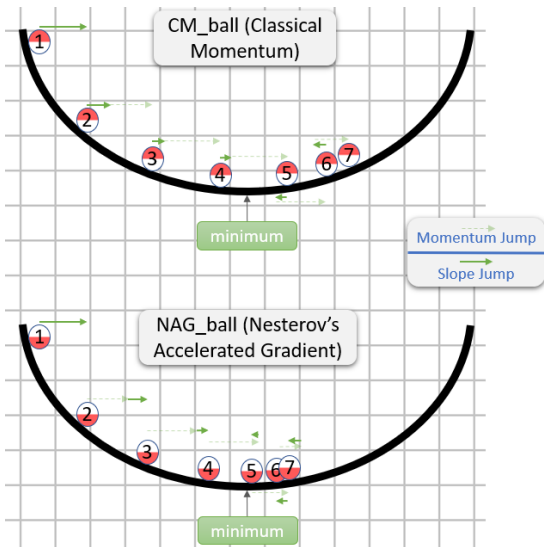
Nesterov Momentum update

```
# Momentum update
v = mu * v - learning_rate * dx # integrate velocity
x += v # integrate position
```

Ordinary momentum update:



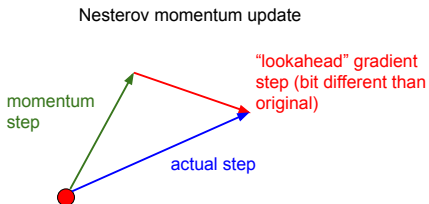
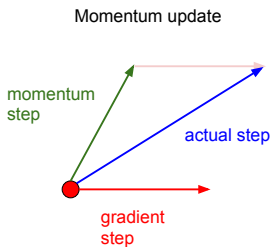
Optimizers



Reference: <https://stats.stackexchange.com/questions/179915/whats-the-difference-between-momentum-based-gradient-descent-and-nesterovs-acc>

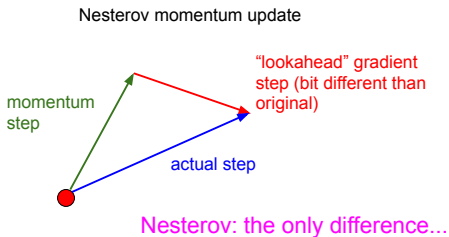
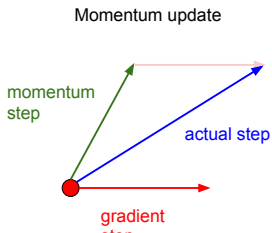
Optimizers

Nesterov Momentum update



Optimizers

Nesterov Momentum update



$$v_t = \mu v_{t-1} - \epsilon \nabla f(x_{t-1} + \mu v_{t-1})$$

$$x_t = x_{t-1} + v_t$$

We want to deal with $\nabla f(x_{t-1})$ instead

Optimizers

$$v_t = \mu v_{t-1} - \epsilon \nabla f(x_{t-1} + \mu v_{t-1})$$

$$x_t = x_{t-1} + v_t$$

Optimizers

$$v_t = \mu v_{t-1} - \epsilon \nabla f(x_{t-1} + \mu v_{t-1})$$

$$x_t = x_{t-1} + v_t$$

Pick $\tilde{x}_t = x_t + \mu v_t$,

$$v_t = \mu v_{t-1} - \epsilon \nabla(\tilde{x}_{t-1})$$

Optimizers

$$v_t = \mu v_{t-1} - \epsilon \nabla f(x_{t-1} + \mu v_{t-1})$$

$$x_t = x_{t-1} + v_t$$

Pick $\tilde{x}_t = x_t + \mu v_t$,

$$v_t = \mu v_{t-1} - \epsilon \nabla(\tilde{x}_{t-1})$$

$$\tilde{x}_t = x_t + \mu v_t = x_{t-1} + v_t + \mu v_t$$

Optimizers

$$v_t = \mu v_{t-1} - \epsilon \nabla f(x_{t-1} + \mu v_{t-1})$$

$$x_t = x_{t-1} + v_t$$

Pick $\tilde{x}_t = x_t + \mu v_t$,

$$v_t = \mu v_{t-1} - \epsilon \nabla(\tilde{x}_{t-1})$$

$$\tilde{x}_t = x_t + \mu v_t = x_{t-1} + v_t + \mu v_t$$

$$= \tilde{x}_{t-1} - \mu v_{t-1} + v_t + \mu v_t$$

Optimizers

$$v_t = \mu v_{t-1} - \epsilon \nabla f(x_{t-1} + \mu v_{t-1})$$

$$x_t = x_{t-1} + v_t$$

Pick $\tilde{x}_t = x_t + \mu v_t$,

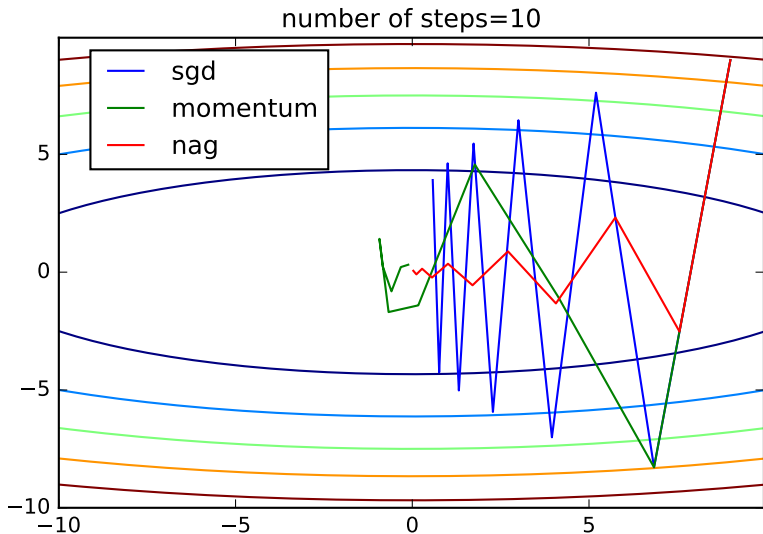
$$v_t = \mu v_{t-1} - \epsilon \nabla(\tilde{x}_{t-1})$$

$$\tilde{x}_t = x_t + \mu v_t = x_{t-1} + v_t + \mu v_t$$

$$= \tilde{x}_{t-1} - \mu v_{t-1} + v_t + \mu v_t$$

$$= \tilde{x}_{t-1} + v_t + \mu(v_t - v_{t-1})$$

Optimizers



Optimizers

AdaGrad update

[Duchi et al., 2011]

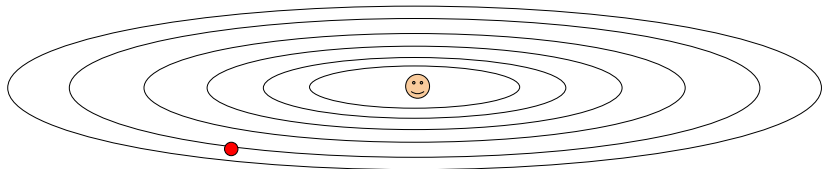
```
# Adagrad update  
cache += dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

Added element-wise scaling of the gradient based on the historical sum of squares in each dimension

Optimizers

AdaGrad update

```
# Adagrad update  
cache += dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

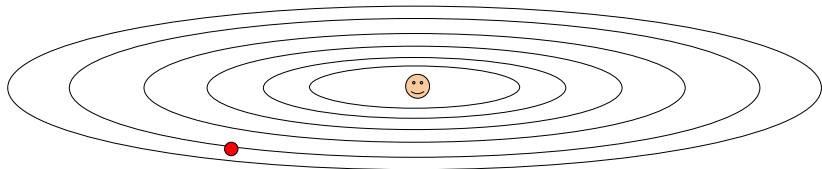


Q: What happens with AdaGrad?

Optimizers

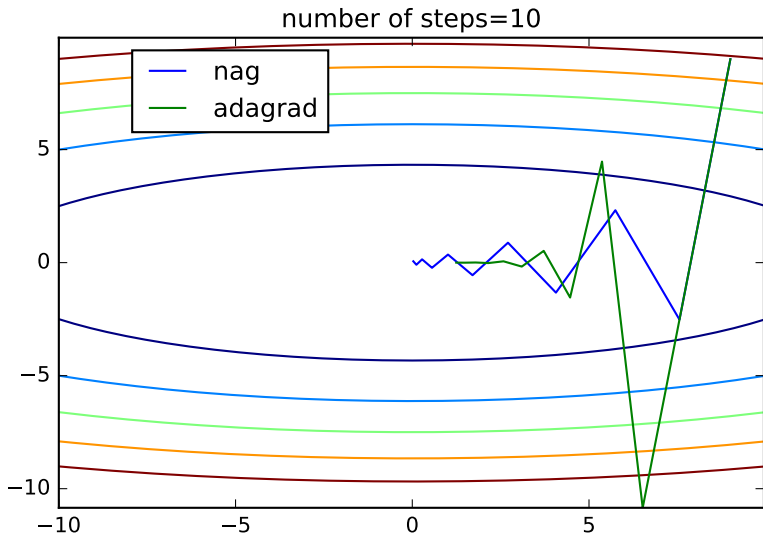
AdaGrad update

```
# Adagrad update  
cache += dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```



Q2: What happens to the step size over long time?

Optimizers



Optimizers

RMSProp update

[Tieleman and Hinton, 2012]

```
# Adagrad update  
cache += dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```



```
# RMSProp  
cache = decay_rate * cache + (1 - decay_rate) * dx**2  
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

Optimizers

rmsprop: A mini-batch version of rprop

- rprop is equivalent to using the gradient but also dividing by the size of the gradient.
 - The problem with mini-batch rprop is that we divide by a different number for each mini-batch. So why not force the number we divide by to be very similar for adjacent mini-batches?
- rmsprop: Keep a moving average of the squared gradient for each weight
$$MeanSquare(w, t) = 0.9 MeanSquare(w, t-1) + 0.1 \left(\frac{\partial E}{\partial w}(t) \right)^2$$
- Dividing the gradient by $\sqrt{MeanSquare(w, t)}$ makes the learning work much better (Tijmen Tieleman, unpublished).

Introduced in a slide in Geoff Hinton's Coursera class, lecture 6

Optimizers

rmsprop: A mini-batch version of rprop

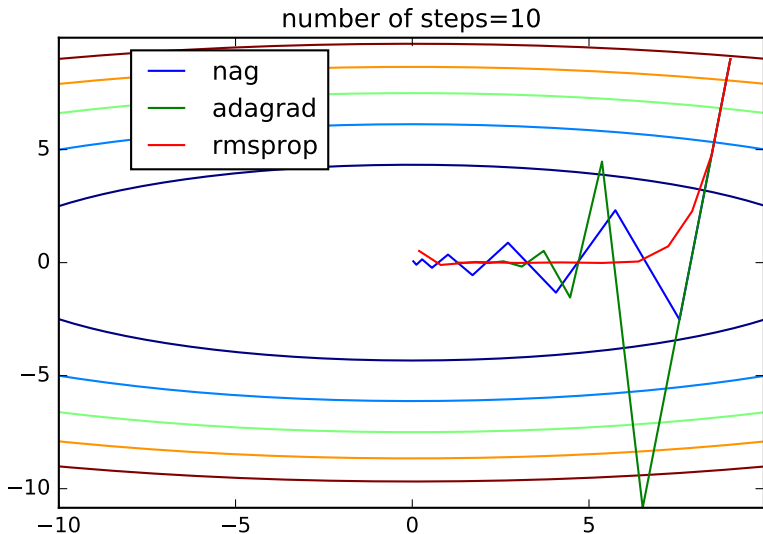
- rprop is equivalent to using the gradient but also dividing by the size of the gradient.
 - The problem with mini-batch rprop is that we divide by a different number for each mini-batch. So why not force the number we divide by to be very similar for adjacent mini-batches?
- rmsprop: Keep a moving average of the squared gradient for each weight
$$MeanSquare(w, t) = 0.9 MeanSquare(w, t-1) + 0.1 \left(\frac{\partial E}{\partial w}(t) \right)^2$$
- Dividing the gradient by $\sqrt{MeanSquare(w, t)}$ makes the learning work much better (Tijmen Tieleman, unpublished).

Introduced in a slide in Geoff Hinton's Coursera class, lecture 6

Cited by several papers as:

[52] T. Tieleman and G. E. Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude., 2012.

Optimizers



Optimizers

Adam update

[Kingma and Ba, 2014]

(incomplete, but close)

```
# Adam
m = beta1*m + (1-beta1)*dx # update first moment
v = beta2*v + (1-beta2)*(dx**2) # update second moment
x += - learning_rate * m / (np.sqrt(v) + 1e-7)
```

Optimizers

Adam update

[Kingma and Ba, 2014]

(incomplete, but close)

```
# Adam
m = beta1*m + (1-beta1)*dx # update first moment
v = beta2*v + (1-beta2)*(dx**2) # update second moment
x += - learning_rate * m / (np.sqrt(v) + 1e-7)
```

momentum

RMSProp-like

Looks a bit like RMSProp with momentum

Optimizers

Adam update

[Kingma and Ba, 2014]

(incomplete, but close)

```
# Adam
m = beta1*m + (1-beta1)*dx # update first moment
v = beta2*v + (1-beta2)*(dx**2) # update second moment
x += - learning_rate * m / (np.sqrt(v) + 1e-7)
```

momentum

RMSProp-like

Looks a bit like RMSProp with momentum

```
# RMSProp
cache = decay_rate * cache + (1 - decay_rate) * dx**2
x += - learning_rate * dx / (np.sqrt(cache) + 1e-7)
```

Optimizers

Adam (full form)

```

first_moment = 0
second_moment = 0
for t in range(1, num_iterations):
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    first_unbias = first_moment / (1 - beta1 ** t)
    second_unbias = second_moment / (1 - beta2 ** t)
    x -= learning_rate * first_unbias / (np.sqrt(second_unbias) + 1e-7)

```

Momentum

Bias correction

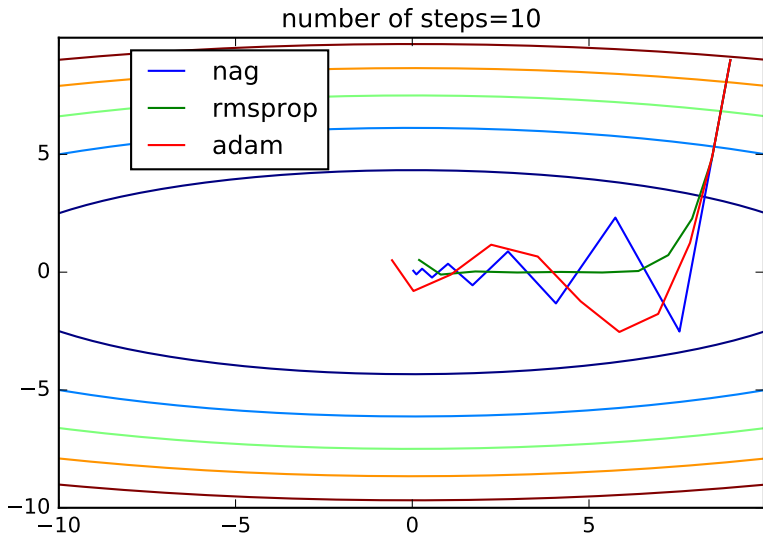
AdaGrad / RMSProp

Bias correction for the fact that first and second moment estimates start at zero

Adam with $\beta_1 = 0.9$, $\beta_2 = 0.999$, and $\text{learning_rate} = 1e-3$ or $5e-4$ is a great starting point for many models!

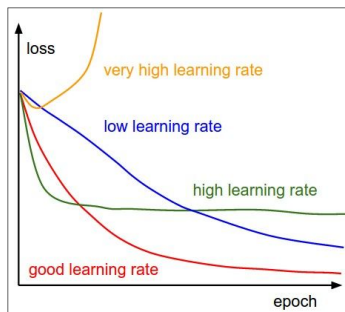
Kingma and Ba, "Adam: A method for stochastic optimization", ICLR 2015

Optimizers



Optimizers

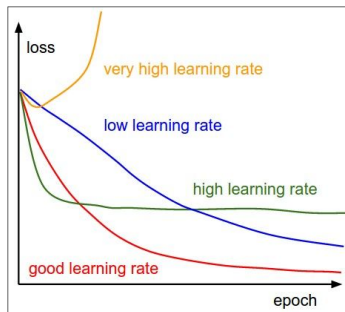
SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



Q: Which one of these learning rates is best to use?

Optimizers

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



=> **Learning rate decay over time!**

step decay:

e.g. decay learning rate by half every few epochs.

exponential decay:

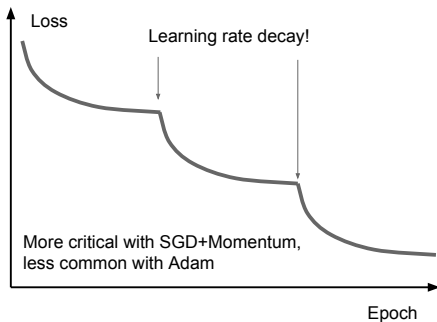
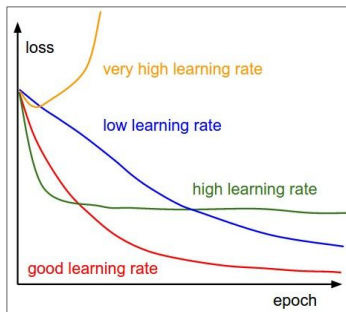
$$\alpha = \alpha_0 e^{-kt}$$

1/t decay:

$$\alpha = \alpha_0 / (1 + kt)$$

Optimizers

SGD, SGD+Momentum, Adagrad, RMSProp, Adam all have **learning rate** as a hyperparameter.



Optimizers

Second order optimization methods

second-order Taylor expansion:

$$J(\boldsymbol{\theta}) \approx J(\boldsymbol{\theta}_0) + (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)^\top \mathbf{H} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$$

Solving for the critical point we obtain the Newton parameter update:

$$\boldsymbol{\theta}^* = \boldsymbol{\theta}_0 - \mathbf{H}^{-1} \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}_0)$$

Q: what is nice about this update?

Optimizers

Second order optimization methods

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Inverting Hessian is very expensive ($O(N^3)$). Avoiding that resulting in so-called Quasi-Newton methods

Optimizers

Second order optimization methods

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Inverting Hessian is very expensive ($O(N^3)$). Avoiding that resulting in so-called Quasi-Newton methods

- Rank-1 inverse Hessian update (simple but not too commonly used)

Optimizers

Second order optimization methods

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Inverting Hessian is very expensive ($O(N^3)$). Avoiding that resulting in so-called Quasi-Newton methods

- Rank-1 inverse Hessian update (simple but not too commonly used)
- Rank-2 inverse Hessian update

Optimizers

Second order optimization methods

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Inverting Hessian is very expensive ($O(N^3)$). Avoiding that resulting in so-called Quasi-Newton methods

- Rank-1 inverse Hessian update (simple but not too commonly used)
- Rank-2 inverse Hessian update
 - BFGS (most popular) and DFS

Optimizers

Second order optimization methods

$$\theta^* = \theta_0 - H^{-1} \nabla_{\theta} J(\theta_0)$$

Inverting Hessian is very expensive ($O(N^3)$). Avoiding that resulting in so-called Quasi-Newton methods

- Rank-1 inverse Hessian update (simple but not too commonly used)
- Rank-2 inverse Hessian update
 - BFGS (most popular) and DFS
 - LBFGS
 - Does not store the entire inverse Hessian
 - Tradeoff space with time and accuracy

Quasi-Newton methods (watch this)

- Ref:
 - ① https://www.youtube.com/watch?v=uo2z0AT_83k
 - ② Nocedal & Wright - Numerical Optimization ($B \leftrightarrow H$)
 - ③ http://users.ece.utexas.edu/~cmcaram/EE381V_2012F/Lecture_10_Scribe_Notes.final.pdf
- The inverse of Hessian H is expensive to compute. Want to approximate it iteratively instead

Quasi-Newton methods (watch this)

- Ref:
 - ① https://www.youtube.com/watch?v=uo2z0AT_83k
 - ② Nocedal & Wright - Numerical Optimization ($B \leftrightarrow H$)
 - ③ http://users.ece.utexas.edu/~cmcaram/EE381V_2012F/Lecture_10_Scribe_Notes.final.pdf
- The inverse of Hessian H is expensive to compute. Want to approximate it iteratively instead
- Quasi-Newton methods:
 - ① Approximate Newton direction

$$d_k = -B_k g_k,$$

where $B_k \approx H^{-1}$ and $g_k = \nabla J(\theta_k)$

Quasi-Newton methods (watch this)

- Ref:
 - ① https://www.youtube.com/watch?v=uo2z0AT_83k
 - ② Nocedal & Wright - Numerical Optimization ($B \leftrightarrow H$)
 - ③ http://users.ece.utexas.edu/~cmcaram/EE381V_2012F/Lecture_10_Scribe_Notes.final.pdf
- The inverse of Hessian H is expensive to compute. Want to approximate it iteratively instead
- Quasi-Newton methods:
 - ① Approximate Newton direction

$$d_k = -B_k g_k,$$

where $B_k \approx H^{-1}$ and $g_k = \nabla J(\theta_k)$

- ② Line search: $\theta_{k+1} = \theta_k + \alpha_k d_k$

Quasi-Newton methods (watch this)

- Ref:
 - ① https://www.youtube.com/watch?v=uo2z0AT_83k
 - ② Nocedal & Wright - Numerical Optimization ($B \leftrightarrow H$)
 - ③ http://users.ece.utexas.edu/~cmcaram/EE381V_2012F/Lecture_10_Scribe_Notes.final.pdf
- The inverse of Hessian H is expensive to compute. Want to approximate it iteratively instead
- Quasi-Newton methods:
 - ① Approximate Newton direction

$$d_k = -B_k g_k,$$

where $B_k \approx H^{-1}$ and $g_k = \nabla J(\theta_k)$

- ② Line search: $\theta_{k+1} = \theta_k + \alpha_k d_k$
- ③ Update $g_{k+1} = \nabla J(\theta_{k+1})$

Quasi-Newton methods (watch this)

- Ref:
 - ① https://www.youtube.com/watch?v=uo2z0AT_83k
 - ② Nocedal & Wright - Numerical Optimization ($B \leftrightarrow H$)
 - ③ http://users.ece.utexas.edu/~cmcaram/EE381V_2012F/Lecture_10_Scribe_Notes.final.pdf
- The inverse of Hessian H is expensive to compute. Want to approximate it iteratively instead
- Quasi-Newton methods:
 - ① Approximate Newton direction

$$d_k = -B_k g_k,$$

where $B_k \approx H^{-1}$ and $g_k = \nabla J(\theta_k)$

- ② Line search: $\theta_{k+1} = \theta_k + \alpha_k d_k$
- ③ Update $g_{k+1} = \nabla J(\theta_{k+1})$
- ④ Approximate inverse Hessian

$$B_{k+1} = \text{update_formula}(B_k, \theta_{k+1} - \theta_k, g_{k+1} - g_k)$$

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$
- We may assume the above is satisfied and use this to iteratively approximate H .

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$
- We may assume the above is satisfied and use this to iteratively approximate H . That is (known as secant equation) $H p_k = q_k$, where $p_k = \theta_{k+1} - \theta_k$ and $q_k = \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$
- We may assume the above is satisfied and use this to iteratively approximate H . That is (known as secant equation) $H p_k = q_k$, where $p_k = \theta_{k+1} - \theta_k$ and $q_k = \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$
- Let $H_{k+1} = H_k + uv^T$

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$
- We may assume the above is satisfied and use this to iteratively approximate H . That is (known as secant equation) $H p_k = q_k$, where $p_k = \theta_{k+1} - \theta_k$ and $q_k = \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$
- Let $H_{k+1} = H_k + uv^T \Rightarrow (H_k + uv^T)p_k = q_k$

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$
- We may assume the above is satisfied and use this to iteratively approximate H . That is (known as secant equation) $H p_k = q_k$, where $p_k = \theta_{k+1} - \theta_k$ and $q_k = \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$
- Let $H_{k+1} = H_k + uv^T \Rightarrow (H_k + uv^T)p_k = q_k$
 $\Rightarrow u(v^T p_k) = q_k - H_k p_k$

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$
- We may assume the above is satisfied and use this to iteratively approximate H . That is (known as secant equation) $H p_k = q_k$, where $p_k = \theta_{k+1} - \theta_k$ and $q_k = \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$
- Let $H_{k+1} = H_k + uv^T \Rightarrow (H_k + uv^T)p_k = q_k$
 $\Rightarrow u(v^T p_k) = q_k - H_k p_k \Rightarrow u = \frac{1}{v^T p_k}(q_k - H_k p_k)$

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have
$$\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$$
- We may assume the above is satisfied and use this to iteratively approximate H . That is (known as secant equation) $H p_k = q_k$, where $p_k = \theta_{k+1} - \theta_k$ and $q_k = \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$
- Let $H_{k+1} = H_k + uv^T \Rightarrow (H_k + uv^T)p_k = q_k$
 $\Rightarrow u(v^T p_k) = q_k - H_k p_k \Rightarrow u = \frac{1}{v^T p_k}(q_k - H_k p_k)$
- We are free to pick v . But since we know H has to be symmetric, let's pick $v = q_k - H_k p_k$.

Approximation with rank-1 update

- As Hessian is essentially the “derivative” of ∇J , we have $\nabla J(\theta_{k+1}) \approx \nabla J(\theta_k) + H(\theta_{k+1} - \theta_k)$
- We may assume the above is satisfied and use this to iteratively approximate H . That is (known as secant equation) $H p_k = q_k$, where $p_k = \theta_{k+1} - \theta_k$ and $q_k = \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$
- Let $H_{k+1} = H_k + uv^T \Rightarrow (H_k + uv^T)p_k = q_k$
 $\Rightarrow u(v^T p_k) = q_k - H_k p_k \Rightarrow u = \frac{1}{v^T p_k} (q_k - H_k p_k)$
- We are free to pick v . But since we know H has to be symmetric, let's pick $v = q_k - H_k p_k$. Thus

$$H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T$$

with $v = q_k - H_k p_k$

Updating B

- Recall that we need $B_k = H_k^{-1}$ to approximate the Newton direction ($d_k = -B_k g_k$)

Updating B

- Recall that we need $B_k = H_k^{-1}$ to approximate the Newton direction ($d_k = -B_k g_k$)
- We don't need to invert the matrix H_k directly. Note that $H p_k = q_k$ give us $H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T$ and $v = q_k - H_k p_k$

Updating B

- Recall that we need $B_k = H_k^{-1}$ to approximate the Newton direction ($d_k = -B_k g_k$)
- We don't need to invert the matrix H_k directly. Note that $H p_k = q_k$ give us $H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T$ and $v = q_k - H_k p_k$
- Similarly, given $B q_k = p_k$, we have

$$B_{k+1} = B_k + \frac{1}{w^T q_k} w w^T$$

with $w = p_k - B_k q_k$

Rank-2 approximation

- BFGS utilizes rank-2 approximation update for H . There are other variations (such as DFP). But BFGS is considered the state of the art

Rank-2 approximation

- BFGS utilizes rank-2 approximation update for H . There are other variations (such as DFP). But BFGS is considered the state of the art
- Recall our rank-1 approximation that
$$H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T \text{ and } v = q_k - H_k p_k$$

Rank-2 approximation

- BFGS utilizes rank-2 approximation update for H . There are other variations (such as DFP). But BFGS is considered the state of the art
- Recall our rank-1 approximation that
$$H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T \text{ and } v = q_k - H_k p_k$$
- Consider update $H_{k+1} = H_k + \frac{1}{\alpha} u u^T + \frac{1}{\beta} w w^T$ instead.

Rank-2 approximation

- BFGS utilizes rank-2 approximation update for H . There are other variations (such as DFP). But BFGS is considered the state of the art
- Recall our rank-1 approximation that
$$H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T \text{ and } v = q_k - H_k p_k$$
- Consider update $H_{k+1} = H_k + \frac{1}{\alpha} u u^T + \frac{1}{\beta} w w^T$ instead.
 - Need to pick u and w , q_k and $H_k p_k$ are reasonable choice

Rank-2 approximation

- BFGS utilizes rank-2 approximation update for H . There are other variations (such as DFP). But BFGS is considered the state of the art
- Recall our rank-1 approximation that $H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T$ and $v = q_k - H_k p_k$
- Consider update $H_{k+1} = H_k + \frac{1}{\alpha} u u^T + \frac{1}{\beta} w w^T$ instead.
 - Need to pick u and w , q_k and $H_k p_k$ are reasonable choice
- Again, we want $H_{k+1} p_k = q_k$
 $\Rightarrow H_k p_k + \frac{1}{\alpha} q_k q_k^T p_k + \frac{1}{\beta} H_k p_k p_k^T H_k^T p_k = q_k.$

Rank-2 approximation

- BFGS utilizes rank-2 approximation update for H . There are other variations (such as DFP). But BFGS is considered the state of the art
- Recall our rank-1 approximation that $H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T$ and $v = q_k - H_k p_k$
- Consider update $H_{k+1} = H_k + \frac{1}{\alpha} u u^T + \frac{1}{\beta} w w^T$ instead.
 - Need to pick u and w , q_k and $H_k p_k$ are reasonable choice
- Again, we want $H_{k+1} p_k = q_k$
 $\Rightarrow H_k p_k + \frac{1}{\alpha} q_k q_k^T p_k + \frac{1}{\beta} H_k p_k p_k^T H_k^T p_k = q_k$. By inspection, this can be satisfied if we pick $\alpha = q_k^T p_k$ and $\beta = -p_k^T H_k^T p_k$.

Rank-2 approximation

- BFGS utilizes rank-2 approximation update for H . There are other variations (such as DFP). But BFGS is considered the state of the art
- Recall our rank-1 approximation that $H_{k+1} = H_k + \frac{1}{v^T p_k} v v^T$ and $v = q_k - H_k p_k$
- Consider update $H_{k+1} = H_k + \frac{1}{\alpha} u u^T + \frac{1}{\beta} w w^T$ instead.
 - Need to pick u and w , q_k and $H_k p_k$ are reasonable choice
- Again, we want $H_{k+1} p_k = q_k$
 $\Rightarrow H_k p_k + \frac{1}{\alpha} q_k q_k^T p_k + \frac{1}{\beta} H_k p_k p_k^T H_k^T p_k = q_k$. By inspection, this can be satisfied if we pick $\alpha = q_k^T p_k$ and $\beta = -p_k^T H_k^T p_k$. Thus we have

$$H_{k+1} = H_k + \frac{q_k q_k^T}{q_k^T p_k} - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}$$

Sherman-Morrison-formula

- But we are interested in $B_k = H_k^{-1}$
- Sherman-Morrison-formula:

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

Proof.

$$(A + uv^T) \left(A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right)$$

Sherman-Morrison-formula

- But we are interested in $B_k = H_k^{-1}$
- Sherman-Morrison-formula:

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

Proof.

$$\begin{aligned} & (A + uv^T) \left(A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right) \\ &= AA^{-1} + uv^T A^{-1} \end{aligned}$$

Sherman-Morrison-formula

- But we are interested in $B_k = H_k^{-1}$
- Sherman-Morrison-formula:

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

Proof.

$$\begin{aligned} & (A + uv^T) \left(A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right) \\ &= AA^{-1} + uv^T A^{-1} - \frac{AA^{-1}uv^T A^{-1} + uv^T A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \end{aligned}$$

Sherman-Morrison-formula

- But we are interested in $B_k = H_k^{-1}$
- Sherman-Morrison-formula:

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

Proof.

$$\begin{aligned} & (A + uv^T) \left(A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right) \\ &= AA^{-1} + uv^T A^{-1} - \frac{AA^{-1}uv^T A^{-1} + uv^T A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \\ &= I + uv^T A^{-1} - \frac{uv^T A^{-1} + uv^T A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \end{aligned}$$

Sherman-Morrison-formula

- But we are interested in $B_k = H_k^{-1}$
- Sherman-Morrison-formula:

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

Proof.

$$\begin{aligned} & (A + uv^T) \left(A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right) \\ &= AA^{-1} + uv^T A^{-1} - \frac{AA^{-1}uv^T A^{-1} + uv^T A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \\ &= I + uv^T A^{-1} - \frac{uv^T A^{-1} + uv^T A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \\ &= I + uv^T A^{-1} - \frac{u(1 + v^T A^{-1}u)v^T A^{-1}}{1 + v^T A^{-1}u} \end{aligned}$$

Sherman-Morrison-formula

- But we are interested in $B_k = H_k^{-1}$
- Sherman-Morrison-formula:

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

Proof.

$$\begin{aligned} & (A + uv^T) \left(A^{-1} - \frac{A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \right) \\ &= AA^{-1} + uv^T A^{-1} - \frac{AA^{-1}uv^T A^{-1} + uv^T A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \\ &= I + uv^T A^{-1} - \frac{uv^T A^{-1} + uv^T A^{-1}uv^T A^{-1}}{1 + v^T A^{-1}u} \\ &= I + uv^T A^{-1} - \frac{u(1 + v^T A^{-1}u)v^T A^{-1}}{1 + v^T A^{-1}u} = I + uv^T A^{-1} - uv^T A^{-1} = I \quad \square \end{aligned}$$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k}}_D - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k}}_D - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k} - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}}_D$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T H D^{-1}}{p^T H^T p(1 - p^T H D^{-1} H p/(p^T H^T p))}$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k} - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}}_D$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T H^T p(1 - p^T HD^{-1}Hp/(p^T H^T p))}$
 $= D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T Hp - p^T HD^{-1}Hp}$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k}}_D - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T H^T p(1 - p^T HD^{-1}Hp/(p^T H^T p))}$
 $= D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T Hp - p^T HD^{-1}Hp}$

- $D^{-1}Hp = (BHp + \frac{Bqq^T BHp}{q^T p - q^T Bq}) = (p + \frac{Bqq^T p}{q^T p - q^T Bq})$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k}}_D - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T H^T p(1 - p^T HD^{-1}Hp/(p^T H^T p))}$
 $= D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T Hp - p^T HD^{-1}Hp}$

- $D^{-1}Hp = (BHp + \frac{Bqq^T BHp}{q^T p - q^T Bq}) = (p + \frac{Bqq^T p}{q^T p - q^T Bq})$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T qq^T p(q^T p - q^T Bq)}$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k}}_D - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T H^T p(1 - p^T HD^{-1}Hp/(p^T H^T p))}$
 $= D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T Hp - p^T HD^{-1}Hp}$

- $D^{-1}Hp = (BHp + \frac{Bqq^T BHp}{q^T p - q^T Bq}) = (p + \frac{Bqq^T p}{q^T p - q^T Bq})$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T qq^T p(q^T p - q^T Bq)} \dots$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k}}_D - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T H^T p(1 - p^T HD^{-1}Hp/(p^T H^T p))}$
 $= D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T Hp - p^T HD^{-1}Hp}$

- $D^{-1}Hp = (BHp + \frac{Bqq^T BHp}{q^T p - q^T Bq}) = (p + \frac{Bqq^T p}{q^T p - q^T Bq})$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T qq^T p(q^T p - q^T Bq)} \dots$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = \left(I - \frac{pq^T}{q^T p}\right) B \left(I - \frac{qp^T}{q^T p}\right) + \frac{pp^T}{q^T p}$
 $\Rightarrow B_{k+1} = \left(I - \frac{p_k q_k^T}{q_k^T p_k}\right) B_k \left(I - \frac{q_k p_k^T}{q_k^T p_k}\right) + \frac{p_k p_k^T}{q_k^T p_k}$

Inverse Hessian update for BFGS

- Recall $H_{k+1} = H_k + \underbrace{\frac{q_k q_k^T}{q_k^T p_k} - \frac{H_k p_k p_k^T H_k}{p_k^T H_k^T p_k}}_D$ and

$$(A + uv^T)^{-1} = A^{-1} + \frac{A^{-1}uv^T A^{-1}}{1 - v^T A^{-1}u}$$

- $D^{-1} = (H + \frac{qq^T}{q^T p})^{-1} = H^{-1} + \frac{H^{-1}qq^T H^{-1}}{(q^T p)(1 - q^T H^{-1}q/(q^T p))} = B + \frac{Bqq^T B}{q^T p - q^T Bq}$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T H^T p(1 - p^T HD^{-1}Hp/(p^T H^T p))}$
 $= D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T Hp - p^T HD^{-1}Hp}$

- $D^{-1}Hp = (BHp + \frac{Bqq^T BHp}{q^T p - q^T Bq}) = (p + \frac{Bqq^T p}{q^T p - q^T Bq})$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = D^{-1} - \frac{D^{-1}Hpp^T HD^{-1}}{p^T qq^T p(q^T p - q^T Bq)} \dots$

- $(D - \frac{Hpp^T H}{p^T H^T p})^{-1} = \left(I - \frac{pq^T}{q^T p}\right) B \left(I - \frac{qp^T}{q^T p}\right) + \frac{pp^T}{q^T p}$
 $\Rightarrow B_{k+1} = \left(I - \frac{p_k q_k^T}{q_k^T p_k}\right) B_k \left(I - \frac{q_k p_k^T}{q_k^T p_k}\right) + \frac{p_k p_k^T}{q_k^T p_k}$

- Bounty: 3% bonus to complete the algebra

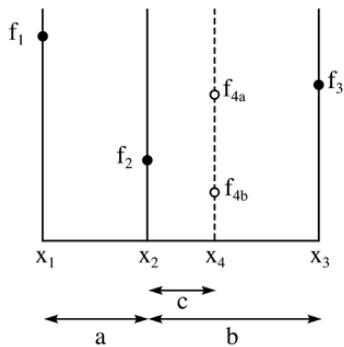
Summary of BFGS

Initialize Initialize inverse Hessian approximation $B \leftarrow B_0$. Can set $B \leftarrow I$ if no initial estimate; $k \leftarrow 0$; Pick a random starting point θ_0

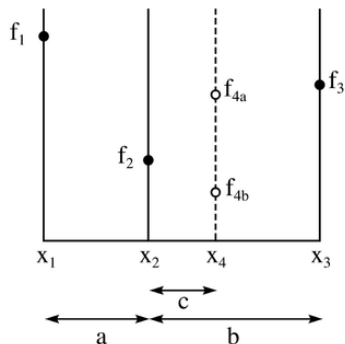
- Loop**
- ① Get search direction $d_k = -B_k \nabla J(\theta_k)$
 - ② Conduct line search to find optimum
 $\theta_{k+1} = \theta_k + \alpha_k d_k$
 - ③ $p_k \leftarrow \theta_{k+1} - \theta_k$; $q_k \leftarrow \nabla J(\theta_{k+1}) - \nabla J(\theta_k)$;

$$B_{k+1} = \left(I - \frac{p_k q_k^T}{q_k^T p_k} \right) B_k \left(I - \frac{q_k p_k^T}{q_k^T p_k} \right) + \frac{p_k p_k^T}{q_k^T p_k}$$
 - ④ $k \leftarrow k + 1$; Exit if $\|\nabla J(\theta_k)\| < \epsilon$

Golden-section search

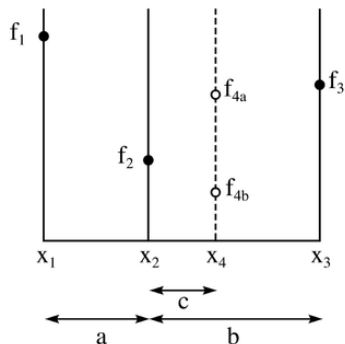


Golden-section search



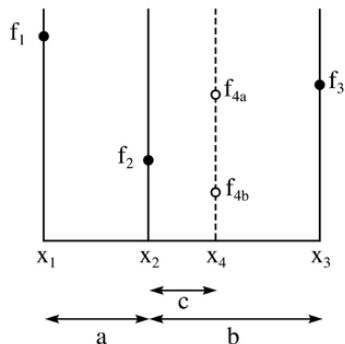
- If we have f_{4a} , minimum is in $[x_1, x_4]$

Golden-section search



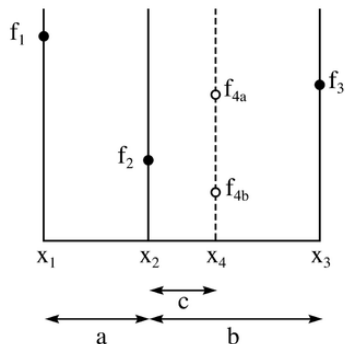
- If we have f_{4a} , minimum is in $[x_1, x_4]$
- If we have f_{4b} , minimum is in $[x_2, x_3]$

Golden-section search



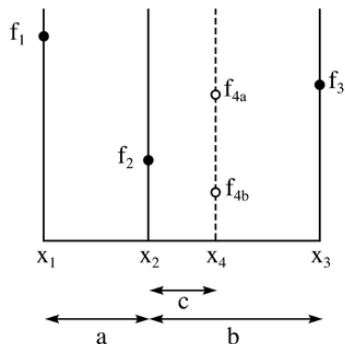
- If we have f_{4a} , minimum is in $[x_1, x_4]$
- If we have f_{4b} , minimum is in $[x_2, x_3]$
- To maximize expected search speed, set $x_4 - x_1 = x_3 - x_2 \Rightarrow a + c = b$

Golden-section search



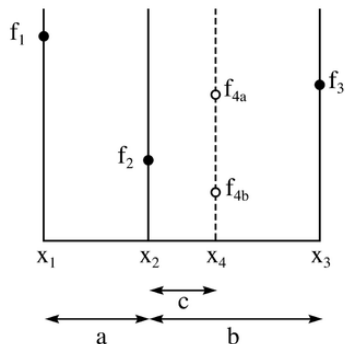
- If we have f_{4a} , minimum is in $[x_1, x_4]$
- If we have f_{4b} , minimum is in $[x_2, x_3]$
- To maximize expected search speed, set $x_4 - x_1 = x_3 - x_2 \Rightarrow a + c = b$
 - Given x_1, x_2, x_3 , we know how to pick x_4

Golden-section search



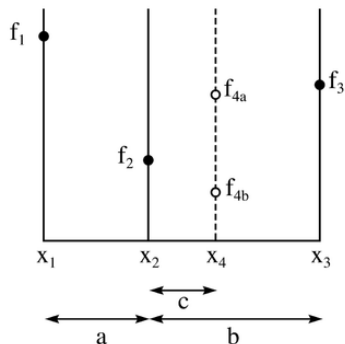
- If we have f_{4a} , minimum is in $[x_1, x_4]$
- If we have f_{4b} , minimum is in $[x_2, x_3]$
- To maximize expected search speed, set $x_4 - x_1 = x_3 - x_2 \Rightarrow a + c = b$
 - Given x_1, x_2, x_3 , we know how to pick x_4
 - How to pick x_2 given x_1 and x_3 ?

Golden-section search



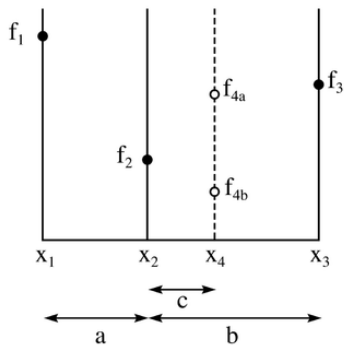
- If we have f_{4a} , minimum is in $[x_1, x_4]$
- If we have f_{4b} , minimum is in $[x_2, x_3]$
- To maximize expected search speed, set $x_4 - x_1 = x_3 - x_2 \Rightarrow a + c = b$
 - Given x_1, x_2, x_3 , we know how to pick x_4
 - How to pick x_2 given x_1 and x_3 ?
- Golden-section search simply assume the “spacing” of each iteration is proportional

Golden-section search



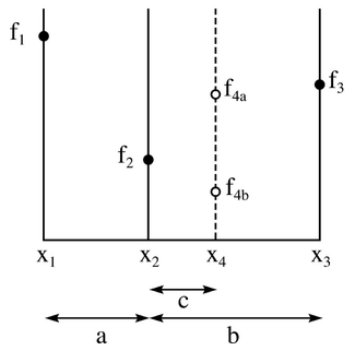
- If we have f_{4a} , minimum is in $[x_1, x_4]$
- If we have f_{4b} , minimum is in $[x_2, x_3]$
- To maximize expected search speed, set $x_4 - x_1 = x_3 - x_2 \Rightarrow a + c = b$
 - Given x_1, x_2, x_3 , we know how to pick x_4
 - How to pick x_2 given x_1 and x_3 ?
- Golden-section search simply assume the “spacing” of each iteration is proportional
 - That is, $\frac{c}{a} = \frac{a}{b}$

Golden-section search



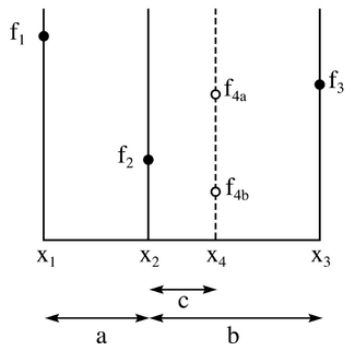
- $a + c = b$ and $\frac{c}{a} = \frac{a}{b}$

Golden-section search



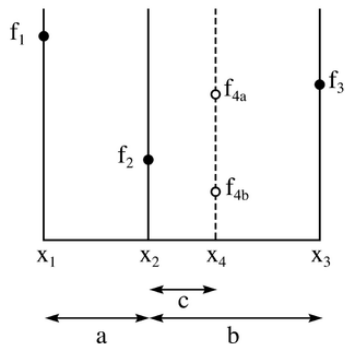
- $a + c = b$ and $\frac{c}{a} = \frac{a}{b}$
 $\Rightarrow \frac{b-a}{a} = \frac{a}{b}$

Golden-section search



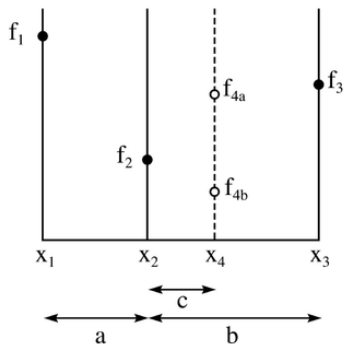
- $a + c = b$ and $\frac{c}{a} = \frac{a}{b}$
 - $\Rightarrow \frac{b-a}{a} = \frac{a}{b}$
 - $\Rightarrow \frac{b}{a} - 1 = \frac{1}{b/a}$

Golden-section search



- $a + c = b$ and $\frac{c}{a} = \frac{a}{b}$
 - $\Rightarrow \frac{b-a}{a} = \frac{a}{b}$
 - $\Rightarrow \frac{b}{a} - 1 = \frac{1}{b/a}$
 - $\Rightarrow \left(\frac{b}{a}\right)^2 - \frac{b}{a} - 1 = 0$

Golden-section search



- $a + c = b$ and $\frac{c}{a} = \frac{a}{b}$
 - $\Rightarrow \frac{b-a}{a} = \frac{a}{b}$
 - $\Rightarrow \frac{b}{a} - 1 = \frac{1}{b/a}$
 - $\Rightarrow \left(\frac{b}{a}\right)^2 - \frac{b}{a} - 1 = 0$

$$\frac{b}{a} = \frac{1 + \sqrt{5}}{2} = 1.618034\dots \triangleq \varphi$$

↑
golden ratio

Inverse Hessian update for BFGS

- Like rank-1 update, we can also rearrange the variables to obtain an update rule for $B = H^{-1}$
- Instead of $H_{k+1}p_k = q_k$, we want $B_{k+1}q_k = p_k$.

Inverse Hessian update for BFGS

- Like rank-1 update, we can also rearrange the variables to obtain an update rule for $B = H^{-1}$
- Instead of $H_{k+1}p_k = q_k$, we want $B_{k+1}q_k = p_k$. Thus we have

$$B_{k+1} = B_k + \frac{p_k p_k^T}{p_k^T q_k} - \frac{B_k q_k q_k^T B_k}{q_k^T B_k^T q_k}$$

- Note that this update rule of B is different from before. Actually this is the update rule of DFP. An older approach that is considered worse compared with BFGS

Some theoretical notes

- A prettier but more technical explanation of BFGS/DFP involves weighted matrix norm

Some theoretical notes

- A prettier but more technical explanation of BFGS/DFP involves weighted matrix norm
- Comparing with rank-1 update, we have more degree of freedom and thus can impose more requirement. Besides
 - ① $B_{k+1}q_k = p_k$ (secant equation)
 - ② $B_{k+1} \succ 0$ (symmetric and positive definite),we also require each update to be small.

Some theoretical notes

- A prettier but more technical explanation of BFGS/DFP involves weighted matrix norm
- Comparing with rank-1 update, we have more degree of freedom and thus can impose more requirement. Besides
 - ① $B_{k+1}q_k = p_k$ (secant equation)
 - ② $B_{k+1} \succ 0$ (symmetric and positive definite),

we also require each update to be small. Namely,

$$\|B_{k+1} - B_k\|_W \rightarrow \min,$$

where $\|A\|_W = \|W^{1/2}AW^{1/2}\|_F$ is the weighted Frobenius norm

Some theoretical notes

- A prettier but more technical explanation of BFGS/DFP involves weighted matrix norm
- Comparing with rank-1 update, we have more degree of freedom and thus can impose more requirement. Besides
 - ① $B_{k+1}q_k = p_k$ (secant equation)
 - ② $B_{k+1} \succ 0$ (symmetric and positive definite),

we also require each update to be small. Namely,

$$\|B_{k+1} - B_k\|_W \rightarrow \min,$$

where $\|A\|_W = \|W^{1/2}AW^{1/2}\|_F$ is the weighted Frobenius norm

- $\Rightarrow \begin{cases} \text{BFGS} & W = H \\ \text{DFP} & W = H^{-1} \end{cases}$

LBFGS

- BFGS requires us to store the complete estimate of the Hessian or inverse Hessian

LBFGS

- BFGS requires us to store the complete estimate of the Hessian or inverse Hessian
- The matrix is too big to be stored in deep learning setting (millions of variables)

LBFGS

- BFGS requires us to store the complete estimate of the Hessian or inverse Hessian
- The matrix is too big to be stored in deep learning setting (millions of variables)
- Recall that $B_{k+1} = \left(I - \frac{p_k q_k^T}{q_k^T p_k} \right) B_k \left(I - \frac{q_k p_k^T}{q_k^T p_k} \right) + \frac{p_k p_k^T}{q_k^T p_k}$, size of p_k and q_k are much smaller

LBFGS

- BFGS requires us to store the complete estimate of the Hessian or inverse Hessian
- The matrix is too big to be stored in deep learning setting (millions of variables)
- Recall that $B_{k+1} = \left(I - \frac{p_k q_k^T}{q_k^T p_k} \right) B_k \left(I - \frac{q_k p_k^T}{q_k^T p_k} \right) + \frac{p_k p_k^T}{q_k^T p_k}$, size of p_k and q_k are much smaller
- Instead of storing B_k , we can store the previous last several p and q to estimate B_{k+1}

LBFGS

- BFGS requires us to store the complete estimate of the Hessian or inverse Hessian
- The matrix is too big to be stored in deep learning setting (millions of variables)
- Recall that $B_{k+1} = \left(I - \frac{p_k q_k^T}{q_k^T p_k} \right) B_k \left(I - \frac{q_k p_k^T}{q_k^T p_k} \right) + \frac{p_k p_k^T}{q_k^T p_k}$, size of p_k and q_k are much smaller
- Instead of storing B_k , we can store the previous last several p and q to estimate B_{k+1}
 - Let say we store the last r pairs, we need to iterate r times (instead of just once) and the estimate is less accurate

LBFGS

- BFGS requires us to store the complete estimate of the Hessian or inverse Hessian
- The matrix is too big to be stored in deep learning setting (millions of variables)
- Recall that $B_{k+1} = \left(I - \frac{p_k q_k^T}{q_k^T p_k} \right) B_k \left(I - \frac{q_k p_k^T}{q_k^T p_k} \right) + \frac{p_k p_k^T}{q_k^T p_k}$, size of p_k and q_k are much smaller
- Instead of storing B_k , we can store the previous last several p and q to estimate B_{k+1}
 - Let say we store the last r pairs, we need to iterate r times (instead of just once) and the estimate is less accurate
 - Storage requirement decreases drastically

Optimizers

L-BFGS

- **Usually works very well in full batch, deterministic mode** i.e. if you have a single, deterministic $f(x)$ then L-BFGS will probably work very nicely
- **Does not transfer very well to mini-batch setting.** Gives bad results. Adapting L-BFGS to large-scale, stochastic setting is an active area of research.

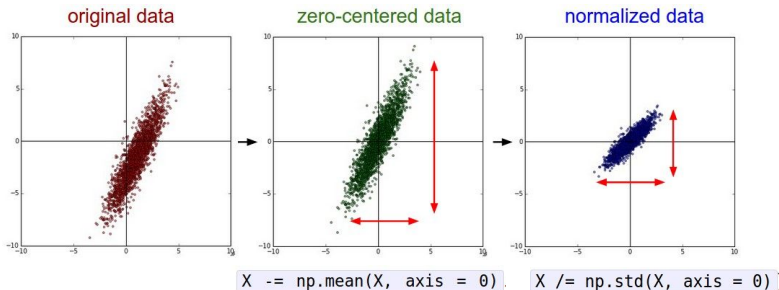
Optimizers

In practice:

- **Adam** is a good default choice in most cases
- If you can afford to do full batch updates then try out **L-BFGS** (and don't forget to disable all sources of noise)

Babysitting learning process

Step 1: Preprocess the data

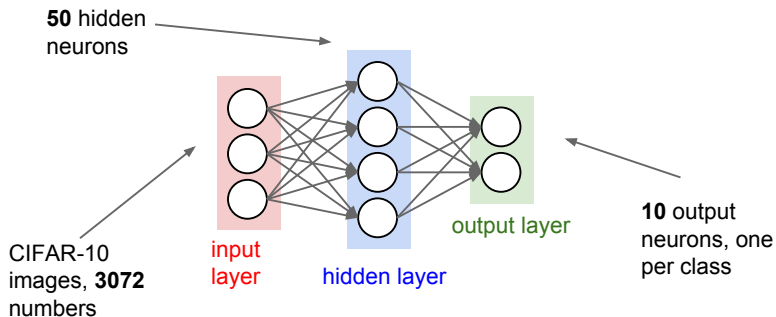


(Assume X [NxD] is data matrix,
each example in a row)

Babysitting learning process

Step 2: Choose the architecture:

say we start with one hidden layer of 50 neurons:



Babysitting learning process

Double check that the loss is reasonable:

```
def init_two_layer_model(input_size, hidden_size, output_size):
    # initialize a model
    model = {}
    model['W1'] = 0.0001 * np.random.randn(input_size, hidden_size)
    model['b1'] = np.zeros(hidden_size)
    model['W2'] = 0.0001 * np.random.randn(hidden_size, output_size)
    model['b2'] = np.zeros(output_size)
    return model
```

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
loss, grad = two_layer_net(X_train, model, y_train, 0.0) # disable regularization
print loss
```

2.30261216167

loss ~2.3.
"correct" for
10 classes

returns the loss and the
gradient for all parameters

Debugging optimizer

Double check that the loss is reasonable:

```
def init_two_layer_model(input_size, hidden_size, output_size):  
    # initialize a model  
    model = {}  
    model['W1'] = 0.0001 * np.random.randn(input_size, hidden_size)  
    model['b1'] = np.zeros(hidden_size)  
    model['W2'] = 0.0001 * np.random.randn(hidden_size, output_size)  
    model['b2'] = np.zeros(output_size)  
    return model
```

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes  
loss, grad = two_layer_net(X_train, model, y_train, 1e3) # crank up regularization  
print loss
```

3.06859716482

loss went up, good. (sanity check)

Debugging optimizer

Lets try to train now...

Tip: Make sure that you can overfit very small portion of the training data

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
X_tiny = X_train[:20] # take 20 examples
y_tiny = y_train[:20]
best_model, stats = trainer.train(X_tiny, y_tiny, X_tiny, y_tiny,
                                  model, two_layer_net,
                                  num_epochs=200, reg=0.0,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = False,
                                  learning_rate=1e-3, verbose=True)
```

The above code:

- take the first 20 examples from CIFAR-10
- turn off regularization (reg = 0.0)
- use simple vanilla 'sgd'

Debugging optimizer

Lets try to train now...

Tip: Make sure that you can overfit very small portion of the training data

Very small loss, train accuracy 1.00, nice!

```

model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
X_tiny = X_train[:20] # take 20 examples
y_tiny = y_train[:20]
best_model, stats = trainer.train(X_tiny, y_tiny, X_tiny, y_tiny,
    model, two layer net,
    num_epochs=200, reg=0.0,
    update='sgd', learning_rate_decay=1,
    sample_batches = False,
    learning_rate=1e-3, verbose=True)

```

```

Finished epoch 1 / 200: cost 2.302603, train: 0.400000, val 0.400000, lr 1.000000e-03
Finished epoch 2 / 200: cost 2.302258, train: 0.450000, val 0.450000, lr 1.000000e-03
Finished epoch 3 / 200: cost 2.301849, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 4 / 200: cost 2.301196, train: 0.650000, val 0.650000, lr 1.000000e-03
Finished epoch 5 / 200: cost 2.300044, train: 0.650000, val 0.650000, lr 1.000000e-03
Finished epoch 6 / 200: cost 2.297864, train: 0.550000, val 0.550000, lr 1.000000e-03
Finished epoch 7 / 200: cost 2.293595, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 8 / 200: cost 2.285096, train: 0.550000, val 0.550000, lr 1.000000e-03
Finished epoch 9 / 200: cost 2.268094, train: 0.550000, val 0.550000, lr 1.000000e-03
Finished epoch 10 / 200: cost 2.234787, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 11 / 200: cost 2.173187, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 12 / 200: cost 2.076862, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 13 / 200: cost 1.974090, train: 0.400000, val 0.400000, lr 1.000000e-03
Finished epoch 14 / 200: cost 1.895885, train: 0.400000, val 0.400000, lr 1.000000e-03
Finished epoch 15 / 200: cost 1.820876, train: 0.450000, val 0.450000, lr 1.000000e-03
Finished epoch 16 / 200: cost 1.737430, train: 0.450000, val 0.450000, lr 1.000000e-03
Finished epoch 17 / 200: cost 1.642356, train: 0.500000, val 0.500000, lr 1.000000e-03
Finished epoch 18 / 200: cost 1.535239, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 19 / 200: cost 1.421527, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 20 / 200: cost 1.295760, train: 0.600000, val 0.600000, lr 1.000000e-03
Finished epoch 195 / 200: cost 0.002694, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 196 / 200: cost 0.002674, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 197 / 200: cost 0.002655, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 198 / 200: cost 0.002635, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 199 / 200: cost 0.002617, train: 1.000000, val 1.000000, lr 1.000000e-03
Finished epoch 200 / 200: cost 0.002597, train: 1.000000, val 1.000000, lr 1.000000e-03
finished optimization. best validation accuracy: 1.000000

```

Debugging optimizer

Lets try to train now...

I like to start with small regularization and find learning rate that makes the loss go down.

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = True,
                                  learning_rate=1e-6, verbose=True)
```

Debugging optimizer

Lets try to train now...

I like to start with small regularization and find learning rate that makes the loss go down.

```

model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches=True,
                                  learning_rate=1e-6, verbose=True)
Finished epoch 1 / 10: cost 2.302576, train: 0.080000, val 0.103000, lr 1.000000e-06
Finished epoch 2 / 10: cost 2.302582, train: 0.121000, val 0.124000, lr 1.000000e-06
Finished epoch 3 / 10: cost 2.302558, train: 0.119000, val 0.138000, lr 1.000000e-06
Finished epoch 4 / 10: cost 2.302519, train: 0.127000, val 0.151000, lr 1.000000e-06
Finished epoch 5 / 10: cost 2.302517, train: 0.158000, val 0.171000, lr 1.000000e-06
Finished epoch 6 / 10: cost 2.302518, train: 0.179000, val 0.172000, lr 1.000000e-06
Finished epoch 7 / 10: cost 2.302466, train: 0.180000, val 0.176000, lr 1.000000e-06
Finished epoch 8 / 10: cost 2.302452, train: 0.175000, val 0.185000, lr 1.000000e-06
Finished epoch 9 / 10: cost 2.302459, train: 0.206000, val 0.192000, lr 1.000000e-06
Finished epoch 10 / 10: cost 2.302420, train: 0.190000, val 0.192000, lr 1.000000e-06
finished optimization. best validation accuracy: 0.192000
  
```

Loss barely changing

Debugging optimizer

Lets try to train now...

I like to start with small regularization and find learning rate that makes the loss go down.

loss not going down:
learning rate too low

```

model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches=True,
                                  learning_rate=1e-6, verbose=True)
Finished epoch 1 / 10: cost 2.302576, train: 0.080000, val 0.103000, lr 1.000000e-06
Finished epoch 2 / 10: cost 2.302582, train: 0.121000, val 0.124000, lr 1.000000e-06
Finished epoch 3 / 10: cost 2.302558, train: 0.119000, val 0.138000, lr 1.000000e-06
Finished epoch 4 / 10: cost 2.302519, train: 0.127000, val 0.151000, lr 1.000000e-06
Finished epoch 5 / 10: cost 2.302517, train: 0.158000, val 0.171000, lr 1.000000e-06
Finished epoch 6 / 10: cost 2.302518, train: 0.179000, val 0.172000, lr 1.000000e-06
Finished epoch 7 / 10: cost 2.302466, train: 0.180000, val 0.176000, lr 1.000000e-06
Finished epoch 8 / 10: cost 2.302452, train: 0.175000, val 0.185000, lr 1.000000e-06
Finished epoch 9 / 10: cost 2.302459, train: 0.206000, val 0.192000, lr 1.000000e-06
Finished epoch 10 / 10: cost 2.302420, train: 0.190000, val 0.192000, lr 1.000000e-06
finished optimization. best_validation accuracy: 0.192000
  
```

Loss barely changing: Learning rate is probably too low

Debugging optimizer

Lets try to train now...

I like to start with small regularization and find learning rate that makes the loss go down.

loss not going down:
learning rate too low

```

model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches=True,
                                  learning_rate=1e-6, verbose=True)
Finished epoch 1 / 10: cost 2.302576, train: 0.080000, val 0.103000, lr 1.000000e-06
Finished epoch 2 / 10: cost 2.302582, train: 0.121000, val 0.124000, lr 1.000000e-06
Finished epoch 3 / 10: cost 2.302558, train: 0.119000, val 0.138000, lr 1.000000e-06
Finished epoch 4 / 10: cost 2.302519, train: 0.127000, val 0.151000, lr 1.000000e-06
Finished epoch 5 / 10: cost 2.302517, train: 0.158000, val 0.171000, lr 1.000000e-06
Finished epoch 6 / 10: cost 2.302518, train: 0.179000, val 0.172000, lr 1.000000e-06
Finished epoch 7 / 10: cost 2.302466, train: 0.180000, val 0.176000, lr 1.000000e-06
Finished epoch 8 / 10: cost 2.302452, train: 0.175000, val 0.185000, lr 1.000000e-06
Finished epoch 9 / 10: cost 2.302459, train: 0.206000, val 0.192000, lr 1.000000e-06
Finished epoch 10 / 10: cost 2.302420, train: 0.190000, val 0.192000, lr 1.000000e-06
finished optimization. best_validation accuracy: 0.192000
  
```

Loss barely changing: Learning rate is probably too low

Notice train/val accuracy goes to 20% though, what's up with that? (remember this is softmax)

Debugging optimizer

Lets try to train now...

I like to start with small regularization and find learning rate that makes the loss go down.

loss not going down:
learning rate too low

```
model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = True,
                                  learning_rate=1e6, verbose=True)
```

Okay now lets try learning rate 1e6. What could possibly go wrong?

Debugging optimizer

Lets try to train now...

I like to start with small regularization and find learning rate that makes the loss go down.

loss not going down:
learning rate too low
loss exploding:
learning rate too high

```

model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = True,
                                  learning_rate=1e6, verbose=True)
/home/karpathy/cs231n/code/cs231n/classifiers/neural_net.py:50: RuntimeWarning: divide by zero en
countered in log
  data_loss = -np.sum(np.log(probs[range(N), y])) / N
/home/karpathy/cs231n/code/cs231n/classifiers/neural_net.py:48: RuntimeWarning: invalid value enc
ountered in subtract
  probs = np.exp(scores - np.max(scores, axis=1, keepdims=True))
Finished epoch 1 / 10: cost nan, train: 0.091000, val 0.087000, lr 1.000000e+06
Finished epoch 2 / 10: cost nan, train: 0.095000, val 0.087000, lr 1.000000e+06
Finished epoch 3 / 10: cost nan, train: 0.100000, val 0.087000, lr 1.000000e+06

```

cost: NaN almost
always means high
learning rate...

Debugging optimizer

Lets try to train now...

I like to start with small regularization and find learning rate that makes the loss go down.

loss not going down:
learning rate too low
loss exploding:
learning rate too high

```

model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
trainer = ClassifierTrainer()
best_model, stats = trainer.train(X_train, y_train, X_val, y_val,
                                  model, two_layer_net,
                                  num_epochs=10, reg=0.000001,
                                  update='sgd', learning_rate_decay=1,
                                  sample_batches = True,
                                  learning_rate=3e-3, verbose=True)
Finished epoch 1 / 10: cost 2.186654, train: 0.308000, val 0.306000, lr 3.000000e-03
Finished epoch 2 / 10: cost 2.176230, train: 0.330000, val 0.350000, lr 3.000000e-03
Finished epoch 3 / 10: cost 1.942257, train: 0.376000, val 0.352000, lr 3.000000e-03
Finished epoch 4 / 10: cost 1.827868, train: 0.329000, val 0.310000, lr 3.000000e-03
Finished epoch 5 / 10: cost inf, train: 0.128000, val 0.128000, lr 3.000000e-03
Finished epoch 6 / 10: cost inf, train: 0.144000, val 0.147000, lr 3.000000e-03

```

3e-3 is still too high. Cost explodes....

=> Rough range for learning rate we should be cross-validating is somewhere [1e-3 ... 1e-5]

Hyperparameter optimization

Hyperparameter Optimization

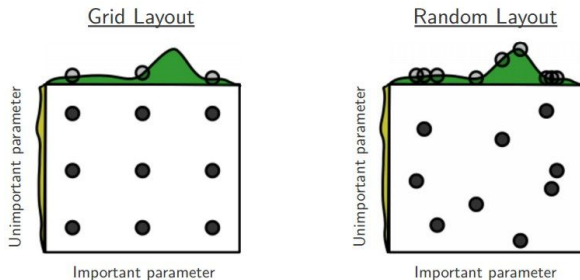
Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 85

20 Jan 2016

Hyperparameter optimization

Random Search vs. Grid Search



Random Search for Hyper-Parameter Optimization
Bergstra and Bengio, 2012

Hyperparameter optimization

Cross-validation strategy

I like to do **coarse** -> **fine** cross-validation in stages

First stage: only a few epochs to get rough idea of what params work

Second stage: longer running time, finer search

... (repeat as necessary)

Tip for detecting explosions in the solver:

If the cost is ever $> 3 * \text{original cost}$, break out early

Hyperparameter optimization

For example: run coarse search for 5 epochs

```

max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-5, 5)
    lr = 10**uniform(-3, -6)

    trainer = ClassifierTrainer()
    model = init_two_layer_model(32*32*3, 50, 10) # input size, hidden size, number of classes
    trainer = ClassifierTrainer()
    best_model_local, stats = trainer.train(X_train, y_train, X_val, y_val,
                                           model, two_layer_net,
                                           num_epochs=5, reg=reg,
                                           update='momentum', learning_rate_decay=0.9,
                                           sample_batches = True, batch_size = 100,
                                           learning_rate=lr, verbose=False)
  
```

note it's best to optimize
in log space!

```

val_acc: 0.412000, lr: 1.405206e-04, reg: 4.793564e-01, (1 / 100)
val_acc: 0.214000, lr: 7.231888e-06, reg: 2.321281e-04, (2 / 100)
val_acc: 0.208000, lr: 2.119571e-06, reg: 8.011857e+01, (3 / 100)
val_acc: 0.196000, lr: 1.551131e-05, reg: 4.374936e-05, (4 / 100)
val_acc: 0.079000, lr: 1.753300e-05, reg: 1.200424e+03, (5 / 100)
val_acc: 0.223000, lr: 4.215128e-05, reg: 4.196174e+01, (6 / 100)
val_acc: 0.441000, lr: 1.750259e-04, reg: 2.110807e-04, (7 / 100)
val_acc: 0.241000, lr: 6.749231e-05, reg: 4.226413e+01, (8 / 100)
val_acc: 0.482000, lr: 4.296863e-04, reg: 6.642555e-01, (9 / 100)
val_acc: 0.079000, lr: 5.401602e-06, reg: 1.599828e+04, (10 / 100)
val_acc: 0.154000, lr: 1.618508e-06, reg: 4.925252e-01, (11 / 100)
  
```

nice

Hyperparameter optimization

Now run finer search...

```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-5, 5)
    lr = 10**uniform(-3, -6)
```

adjust range



```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-4, 0)
    lr = 10**uniform(-3, -4)
```

```
val_acc: 0.527000, lr: 5.340517e-04, reg: 4.097824e-01, (0 / 100)
val_acc: 0.492000, lr: 2.279484e-04, reg: 9.991345e-04, (1 / 100)
val_acc: 0.512000, lr: 8.680827e-04, reg: 1.349727e-02, (2 / 100)
val_acc: 0.461000, lr: 1.028377e-04, reg: 1.220193e-02, (3 / 100)
val_acc: 0.460000, lr: 1.113730e-04, reg: 5.244309e-02, (4 / 100)
val_acc: 0.498000, lr: 9.477776e-04, reg: 2.001293e-03, (5 / 100)
val_acc: 0.469000, lr: 1.484369e-04, reg: 4.328313e-01, (6 / 100)
val_acc: 0.522000, lr: 5.586261e-04, reg: 2.312685e-04, (7 / 100)
val_acc: 0.530000, lr: 5.808183e-04, reg: 8.259964e-02, (8 / 100)
val_acc: 0.489000, lr: 1.979168e-04, reg: 1.010889e-04, (9 / 100)
val_acc: 0.490000, lr: 2.036031e-04, reg: 2.406271e-03, (10 / 100)
val_acc: 0.475000, lr: 2.021162e-04, reg: 2.287807e-01, (11 / 100)
val_acc: 0.460000, lr: 1.135527e-04, reg: 3.905040e-02, (12 / 100)
val_acc: 0.515000, lr: 6.947668e-04, reg: 1.562808e-02, (13 / 100)
val_acc: 0.531000, lr: 9.471549e-04, reg: 1.433895e-03, (14 / 100)
val_acc: 0.509000, lr: 3.140888e-04, reg: 2.857518e-01, (15 / 100)
val_acc: 0.514000, lr: 6.438349e-04, reg: 3.033781e-01, (16 / 100)
val_acc: 0.502000, lr: 3.921784e-04, reg: 2.707126e-04, (17 / 100)
val_acc: 0.509000, lr: 9.752279e-04, reg: 2.850865e-03, (18 / 100)
val_acc: 0.500000, lr: 2.412048e-04, reg: 4.997821e-04, (19 / 100)
val_acc: 0.466000, lr: 1.319314e-04, reg: 1.189915e-02, (20 / 100)
val_acc: 0.516000, lr: 8.039527e-04, reg: 1.528291e-02, (21 / 100)
```

53% - relatively good
for a 2-layer neural net
with 50 hidden neurons.

Hyperparameter optimization

Now run finer search...

```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-5, 5)
    lr = 10**uniform(-3, -6)
```

adjust range



```
max_count = 100
for count in xrange(max_count):
    reg = 10**uniform(-4, 0)
    lr = 10**uniform(-3, -4)
```

```
val_acc: 0.527000, lr: 5.340517e-04, reg: 4.097824e-01, (0 / 100)
val_acc: 0.492000, lr: 2.279484e-04, reg: 9.991345e-04, (1 / 100)
val_acc: 0.512000, lr: 8.680827e-04, reg: 1.349727e-02, (2 / 100)
val_acc: 0.461000, lr: 1.028377e-04, reg: 1.220193e-02, (3 / 100)
val_acc: 0.460000, lr: 1.113730e-04, reg: 5.244309e-02, (4 / 100)
val_acc: 0.498000, lr: 9.477776e-04, reg: 2.001293e-03, (5 / 100)
val_acc: 0.469000, lr: 1.484369e-04, reg: 4.328313e-01, (6 / 100)
val_acc: 0.522000, lr: 5.586261e-04, reg: 2.312685e-04, (7 / 100)
val_acc: 0.530000, lr: 5.808183e-04, reg: 8.259964e-02, (8 / 100)
val_acc: 0.489000, lr: 1.979168e-04, reg: 1.010889e-04, (9 / 100)
val_acc: 0.490000, lr: 2.036031e-04, reg: 2.406271e-03, (10 / 100)
val_acc: 0.475000, lr: 2.021162e-04, reg: 2.287807e-01, (11 / 100)
val_acc: 0.460000, lr: 1.135527e-04, reg: 3.905040e-02, (12 / 100)
val_acc: 0.515000, lr: 6.947668e-04, reg: 1.562808e-02, (13 / 100)
val_acc: 0.531000, lr: 9.471549e-04, reg: 1.433895e-03, (14 / 100)
val_acc: 0.509000, lr: 3.140888e-04, reg: 2.857518e-01, (15 / 100)
val_acc: 0.514000, lr: 6.438349e-04, reg: 3.033781e-01, (16 / 100)
val_acc: 0.502000, lr: 3.921784e-04, reg: 2.707126e-04, (17 / 100)
val_acc: 0.509000, lr: 9.752279e-04, reg: 2.850865e-03, (18 / 100)
val_acc: 0.500000, lr: 2.412048e-04, reg: 4.997821e-04, (19 / 100)
val_acc: 0.466000, lr: 1.319314e-04, reg: 1.189915e-02, (20 / 100)
val_acc: 0.516000, lr: 8.039527e-04, reg: 1.528291e-02, (21 / 100)
```

53% - relatively good
for a 2-layer neural net
with 50 hidden neurons.

But this best cross-
validation result is
worrying. Why?

Hyperparameter optimization

Hyperparameters to play with:

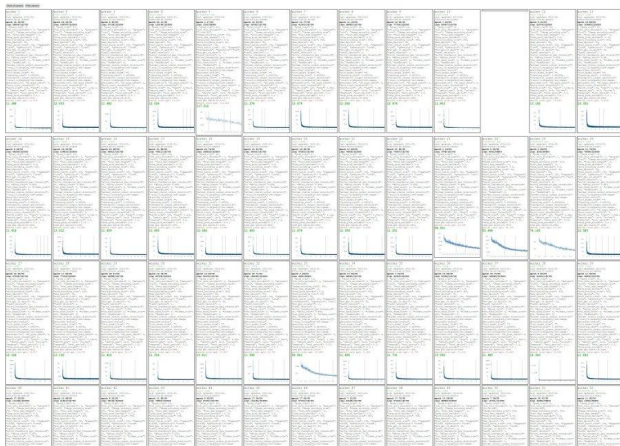
- network architecture
- learning rate, its decay schedule, update type
- regularization (L2/Dropout strength)

neural networks practitioner
music = loss function



Hyperparameter optimization

My cross-validation
“command center”



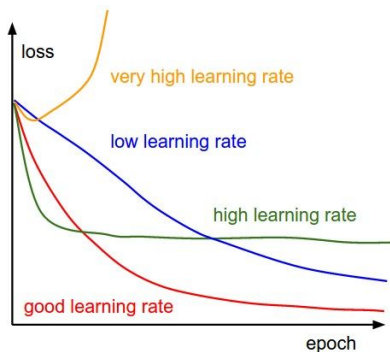
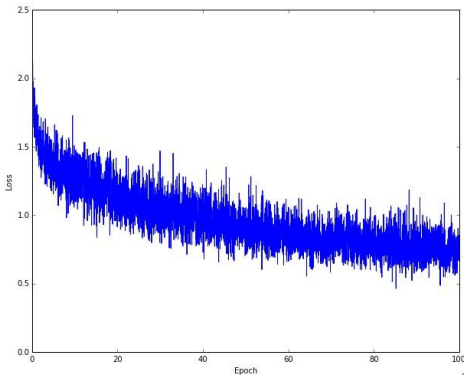
Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 92

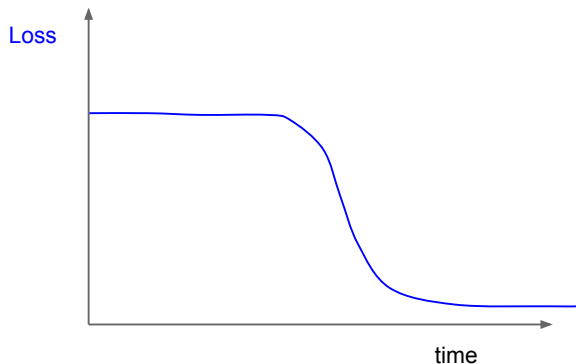
20 Jan 2016

Hyperparameter optimization

Monitor and visualize the loss curve



Hyperparameter optimization

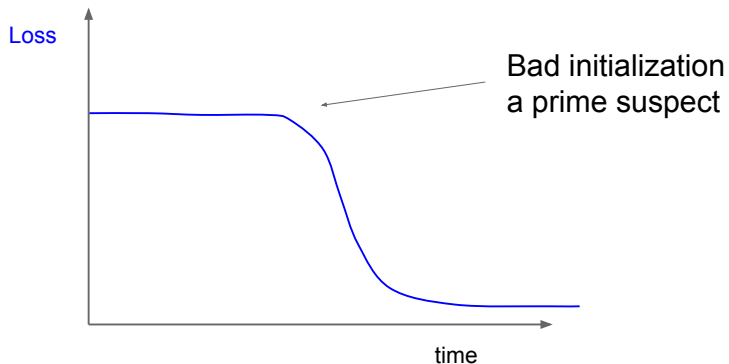


Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 94

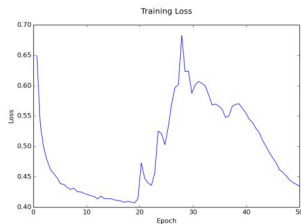
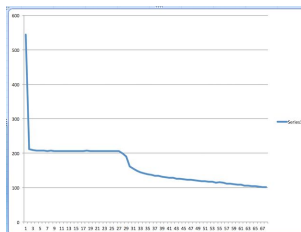
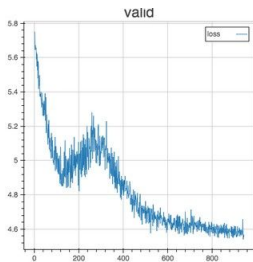
20 Jan 2016

Hyperparameter optimization



Hyperparameter optimization

lossfunctions.tumblr.com Loss function specimen



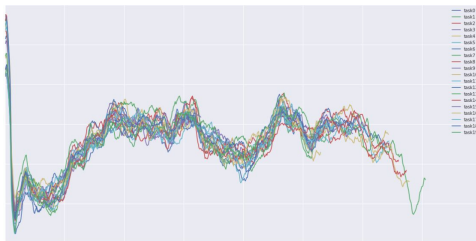
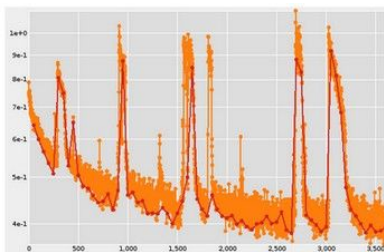
Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 96

20 Jan 2016

Hyperparameter optimization

lossfunctions.tumblr.com

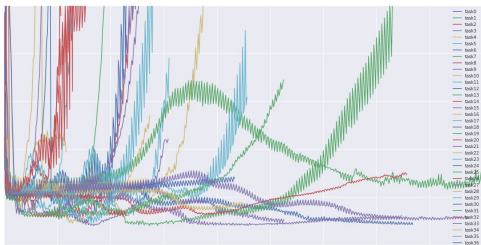
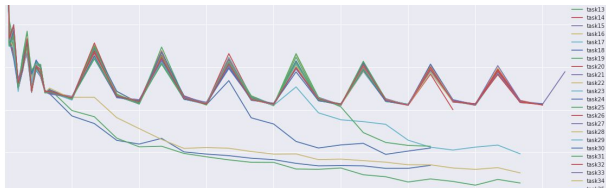


Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 97

20 Jan 2016

Hyperparameter optimization



lossfunctions.tumblr.com

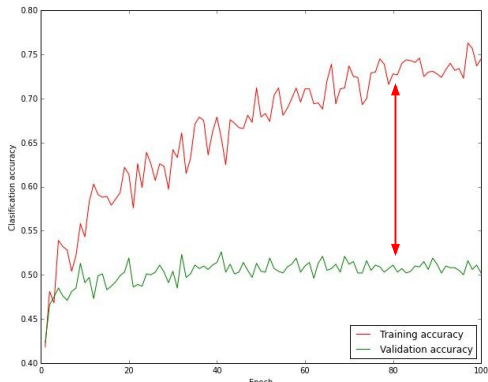
Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 - 98

20 Jan 2016

Hyperparameter optimization

Monitor and visualize the accuracy:



big gap = overfitting

=> increase regularization strength?

no gap

=> increase model capacity?

Hyperparameter optimization

Track the ratio of weight updates / weight magnitudes:

```
# assume parameter vector W and its gradient vector dW
param_scale = np.linalg.norm(W.ravel())
update = -learning_rate*dW # simple SGD update
update_scale = np.linalg.norm(update.ravel())
W += update # the actual update
print update_scale / param_scale # want ~1e-3
```

ratio between the values and updates: $\sim 0.0002 / 0.02 = 0.01$ (about okay)
want this to be somewhere around 0.001 or so

Conclusions (What we know in 2017)

- BP is just chain rule in calculus
- Use ReLU. Never use Sigmoid (use Tanh instead)
- Input preprocessing is no longer very important
 - Do subtract mean
 - Whitening and normalizing are not much needed
- Weight initialization on the other hand is extremely important for deep networks
- Use batch normalization if you can
- Use dropout
- Use Adam (or maybe RMSprop) for optimizer. If you don't have much data, can consider LBFGS
- Need to babysit your learning for real-world problems
- Never use grid search for tuning your hyperparameters