Deep Neural Networks

Computer Vision (UW EE/CSE 576)



Richard Szeliski Facebook & UW Lecture 7 – Apr 21, 2020



Class calendar

Date	Торіс	Slides	Reading	Homework
April 9	Filters and convolutions	Google Slides	<u>Szeliski</u> , Chapter 3	HW1 due, <u>HW2</u> assigned
April 14	Interpolation and Optimization	<u>pdf, pptx</u>	Szeliski, Chapter 4	
April 16	Machine Learning	<u>pdf, pptx</u>	Szeliski, Chapter 5.1-5.2	
April 21	Deep Neural Networks		Szeliski, Chapter 5.3	
April 23	Convolutional Neural Networks		Szeliski, Chapter 5.4	HW2 due, HW3 assigned
April 28	Network Architectures		Szeliski, Chapter 5.4-5.5	
April 30	Object Detection		Szeliski, Chapter 6.3	
May 5	Detection and Instance Segmentation		<u>Szeliski</u> , Chapter 6.4	
May 7	Edges, features, matching, RANSAC		<u>Szeliski</u> , Chapter 7.1-7.2, 8.1-8.2	HW3 due, HW4 assigned

References



(Thanks, Matt Dietke, for the very helpful comments.)





https://d2l.ai/

Chapter 5 **Deep Learning**

5.1	Superv	vised learning	239
	5.1.1	Nearest neighbors	241
	5.1.2	Bayesian classification	243
	5.1.3	Logistic regression	247
	5.1.4	Support vector machines	249
	5.1.5	Decision trees and forests	251
5.2	Unsup	ervised learning	254
	5.2.1	Clustering	254
	5.2.2	K-means and mixtures of Gaussians	256
	5.2.3	Principal component analysis	258
	5.2.4	Semi-supervised learning	268
5.3	Deep 1	neural networks	268
	5.3.1	Weights and layers	270
	5.3.2	Activation functions	272
	5.3.3	Regularization and normalization	274
	5.3.4	Loss functions	279
	5.3.5	Backpropagation	281
	5.3.6	Training and optimization	284
54	Convo	lutional neural networks	288

Readings

Deep neural networks (today's lecture)

- Loss functions
- Regularization
- Weights and layers
- Activation functions
- Backpropagation
- Training and optimization (?)





Machine learning (previous lecture)

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines
- Clustering
- Principal component analysis





Supervised learning



- Goal: provide best output predictions for novel inputs
- How can we predict future performance?
- Placeholder answer: minimize *empirical risk*

$$E_{\text{Risk}}(\mathbf{w}) = \frac{1}{N} \sum L(\mathbf{y}_i, \mathbf{f}(\mathbf{x}_i; \mathbf{w})).$$
(5.1)

The loss function L measures the "cost" of predicting an output $f(x_i; w)$ for input x_i and model parameters w when the corresponding target is y_i .

• What are potential *models* and *loss functions*?

Traditional and deep learning



Bayesian classification





Logistic regression

$$p(\mathcal{C}_0|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b), \tag{5.9}$$

Equation (5.9), which we will revist shortly in the context of non-generative (discriminative) classification (5.18), is called *logistic regression*, since we pass the output of a linear regression formula

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b \tag{5.12}$$

through the logistic function to obtain a class probability. Figure 5.7 illustrates this in two dimensions, there the posterior likelihood of the red class $p(C_0|\mathbf{x})$ is shown on the right side.

In linear regression (5.12), w plays the role a the *weight* vector along which we project the feature vector \mathbf{x} , and b plays the role of the *bias*, which determines where to set the





As before, I'm borrowing slides from

EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Course Description

UNIVERSITY OF MICHIGAN

Computer Vision has become ubiquitous in our society, with applications in search, image understanding, apps, mapping, medicine, drones, and self-driving cars. Core to many of these applications are visual recognition tasks such as image classification and object detection. Recent developments in neural network approaches have greatly advanced the performance of these state-of-the-art visual recognition systems. This course is a deep dive into details of neural-network based deep learning methods for computer vision. During this course, students will learn to implement, train and debug their own neural networks and gain a detailed understanding of cutting-edge research in computer vision. We will cover learning algorithms, neural network architectures, and practical engineering tricks for training and fine-tuning networks for visual recognition tasks.

Instructor Graduate Student Instructors





EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 3: Linear Classifiers



Justin Johnson



Neural Network



This image is CC0 1.0 public domain





Linear Classifier: Three Viewpoints

Algebraic Viewpoint

f(x,W) = Wx



Visual Viewpoint

One template per class



Geometric Viewpoint

Hyperplanes cutting up space



Justin Johnson



So Far: Defined a linear <u>score function</u> f(x,W) = Wx + b







airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

Given a W, we can compute class scores for an image x.

But how can we actually choose a good W?

Cat image by Nikita is licensed under CC-BY 2.0; Car image is CCO 1.0 public domain; Frog image is in the public domain

Justin Johnson



Choosing a good W







airplane	-3.45	-0.51	3.42
automobile	-8.87	6.04	4.64
bird	0.09	5.31	2.65
cat	2.9	-4.22	5.1
deer	4.48	-4.19	2.64
dog	8.02	3.58	5.55
frog	3.78	4.49	-4.34
horse	1.06	-4.37	-1.5
ship	-0.36	-2.09	-4.79
truck	-0.72	-2.93	6.14

TODO:

- Use a loss function to quantify how good a value of W is
- Find a W that minimizes the loss function (optimization)

Justin Johnson

Lecture 3 - 16



f(x,W) = Wx + b

A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)



A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc

Justin Johnson



A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where $oldsymbol{x_i}$ is image and $oldsymbol{y_i}$ is (integer) label

Justin Johnson



A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where $oldsymbol{x_i}$ is image and $oldsymbol{y_i}$ is (integer) label

Loss for a single example is $L_i(f(x_i, W), y_i)$

Justin Johnson



A **loss function** tells how good our current classifier is

Low loss = good classifier High loss = bad classifier

(Also called: **objective function**; **cost function**)

Negative loss function sometimes called **reward function**, **profit function**, **utility function**, **fitness function**, etc Given a dataset of examples

$$\{(x_i, y_i)\}_{i=1}^N$$

Where $oldsymbol{x_i}$ is image and $oldsymbol{y_i}$ is (integer) label

Loss for a single example is $L_i(f(x_i, W), y_i)$ Loss for the dataset is average of

per-example losses:

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

Regularization: Beyond Training Error

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

Data loss: Model predictions should match training data





Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W) \qquad \begin{array}{l} \lambda_i = \text{regularization strength} \\ \text{(hyperparameter)} \end{array}$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too* well on training data

Simple examplesMore complex:L2 regularization: $R(W) = \sum_k \sum_l W_{k,l}^2$ DropoutL1 regularization: $R(W) = \sum_k \sum_l |W_{k,l}|$ Batch normalizationElastic net (L1 + L2): $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$ Cutout, Mixup, Stochastic depth, etc...

Justin Johnson



Regularization: Beyond Training Error

$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(f(x_i, W), y_i) + \lambda R(W) \qquad \begin{array}{l} \lambda_i = \text{regularization strength} \\ \text{(hyperparameter)} \end{array}$$

Data loss: Model predictions should match training data

Regularization: Prevent the model from doing *too* well on training data

Purpose of Regularization:

- Express preferences in among models beyond "minimize training error"
- Avoid **overfitting**: Prefer simple models that generalize better
- Improve optimization by adding curvature

Cross-Entropy Loss (Multinomial Logistic Regression)

Want to interpret raw classifier scores as probabilities



Justin Johnson



Cross-Entropy Loss (Multinomial Logistic Regression)

Want to interpret raw classifier scores as probabilities

Fall 2019

			s = f(x)	$e_i; W)$	P(Y =	$=k X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$	Softmax function
		P m	Probabilities nust be >= C	Pro mu	obabilities st sum to 1	$L_i = -\log P(Y=y)$	$_i X=x_i)$
cat	3.2		24.5		0.13	> Compare	1.00
car	5.1	exp	164.0	normalize	0.87	Kullback–Leibler divergence	0.00
frog	-1.7		0.18		0.00	$D_{KL}(P \ Q) =$	0.00
Un pro	normalized lo babilities / lo	og- U gits J	nnormalize probabilities	d p	robabilities	$\sum_{y} P(y) \log \frac{P(y)}{Q(y)}$	Correct probs

Lecture 3 - 26

Justin Johnson

Cross-Entropy Loss (Multinomial Logistic Regression)

Want to interpret raw classifier scores as probabilities

			s = f(x)	$e_i; W)$	P(Y =	$=k X=x_i)=rac{e^{s_k}}{\sum_j e^{s_j}}$	Softmax function
		P m	Probabilities nust be >= 0	Pro mu	obabilities st sum to 1	$L_i = -\log P(Y=y_i)$	$_i X=x_i)$
cat	3.2		24.5		0.13	🔶 Compare ←	1.00
car	5.1	exp	164.0	normalize	0.87	Cross Entropy	0.00
frog	-1.7		0.18		0.00	H(P,Q) =	0.00
Un pro	normalized lo babilities / lo	og- u gits l	nnormalize probabilities	d s	robabilities	$H(p) + D_{KL}(P \ Q)$	Correct probs

Lecture 3 - 27

Justin Johnson

Recap: Three ways to think about linear classifiers

Algebraic Viewpoint

f(x,W) = Wx



Visual Viewpoint

One template per class



Geometric Viewpoint

Hyperplanes cutting up space



Fall 2019

Justin Johnson

Recap: Loss Functions quantify preferences

- We have some dataset of (x, y)
- We have a **score function**:
- We have a **loss function**:

$$s = f(x;W) = Wx$$
Linear classifier





Recap: Loss Functions quantify preferences

- We have some dataset of (x, y)
- We have a **score function**:
- We have a **loss function**:

Q: How do we find the best W?

$$s = f(x; W) = Wx$$
Linear classifier

A: Later in this lecture



Justin Johnson



Machine learning

- Why learning?
- Nearest neighbors
- Bayesian classification
- Logistic regression
- Support vector machines
- Clustering
- Principal component analysis

Support vector machines (SVMs)

- Maximize the *margin* between the decision surfaces
- The only points that matter are the circled *support vectors*



n margin classifier, we need to $l_i = \mathbf{w} \cdot \mathbf{x}_i + b$ (5.17) have an note this more compactly, let

$$\hat{t}_i = 2t_i - 1, \quad \hat{t}_i \in \{-1, 1\}$$

an now re-write the inequality c

$$\hat{t}_i(\mathbf{w} \cdot \mathbf{x}_i + b) \ge 1.$$

simply find the smallest norm ization problem

 $\arg\min_{\mathbf{w},b}\|\mathbf{w}\|^2$

assic quadratic programming p

Kernelized support vector machines

 Replace linear function with a sum of *kernel functions* (e.g., Gaussian bumps)

$$l = f(\mathbf{x}) = \sum_{k} \mathbf{w}_{k} \phi(\|\mathbf{x} - \mathbf{x}_{k}\|).$$

 Circled points are support vectors, lie on f = ± 1 surfaces (f = 0 is the dark curve)



Hinge loss vs. logistic regression

Figure 7.5 Plot of the 'hinge' error function used in support vector machines, shown in blue, along with the error function for logistic regression, rescaled by a factor of $1/\ln(2)$ so that it passes through the point (0, 1), shown in red. Also shown are the misclassification error in black and the squared error in green.

remaining points we have $\xi_n = 1 - y_n t_n$. Thus the objective function (7.21) can be written (up to an overall multiplicative constant) in the form

$$\sum_{n=1}^{N} E_{\rm SV}(y_n t_n) + \lambda \|\mathbf{w}\|^2$$
(7.44)

where $\lambda = (2C)^{-1}$, and $E_{SV}(\cdot)$ is the *hinge* error function defined by

$$E_{\rm SV}(y_n t_n) = [1 - y_n t_n]_+ \tag{7.45}$$

where $[\cdot]_+$ denotes the positive part. The hinge error function, so-called because of its shape, is plotted in Figure 7.5. It can be viewed as an approximation to the misclassification error, i.e., the error function that ideally we would like to minimize, which is also shown in Figure 7.5.



Richard Szeliski

Deep neural networks (today's lecture)

- Loss functions
- Regularization
- Weights and layers
- Activation functions
- Backpropagation
- Training and optimization







EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 5: Neural Networks

Justin Johnson


Problem: Linear Classifiers aren't that powerful

Geometric Viewpoint



Visual Viewpoint

One template per class: Can't recognize different modes of a class



Justin Johnson









Justin Johnson

Lecture 5 - 45

Fall 2019



Justin Johnson





Image Features: Color Histogram



Frog image is in the public domain





Image Features: Histogram of Oriented Gradients (HoG)



- 1. Compute edge direction / strength at each pixel
- 2. Divide image into 8x8 regions
- Within each region compute a histogram of edge directions weighted by edge strength

Lowe, "Object recognition from local scale-invariant features", ICCV 1999 Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

Justin Johnson



Image Features: Histogram of Oriented Gradients (HoG)



- 1. Compute edge direction / strength at each pixel
- 2. Divide image into 8x8 regions
- Within each region compute a histogram of edge directions weighted by edge strength



Example: 320x240 image gets divided into 40x30 bins; 8 directions per bin; feature vector has 30*40*9 = 10,800 numbers

> Lowe, "Object recognition from local scale-invariant features", ICCV 1999 Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

Justin Johnson



Image Features: Histogram of Oriented Gradients (HoG) Weak edges

Strong diagonal



- 1. Compute edge direction / strength at each pixel
- 2. Divide image into 8x8 regions
- Within each region compute a histogram of edge directions weighted by edge strength

edges **Edges in all** directions Captures texture and position, robust to small image changes



Example: 320x240 image gets divided into 40x30 bins; 8 directions per bin; feature vector has 30*40*9 = 10,800 numbers

> Lowe, "Object recognition from local scale-invariant features", ICCV 1999 Dalal and Triggs, "Histograms of oriented gradients for human detection," CVPR 2005

Justin Johnson



Image Features: Bag of Words (Data-Driven!)

Step 1: Build codebook



Fei-Fei and Perona, "A bayesian hierarchical model for learning natural scene categories", CVPR 2005

Justin Johnson

Car image is CC0 1.0 public domain



Image Features: Bag of Words (Data-Driven!)

Step 1: Build codebook



Step 2: Encode images





Fei-Fei and Perona, "A bayesian hierarchical model for learning natural scene categories", CVPR 2005

Justin Johnson

Image Features



Justin Johnson



Example: Winner of 2011 ImageNet challenge

Low-level feature extraction \approx 10k patches per image

SIFT: 128-dim
color: 96-dim
reduced to 64-dim with PCA

FV extraction and compression:

- N=1,024 Gaussians, R=4 regions ⇒ 520K dim x 2
- compression: G=8, b=1 bit per dimension

One-vs-all SVM learning with SGD

Late fusion of SIFT and color systems

F. Perronnin, J. Sánchez, "Compressed Fisher vectors for LSVRC", PASCAL VOC / ImageNet workshop, ICCV, 2011.





Image Features







Image Features vs Neural Networks





Krizhevsky, Sutskever, and Hinton, "Imagenet classification with deep convolutional neural networks", NIPS 2012. Figure copyright Krizhevsky, Sutskever, and Hinton, 2012. Reproduced with permission.

10 numbers giving scores for classes

training

Justin Johnson

Lecture 5 - 57

Fall 2019

(Before) Linear score function:

$$\boldsymbol{f} = \boldsymbol{W}\boldsymbol{x}$$
$$x \in \mathbb{R}^{D}, W \in \mathbb{R}^{C \times D}$$





(**Before**) Linear score function: f = (Now) 2-layer Neural Network f = f

$$f = Wx$$
$$f = W_2 \max(0, W_1 x)$$
$$\mathbb{D}^{C \times H} = W_2 - \mathbb{D}^{H \times D} = \mathbb{D}^{H \times D}$$

$$W_2 \in \mathbb{R}^{C \times H} \quad W_1 \in \mathbb{R}^{H \times D} \quad x \in \mathbb{R}^D$$

(In practice we will usually add a learnable bias at each layer as well)

Justin Johnson	Lecture 5 - 59	Fall 2019

f = Wx(**Before**) Linear score function: $f = W_2 \max(0, W_1 x)$ (**Now**) 2-layer Neural Network or 3-layer Neural Network $f = W_3 \max(0, W_2 \max(0, W_1 x))$ $W_3 \in \mathbb{R}^{C \times H_2} \quad W_2 \in \mathbb{R}^{H_2 \times H_1} \quad W_1 \in \mathbb{R}^{H_1 \times D} \quad x \in \mathbb{R}^D$

(In practice we will usually add a learnable bias at each layer as well)

Justin Johnson



(Before) Linear score function:

(Now) 2-layer Neural Network

$$egin{aligned} f &= Wx \ f &= W_2 \max(0, W_1 x) \end{aligned}$$





Justin Johnson



Fall 2019

Justin Johnson

Justin Johnson

frog

plane

dog

Neural Networks

Linear classifier: One template per class

horse



truck

ship

(**Before**) Linear score function:

(**Now**) 2-layer Neural Network



 $x \in \mathbb{R}^D, W_1 \in \mathbb{R}^{H \times D}, W_2 \in \mathbb{R}^{C \times H}$



Neural net: first layer is bank of templates; Second layer recombines templates



(Before) Linear score function:

(Now) 2-layer Neural Network



Justin Johnson

Can use different templates to cover multiple modes of a class!



(Before) Linear score function:

(Now) 2-layer Neural Network



Justin Johnson

"Distributed representation": Most templates not interpretable!



(Before) Linear score function:

(Now) 2-layer Neural Network



Justin Johnson

Deep Neural Networks Depth = number of layers Width: Size of h_4 h₁ h_3 W_1 h₂ h_5 W_6 Х W_3 W_4 W_2 W_5 S each layer Output: 10 Input: 3072

 $s = W_6 \max(0, W_6 \max(0, W_5 \max(0, W_4 \max(0, W_3 \max(0, W_2 \max(0, W_1 x)))))))$

Justin Johnson



2-layer Neural Network

The function ReLU(z) = max(0, z) is called "Rectified Linear Unit"

$$f=W_2\max(0,W_1x)$$

This is called the **activation function** of the neural network





2-layer Neural Network

The function ReLU(z) = max(0, z) is called "Rectified Linear Unit"



$$f=W_2\max(0,W_1x)$$

This is called the **activation function** of the neural network

Q: What happens if we build a neural network with no activation function?

$$s = W_2 W_1 x$$

Justin Johnson

2-layer Neural Network

The function ReLU(z) = max(0, z) is called "Rectified Linear Unit"

10

$$f=W_2\max(0,W_1x)$$

This is called the **activation function** of the neural network

Q: What happens if we build a neural network with no activation function?

 $s = W_2 W_1 x$ $W_3 = W_2 W_1 \in \mathbb{R}^{C \times H} \quad s = W_3 x$

A: We end up with a linear classifier!

Justin Johnson

-10

Lecture 5 - 71

10





Leaky ReLU $\max(0.1x, x)$



 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$



Justin Johnson





ReLU is a good default choice for most problems

Leaky ReLU $\max(0.1x, x)$



 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$



Justin Johnson



Justin Johnson

Lecture 5 - 74

Fall 2019





Justin Johnson

Lecture 5 - 76

Fall 2019

Setting the number of layers and their sizes



More hidden units = more capacity

	1 .				
	UST	In J	\mathbf{O}	nns	<u>on</u>
-					• •••



Don't regularize with size; instead use stronger L2 $\lambda = 0.001$ $\lambda = 0.01$ $\lambda = 0.1$



(Web demo with ConvNetJS:

http://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html

Justin Johnson


Summary

Feature transform + Linear classifier allows nonlinear decision boundaries



Neural Networks as learnable feature transforms



Justin Johnson



Summary

From linear classifiers to fully-connected networks



Linear classifier: One template per class



Neural networks: Many reusable templates





Deep neural networks

- Loss functions
- Regularization
- Weights and layers
- Activation functions
- Backpropagation
- Training and optimization





EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 6: Backpropagation

Justin Johnson



Problem: How to compute gradients?

$$\begin{split} s &= f(x; W_1, W_2) = W_2 \max(0, W_1 x) \quad \text{Nonlinear score function} \\ L_i &= \sum_{j \neq y_i} \max(0, s_j - s_{y_i} + 1) \quad \text{SVM Loss on predictions} \\ R(W) &= \sum_k W_k^2 \quad \text{Regularization} \\ L &= \frac{1}{N} \sum_{i=1}^N L_i + \lambda R(W_1) + \lambda R(W_2) \quad \text{Total loss: data loss + regularization} \\ \text{If we can compute } \frac{\partial L}{\partial W_1}, \frac{\partial L}{\partial W_2} \text{ then we can learn } W_1 \text{ and } W_2 \end{split}$$

(Bad) Idea: Derive $abla_W L$ on paper

$$s = f(x; W) = Wx$$

$$L_{i} = \sum_{j \neq y_{i}} \max(0, s_{j} - s_{y_{i}} + 1)$$

$$= \sum_{j \neq y_{i}} \max(0, W_{j,:} \cdot x + W_{y_{i},:} \cdot x + 1)$$

$$loss = \sum_{j \neq y_{i}} \max(0, W_{j,:} \cdot x + W_{y_{i},:} \cdot x + 1)$$

$$L = \frac{1}{N} \sum_{i=1}^{N} L_{i} + \lambda \sum_{k} W_{k}^{2}$$

$$= \frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_{i}} \max(0, W_{j,:} \cdot x + W_{y_{i},:} \cdot x + 1) + \lambda \sum_{k} W_{k}^{2}$$

$$\nabla_{W}L = \nabla_{W} \left(\frac{1}{N} \sum_{i=1}^{N} \sum_{j \neq y_{i}} \max(0, W_{j,:} \cdot x + W_{y_{i},:} \cdot x + 1) + \lambda \sum_{k} W_{k}^{2}\right)$$

Problem: Very tedious: Lots of matrix calculus, need lots of paper

Problem: What if we want to change loss? E.g. use softmax instead of SVM? Need to re-derive from scratch. Not modular!

Problem: Not feasible for very complex models!

Justin Johnson



Better Idea: Computational Graphs



J	U	S	ti	n	0	h	n	S	0	n
<u> </u>	<u> </u>	<u> </u>			 \sim				\sim	•••



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







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

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







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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$



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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Justin Johnson



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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Justin Johnson



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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



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

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

1. Forward pass: Compute outputs q = x + y f = qz

2. Backward pass: Compute derivatives

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



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

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

1. Forward pass: Compute outputs q = x + y f = qz

2. Backward pass: Compute derivatives

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



Justin Johnson



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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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





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

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

1. Forward pass: Compute outputs q = x + y f = qz

2. Backward pass: Compute derivatives

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



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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Justin Johnson



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

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

1. Forward pass: Compute outputs

$$q = x + y \quad f = qz$$

2. Backward pass: Compute derivatives

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



Backpropagation: x -2 Simple Example 3 f(x, y, z) = (x + y)zy 5 -12 e.g. x = -2, y = 5, z = -4 z -4 3 **1. Forward pass**: Compute outputs **Chain Rule** q = x + y | f = qz ∂f $\partial q \ \partial f$ $\frac{1}{\partial y} \frac{1}{\partial q}$ ∂u 2. Backward pass: Compute derivatives Want: $\frac{\partial f}{\partial x}, \ \frac{\partial f}{\partial y}, \ \frac{\partial f}{\partial z}$ Upstream **Downstream** Local Gradient Gradient Gradient Justin Johnson Lecture 6 - 98 Fall 2019

Backpropagation: x -2 Simple Example 3 f(x, y, z) = (x + y)zy <u>5</u> -12 e.g. x = -2, y = 5, z = -4 z -4 3 **1. Forward pass**: Compute outputs **Chain Rule** q = x + y | f = qz ∂f $\partial q \ \partial f$ $\frac{1}{\partial y} \overline{\partial q}$ ∂u 2. Backward pass: Compute derivatives Want: $\frac{\partial f}{\partial x}, \ \frac{\partial f}{\partial y}, \ \frac{\partial f}{\partial z}$ Upstream **Downstream** Local Gradient Gradient Gradient Justin Johnson Lecture 6 - 99 Fall 2019

Backpropagation: x -2 Simple Example 3 f(x, y, z) = (x + y)zy <u>5</u> -12 e.g. x = -2, y = 5, z = -4 z <u>-4</u> 3 **1. Forward pass**: Compute outputs **Chain Rule** q = x + y | f = qz ∂f $\partial q \ \partial f$ $\partial x = \partial x \partial q$ ∂x 2. Backward pass: Compute derivatives Want: $\frac{\partial f}{\partial x}, \ \frac{\partial f}{\partial y}, \ \frac{\partial f}{\partial z}$ Upstream **Downstream** Local Gradient Gradient Gradient Justin Johnson Lecture 6 - 100 Fall 2019

Backpropagation: x -2 Simple Example 3 f(x, y, z) = (x + y)zy 5 -12 e.g. x = -2, y = 5, z = -4 z <u>-4</u> 3 **1. Forward pass**: Compute outputs **Chain Rule** q = x + y | f = qz ∂f $\partial q \ \partial f$ $\partial x = \partial x \partial q$ ∂x 2. Backward pass: Compute derivatives Want: $\frac{\partial f}{\partial x}$, $\frac{\partial f}{\partial y}$, $\frac{\partial f}{\partial z}$ Upstream **Downstream** Local Gradient Gradient Gradient

Justin John<u>son</u>

Lecture 6 - 101



Justin Johnson









Justin Johnson

Lecture 6 - 105



Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$











Lecture 6 - 109



Lecture 6 - 110



Justin Johnson

Lecture 6 - 111

Another Example f(x, w) = $\overline{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$ **Backward pass**: Compute gradients w0 2.00 -2.00 Local Gradient * x0 -1.00 $\cdot e^{x}$ ' $\overline{\partial}x$ 4.00 + w1 -3.00 6.00 * 1.00 -1.00 0.37 1.37 0.73 +1 *_1 x1 -2.00 1.00 -0.20 -0.53-0.53 w2 -3.00 Downstream Upstream Gradient Gradient





	U	S	ti	n		0	h	n	S	\mathbf{O}	n
-	<u> </u>	9	U		2				9	$\mathbf{\nabla}$	





Justin Johnson

Lecture 6 - 114










 $\frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$ Another Example f(x, w) = $= \sigma(w_0 x_0 + w_1 x_1 + w_2)$ **Backward pass**: Compute gradients w0 2.00 -0.20 -2.00 Computational graph is not * 0.20 x0 -1.00 unique: we can use primitives 0.39 that have simple local gradients 4.00 +0.20 w1 -3.00 Sigmoid -0.39 6.00 * 0.73 1.00 -1.00 0.37 1.37 0.20 *_1 x1 -2.00 -0.53 0.20 -0.20-0.531.00 -0.59w2 -3.00 0.20





Lecture 6 - 119

Another Example $f(x, w) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$ $= \sigma(w_0 x_0 + w_1 x_1 + w_2)$ **Backward pass**: Compute gradients w0 2.00 -0.20 -2.00 * Computational graph is not 0.20 $\frac{-}{1 \perp \rho^{-x}}$ $\sigma(x)$ x0 -1.00 unique: we can use primitives 0.39 that have simple local gradients 4.00 +0.20 w1 -3.00 Sigmoid -0.39 6.00 * 1.00 -1.00 0.37 1.37 0.20 x1 -2.00 -0.53 -0.20-0.530.20 -0.59[Downstream] = [Local] * [Upstream] w2 -3.00 0.20 = (1 - 0.73) * 0.73 * 1.0 = 0.2 $\frac{\partial}{\partial x} \left[\sigma(x) \right] = \frac{e^{-x}}{(1+e^{-x})^2} = \left(\frac{1+e^{-x}-1}{1+e^{-x}} \right) \left(\frac{1}{1+e^{-x}} \right) = (1-\sigma(x))\sigma(x)$ Sigmoid local gradient:

Justin Johnson



Gradients of other Activation Functions?



Leaky ReLU $\max(0.1x, x)$



 $\begin{array}{l} \textbf{Maxout} \\ \max(w_1^T x + b_1, w_2^T x + b_2) \end{array}$



Justin Johnson

Gradients of activation functions?

$$\begin{array}{ll} \text{Sigmoid local} & \frac{\partial}{\partial x} \Big[\sigma(x) \Big] = \frac{e^{-x}}{(1+e^{-x})^2} = \left(\frac{1+e^{-x}-1}{1+e^{-x}} \right) \left(\frac{1}{1+e^{-x}} \right) = (1-\sigma(x))\sigma(x)$$
gradient:

Name	Plot	Equation	Derivative
TanH		$f(x) = \tanh(x) = \frac{2}{1 + e^{-2x}} - 1$	$f'(x) = 1 - f(x)^2$
Rectified Linear Unit (ReLU)		$f(x) = \begin{cases} 0 & \text{for } x < 0 \\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} 0 & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$
Parameteric Rectified Linear Unit (PReLU) ^[2]		$f(x) = \begin{cases} \alpha x & \text{for } x < 0\\ x & \text{for } x \ge 0 \end{cases}$	$f'(x) = \begin{cases} \alpha & \text{for } x < 0\\ 1 & \text{for } x \ge 0 \end{cases}$

https://towardsdatascience.com/activation-functions-neural-networks-1cbd9f8d91d6

Cross-Entropy Loss (Multinomial Logistic Regression)

Want to interpret raw classifier scores as probabilities



Justin Johnson



Gradients of softmax?

Derivative of softmax

Let's compute $D_j S_i$ for arbitrary *i* and *j*:

$$D_j S_i = \frac{\partial S_i}{\partial a_j} = \frac{\partial \frac{e^{a_i}}{\sum_{k=1}^N e^{a_k}}}{\partial a_j}$$

We'll be using the quotient rule of derivatives. For $f(x) = \frac{g(x)}{h(x)}$:

$$f'(x) = \frac{g'(x)h(x) - h'(x)g(x)}{[h(x)]^2}$$

$$g_i = e^{a_i}$$
$$h_i = \sum_{k=1}^N e^{a_k}$$

$$D_j S_i = \begin{cases} S_i (1 - S_j) & i = j \\ -S_j S_i & i \neq j \end{cases}$$

es, but if anyone is taking more prid why you'll find various "condensed" s using the Kronecker delta function

$$\delta_{ij} = \begin{cases} 1 & i = j \\ 0 & i \neq j \end{cases}$$

$$D_j S_i = S_i (\delta_{ij} - S_j)$$

https://eli.thegreenplace.net/2016/the-softmax-function-and-its-derivative/

add gate: gradient distributor







add gate: gradient distributor



copy gate: gradient adder



Justin Johnson



add gate: gradient distributor



mul gate: "swap multiplier"



copy gate: gradient adder



Justin Johnson



add gate: gradient distributor



copy gate: gradient adder



mul gate: "swap multiplier"







Justin Johnson



Backprop Implementation: "Flat" gradient code: Forw

Forward pass: Compute output def f(w0, x0, w1, x1, w2):

s0 = w0 * x0
s1 = w1 * x1
s2 = s0 + s1
s3 = s2 + w2
L = sigmoid(s3)



Justin Johnson



Backprop Implementati "Flat" gradient code:	ON: Forward pass: Compute output	<pre>def f(w0, x0, w1, x1, w2): s0 = w0 * x0 s1 = w1 * x1 s2 = s0 + s1 s3 = s2 + w2</pre>
w0 2.00 -0.20 + -2.00		L = sigmoid(s3)
x0 -1.00 + 0.20 + 4.00 + 0.2	Backward pass: Compute grads	<pre>grad_L = 1.0 grad_s3 = grad_L * (1 - L) * L grad_w2 = grad_s3 grad_s2 = grad_s3 grad_s0 = grad_s2 grad_s1 = grad_s2 grad_w1 = grad_s1 * x1 grad_x1 = grad_s1 * w1 grad_w0 = grad_s0 * x0 grad_x0 = grad_s0 * w0</pre>

Lecture 6 - 130

		<pre>def f(w0, x0, w1, x1, w2):</pre>
Backprop Implementat	s0 = w0 * x0	
"Elat" gradiont codo:		s1 = w1 * x1
Flat glaulent coue.	Forward pass:	s2 = s0 + s1
	Compute output	s3 = s2 + w2
w0 2.00		L = sigmoid(s3)
-0.20		
$x_{0} -1.00$ 0.20	Base case	grad_L = 1.0
(1.40)		grad_s3 = grad_L * (1 - L) * L
w1 -3.00 0.20		grad_w2 = grad_s3
		grad_s2 = grad_s3
$x_1 -2.00$ 0.20 + 1.00 0.20		grad_s0 = grad_s2
-0.00		grad_s1 = grad_s2
w2 <u>-3.00</u>		grad_w1 = grad_s1 * x1
0.20		and the second of the second
		grad_x1 = grad_s1 * W1
		grad_x1 = grad_s1 * w1 grad_w0 = grad_s0 * x0
		grad_x1 = grad_s1 * w1 grad_w0 = grad_s0 * x0 grad_x0 = grad_s0 * w0

Lecture 6 - 131

Backprop Implementat "Flat" gradient code:	Tion: Forward pass: Compute output	def f(w0, x0, w1, x1, w2): s0 = w0 * x0 s1 = w1 * x1 s2 = s0 + s1 s3 = s2 + w2
w0 2.00 -0.20 x0 -1.00 * $-2.000.20$		L = sigmoid(s3) grad_L = 1.0
0.40	Sigmoid	grad_s3 = grad_L * (1 - L) * L
w1 -3.00 0.20		grad_w2 = grad_s3
	F	grad_s2 = grad_s3
$x_1 -2.00$ 0.20 + 0.20		grad_s0 = grad_s2
-0.00		grad_s1 = grad_s2
$w_{2} - 3.00$		grad_w1 = grad_s1 * x1
0.20		grad_x1 = grad_s1 * w1
		grad_w0 = grad_s0 * x0
		grad_x0 = grad_s0 * w0

Lecture 6 - 132



def f(w0, x0, w1, x1, w2): s0 = w0 * x0s1 = w1 * x1 s2 = s0 + s1s3 = s2 + w2L = sigmoid(s3)

grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0

Fall 2019

Justin Johnson

Lecture 6 - 133

Add



def f(w0, x0, w1, x1, w2): s0 = w0 * x0s1 = w1 * x1s2 = s0 + s1s3 = s2 + w2L = sigmoid(s3)

grad_L = 1.0
grad_s3 = grad_L * (1 - L) * L
grad_w2 = grad_s3
grad_s2 = grad_s3
grad_s0 = grad_s2
grad_s1 = grad_s2
grad_w1 = grad_s1 * x1
grad_x1 = grad_s1 * w1
grad_w0 = grad_s0 * x0
grad_x0 = grad_s0 * w0

Fall 2019

Justin Johnson

Lecture 6 - 134

Add



):
s0 = w0 * x0	
s1 = w1 * ×1	
s2 = s0 + s1	
s3 = s2 + w2	
L = sigmoid(s3)	

	grad_L = 1.0				
	grad_s3 = grad_L * (1 - L) * L				
	grad_w2 = grad_s3				
	grad_s2 = grad_s3				
	grad_s0 = grad_s2				
	grad_s1 = grad_s2				
Multiply	grad_w1 = grad_s1 * x1				
wuttpry	grad_x1 = grad_s1 * w1				
	grad_w0 = grad_s0 * x0				
	grad_x0 = grad_s0 * w0				



C	lef	<mark>f</mark> (۱	v0,	х¢	0,	w1,	x1	, w2)	:
	s) =	w0	*	X	0			
	s:	L =	w1	*	X	1			
	sź	2 =	s0	+	S	1			
	s	3 =	s2	+	w	2			
	L	= 9	sigr	no:	id	(s3)			

	grad_L = 1.0
	grad_s3 = grad_L * (1 - L) * L
	grad_w2 = grad_s3
	grad_s2 = grad_s3
	grad_s0 = grad_s2
	grad_s1 = grad_s2
	grad_w1 = grad_s1 * x1
	grad_x1 = grad_s1 * w1
,	grad_w0 = grad_s0 * x0
	grad_x0 = grad_s0 * w0

Backprop Implementation: Modular API



Graph (or Net) object (rough pseudo code)



Justin Johnson



Example: PyTorch Autograd Functions



(x,y,z are scalars)

	_
<pre>class Multiply(torch.autograd.Function):</pre>	
@staticmethod	
<pre>def forward(ctx, x, y):</pre>	Need to stash some
ctx.save_for_backward(x, y)	values for use in
z = x * y	backward
return z	
@staticmethod	
<pre>def backward(ctx, grad_z):</pre>	Upstream gradient
<pre>x, y = ctx.saved_tensors</pre>	
grad_x = y * grad_z # dz/dx * dL/dz	Multiply upstream
grad_y = x * grad_z # dz/dy * dL/dz	and local gradients
<mark>return</mark> grad_x, grad_y	

Justin Johnson



Example: PyTorch operators

pytorch / pytorch		⊙ Watch •	1,221	★ Un	star 26,770	¥ Fork	6,340
↔ Code ① Issues 2,286 1	Pull requests 561 💷 Projects 4 💷 V	Viki 🔟 In	sights				
Tree: 517c7c9861 - pytorch / aten	/ src / THNN / generic /		Create n	ew file	Upload files	Find file	History
ezyang and facebook-github-bot Ca	anonicalize all includes in PyTorch. (#14849)			Late	est commit 517	c7c9 on Dec	8, 2018
AbsCriterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
BCECriterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
ClassNLLCriterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
Col2Im.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
ELU.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
FeatureLPPooling.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
GatedLinearUnit.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
HardTanh.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
Im2Col.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
IndexLinear.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
E LeakyReLU.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
LogSigmoid.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
MSECriterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
MultiLabelMarginCriterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
MultiMarginCriterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
RReLU.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago
Sigmoid.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SmoothL1Criterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SoftMarginCriterion.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SoftPlus.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SoftShrink.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SparseLinear.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SpatialAdaptiveAveragePooling.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SpatialAdaptiveMaxPooling.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	nths ago
SpatialAveragePooling.c	Canonicalize all includes in PyTorch. (#1484	9)				4 mor	ths ago

SpatialClassNLLCriterion.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialConvolutionMM.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialDilatedMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialFractionalMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialFullDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialMaxUnpooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialReflectionPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialReplicationPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialUpSamplingBilinear.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
SpatialUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
THNN.h	Canonicalize all includes in PyTorch. (#14849)	4 months ago
Tanh.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalReflectionPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalReplicationPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalRowConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalUpSamplingLinear.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
TemporalUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricAdaptiveAveragePoolin	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricAdaptiveMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricAveragePooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricConvolutionMM.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricDilatedMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricFractionalMaxPooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricFullDilatedConvolution.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricMaxUnpooling.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricReplicationPadding.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricUpSamplingNearest.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
VolumetricUpSamplingTrilinear.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago
linear_upsampling.h	Implement nn.functional.interpolate based on upsample. (#8591)	9 months ago
pooling_shape.h	Use integer math to compute output size of pooling operations (#14405)	4 months ago
infold.c	Canonicalize all includes in PyTorch. (#14849)	4 months ago

Justin Johnson



```
#ifndef TH_GENERIC_FILE
    #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
 2
    #else
 3
 4
 5
     void THNN_(Sigmoid_updateOutput)(
 6
               THNNState *state,
               THTensor *input,
               THTensor *output)
 8
 9
     {
       THTensor_(sigmoid)(output, input);
10
11
12
     void THNN_(Sigmoid_updateGradInput)(
13
14
               THNNState *state,
               THTensor *gradOutput,
15
               THTensor *gradInput,
16
               THTensor *output)
17
18
     {
       THNN_CHECK_NELEMENT(output, gradOutput);
19
       THTensor_(resizeAs)(gradInput, output);
20
       TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
21
         scalar_t z = *output_data;
22
         *gradInput_data = *gradOutput_data * (1. - z) * z;
23
       );
24
25
    }
26
    #endif
27
```

PyTorch sigmoid layer

Source

Fall 2019

Justin Johnson

```
#ifndef TH GENERIC FILE
                                                                                       PyTorch sigmoid layer
    #define TH_GENERIC_FILE "THNN/generic/Sigmoid.c"
 2
    #else
 3
 4
     void THNN_(Sigmoid_updateOutput)(
 5
                                                         Forward
              THNNState *state,
 6
              THTensor *input,
                                                   \sigma(x) =
              THTensor *output)
 8
 9
      THTensor_(sigmoid)(output, input);
10
11
12
    void THNN_(Sigmoid_updateGradInput)(
13
14
              THNNState *state,
              THTensor *gradOutput,
15
              THTensor *gradInput,
16
              THTensor *output)
17
18
19
      THNN_CHECK_NELEMENT(output, gradOutput);
      THTensor_(resizeAs)(gradInput, output);
20
21
      TH_TENSOR_APPLY3(scalar_t, gradInput, scalar_t, gradOutput, scalar_t, output,
        scalar_t z = *output_data;
22
        *gradInput_data = *gradOutput_data * (1. - z) * z;
23
      );
24
25
    }
26
    #endif
27
                                                                                                                        Fall 2019
          Justin Johnson
                                                             Lecture 6 - 143
```

Source





Source

So far: backprop with scalars

What about vector-valued functions?

Justin Johnson



Recap: Vector Derivatives

 $x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:

 $\frac{\partial y}{\partial x} \in \mathbb{R}$

If x changes by a small amount, how much will y change?

Justin Johnson



Recap: Vector Derivatives

$$x \in \mathbb{R}, y \in \mathbb{R}$$

Regular derivative:

$$x \in \mathbb{R}^N, y \in \mathbb{R}$$

Derivative is Gradient:

$$\frac{\partial y}{\partial x} \in \mathbb{R}$$

$$\frac{\partial y}{\partial x} \in \mathbb{R}^N \quad \left(\frac{\partial y}{\partial x}\right)_n = \frac{\partial y}{\partial x_n}$$

If x changes by a small amount, how much will y change? For each element of x, if it changes by a small amount then how much will y change?



Recap: Vector Derivatives

 $x \in \mathbb{R}, y \in \mathbb{R}$

Regular derivative:



 $x \in \mathbb{R}^N, y \in \mathbb{R}$

Derivative is **Gradient**:

$$x \in \mathbb{R}^N, y \in \mathbb{R}^M$$

Derivative is **Jacobian**:



If x changes by a small amount, how much will y change?

For each element of x, if it changes by a small amount then how much will y change?

For each element of x, if it changes by a small amount then how much will each element of y change?

Justin Johnson









Fall 2019

Justin Johnson

Backprop with Vectors

Justin Johnson



Lecture 6 - 152
Backprop with Vectors



Justin Johnson

Lecture 6 - 153

Fall 2019

Backprop with Vectors





Backprop with Vectors



Justin Johnson



Backprop with Vectors



Backprop with Vectors



Justin Johnson

Backprop with Vectors



Backprop with Vectors

Jacobian is **sparse**: off-diagonal entries all zero! Never **explicitly** form Jacobian; instead use **implicit** multiplication



Backprop with Vectors

Jacobian is **sparse**: off-diagonal entries all zero! Never **explicitly** form Jacobian; instead use **implicit** multiplication



Backpropagation: Another View



 $\frac{\partial L}{\partial x_0} = \left(\frac{\partial x_1}{\partial x_0}\right) \left(\frac{\partial x_2}{\partial x_1}\right) \left(\frac{\partial x_3}{\partial x_2}\right) \left(\frac{\partial L}{\partial x_2}\right)$ Chain rule



Backpropagation: Another View



Matrix multiplication is associative: we can compute products in any order

$$\begin{array}{l} {}^{\text{Chain}}_{\text{rule}} & \frac{\partial L}{\partial x_0} = \left(\frac{\partial x_1}{\partial x_0}\right) \left(\frac{\partial x_2}{\partial x_1}\right) \left(\frac{\partial x_3}{\partial x_2}\right) \left(\frac{\partial L}{\partial x_3}\right) \\ & {}^{\text{D}}_{0} \, \text{x} \, \text{D}_{1} & {}^{\text{D}}_{1} \, \text{x} \, \text{D}_{2} & {}^{\text{D}}_{2} \, \text{x} \, \text{D}_{3} & {}^{\text{D}}_{3} \end{array}$$

Reverse-Mode Automatic Differentiation



Matrix multiplication is **associative**: we can compute products in any order Computing products right-to-left avoids matrix-matrix products; only needs matrix-vector

$$\begin{array}{ll} & \underset{\text{rule}}{\text{Chain}} & \frac{\partial L}{\partial x_0} = \left(\frac{\partial x_1}{\partial x_0}\right) \left(\frac{\partial x_2}{\partial x_1}\right) \left(\frac{\partial x_3}{\partial x_2}\right) \left(\frac{\partial L}{\partial x_3}\right) \\ & \quad \mathsf{D}_0 \, \mathsf{x} \, \mathsf{D}_1 \quad \mathsf{D}_1 \, \mathsf{x} \, \mathsf{D}_2 \quad \mathsf{D}_2 \, \mathsf{x} \, \mathsf{D}_3 \quad \mathsf{D}_3 \end{array}$$

Reverse-Mode Automatic Differentiation



Matrix multiplication is **associative**: we can compute products in any order Computing products right-to-left avoids matrix-matrix products; only needs matrix-vector

$$\begin{array}{l} \begin{array}{l} \text{Chain} \\ \text{rule} \end{array} & \frac{\partial L}{\partial x_0} = \left(\frac{\partial x_1}{\partial x_0}\right) \left(\frac{\partial x_2}{\partial x_1}\right) \left(\frac{\partial x_3}{\partial x_2}\right) \left(\frac{\partial L}{\partial x_3}\right) \\ \end{array} \\ \begin{array}{l} \text{Compute grad of scalar output} \\ \text{w/respect to all vector inputs} \end{array} & \textbf{D}_0 \times \textbf{D}_1 \quad \textbf{D}_1 \times \textbf{D}_2 \quad \textbf{D}_2 \times \textbf{D}_3 \quad \textbf{D}_3 \end{array}$$

Justin Johnson

Summary

Represent complex expressions as **computational graphs**



Forward pass computes outputs

Backward pass computes gradients

During the backward pass, each node in the graph receives **upstream gradients** and multiplies them by **local gradients** to compute **downstream gradients**



Justin Johnson



Summary

Backprop can be implemented with "flat" code where the backward pass looks like forward pass reversed (Use this for A2!)

```
def f(w0, x0, w1, x1, w2):
 s0 = w0 * x0
 s1 = w1 * x1
 s2 = s0 + s1
 s3 = s2 + w2
 L = sigmoid(s3)
 grad_L = 1.0
 grad_s3 = grad_L * (1 - L) * L
 grad_w2 = grad_s3
 grad_s2 = grad_s3
 grad_s0 = grad_s2
 grad_s1 = grad_s2
 grad_w1 = grad_s1 * x1
 grad_x1 = grad_s1 * w1
 grad_w0 = grad_s0 * x0
  grad_x0 = grad_s0 * w0
```

Backprop can be implemented with a modular API, as a set of paired forward/backward functions (We will do this on A3!)

<pre>class Multiply(torch.autograd.Function):</pre>
@staticmethod
<pre>def forward(ctx, x, y):</pre>
<pre>ctx.save_for_backward(x, y)</pre>
z = x * y
return z
@staticmethod
<pre>def backward(ctx, grad_z):</pre>
<pre>x, y = ctx.saved_tensors</pre>
grad_x = y * grad_z # dz/dx * dL/dz
grad_y = x * grad_z # dz/dy * dL/dz
<mark>return</mark> grad_x, grad_y

Justin Johnson

Deep neural networks

- Loss functions
- Regularization
- Weights and layers
- Activation functions
- Backpropagation
- Training and optimization





EECS 498-007 / 598-005 Deep Learning for Computer Vision Fall 2019

Lecture 4: Optimization

Justin Johnson



Loss Functions quantify preferences

- We have some dataset of (x, y)
- We have a **score function**:
- We have a **loss function**:

Q: How do we find the best W?

$$s = f(x; W) = Wx$$
Linear classifier



Justin Johnson

Follow the slope

In 1-dimension, the **derivative** of a function gives the slope:

$$rac{df(x)}{dx} = \lim_{h o 0} rac{f(x+h) - f(x)}{h}$$

In multiple dimensions, the **gradient** is the vector of (partial derivatives) along each dimension

The slope in any direction is the **dot product** of the direction with the gradient The direction of steepest descent is the **negative gradient**



Loss is a function of W: Analytic Gradient

$$egin{aligned} L &= rac{1}{N} \sum_{i=1}^N L_i + \sum_k W_k^2 \ L_i &= \sum_{j
eq y_i} \max(0, s_j - s_{y_i} + 1) \ s &= f(x; W) = Wx \end{aligned}$$

want $\nabla_W L$

Use calculus to compute an analytic gradient



This image is in the public domain

This image is in the public domain

Justin Johnson



Numeric gradient: approximate, slow, easy to write Analytic gradient: exact, fast, error-prone

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a **gradient check**.

-

-



Numeric gradient: approximate, slow, easy to write Analytic gradient: exact, fast, error-prone

In practice: Always use analytic gradient, but check implementation with numerical gradient. This is called a gradient check.

```
def grad_check_sparse(f, x, analytic_grad, num_checks=10, h=1e-7):
    """
    sample a few random elements and only return numerical
    in this dimensions.
    """
```

Justin Johnson

-

-



Numeric gradient: approximate, slow, easy to write Analytic gradient: exact, fast, error-prone

torch.autograd.gradcheck(func, inputs, eps=1e-06, atol=1e-05, rtol=0.001, raise_exception=True, check_sparse_nnz=False, nondet_tol=0.0)

[SOURCE] SOURCE]

Check gradients computed via small finite differences against analytical gradients w.r.t. tensors in inputs that are of floating point type and with requires_grad=True.

The check between numerical and analytical gradients uses allclose().

-



Numeric gradient: approximate, slow, easy to write Analytic gradient: exact, fast, error-prone

torch.autograd.gradgradcheck(func, inputs, grad_outputs=None, eps=1e-06, atol=1e-05, rtol=0.001, gen_non_contig_grad_outputs=False, raise_exception=True, [SOURCE] nondet_tol=0.0)

Check gradients of gradients computed via small finite differences against analytical gradients w.r.t. tensors in inputs and grad_outputs that are of floating point type and with requires_grad=True.

This function checks that backpropagating through the gradients computed to the given grad_outputs are correct.

Justin Johnson

-

Gradient Descent

Iteratively step in the direction of the negative gradient (direction of local steepest descent)

Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
 dw = compute_gradient(loss_fn, data, w)
 w -= learning_rate * dw

Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate



Gradient Descent

Iteratively step in the direction of the negative gradient (direction of local steepest descent)

Vanilla gradient descent
w = initialize_weights()
for t in range(num_steps):
 dw = compute_gradient(loss_fn, data, w)
 w -= learning_rate * dw

Hyperparameters:

- Weight initialization method
- Number of steps
- Learning rate





Justin Johnson

Batch Gradient Descent

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

$$7_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

Full sum expensive when N is large!

Justin Johnson



Stochastic Gradient Descent (SGD)

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

$$7_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W) + \lambda \nabla_W R(W)$$

for t in range(num_steps):

minibatch = sample_data(data, batch_size)

w -= learning_rate * dw

Full sum expensive when N is large!

Approximate sum using a **minibatch** of examples 32 / 64 / 128 common

Hyperparameters:

- Weight initialization
- Number of steps
- Learning rate
- Batch size
- Data sampling

Justin Johnson



Stochastic Gradient Descent (SGD)

$$L(W) = \mathbb{E}_{(x,y) \sim p_{data}} \left[L(x, y, W) \right] + \lambda R(W)$$

 $\approx \frac{1}{N} \sum_{i=1}^{N} L(x_i, y_i, W) + \lambda R(W)$

Think of loss as an expectation over the full **data distribution** p_{data}

Approximate expectation via sampling



Stochastic Gradient Descent (SGD)

$$L(W) = \mathbb{E}_{(x,y) \sim p_{data}} \left[L(x, y, W) \right] + \lambda R(W)$$

 $\approx \frac{1}{N} \sum_{i=1}^{N} L(x_i, y_i, W) + \lambda R(W)$

Think of loss as an expectation over the full data distribution p_{data}

$$\nabla_W L(W) = \nabla_W \mathbb{E}_{(x,y) \sim p_{data}} \left[L(x, y, W) \right] + \lambda \nabla_W R(W))$$
$$\approx \sum_{i=1}^N \nabla_W L_W(x_i, y_i, W) + \nabla_W R(W)$$

Justin Johnson



Recall: Reverse-Mode Automatic Differentiation



Matrix multiplication is **associative**: we can compute products in any order Computing products right-to-left avoids matrix-matrix products; only needs matrix-vector

$$\begin{array}{l} \begin{array}{l} \text{Chain} \\ \text{rule} \end{array} & \frac{\partial L}{\partial x_0} = \left(\frac{\partial x_1}{\partial x_0}\right) \left(\frac{\partial x_2}{\partial x_1}\right) \left(\frac{\partial x_3}{\partial x_2}\right) \left(\frac{\partial L}{\partial x_3}\right) \\ \end{array} \\ \begin{array}{l} \text{Compute grad of scalar output} \\ \text{w/respect to all vector inputs} \end{array} & \textbf{D}_0 \times \textbf{D}_1 \quad \textbf{D}_1 \times \textbf{D}_2 \quad \textbf{D}_2 \times \textbf{D}_3 \quad \textbf{D}_3 \end{array}$$

Justin Johnson

Mini-batch evaluation with matrices (HW3)

- DNNs are described as passing vectors between layers
- Why not pass all samples in a mini-batch as a *matrix*?
- What used to be column vectors are now rows
- Need to adjust weight-vector multiplies

 s = W x
 becomes
 S = X W^T
- Need to adjust gradients (Jacobians) as well
- Fix this description after doing homework...





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

Justin Johnson



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

Justin Johnson



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



Justin Johnson



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

Zero gradient, gradient descent gets stuck







Our gradients come from minibatches so they can be noisy!

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

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



Justin Johnson


SGD

SGD $x_{t+1} = x_t - \alpha \nabla f(x_t)$

for t in range(num_steps):
 dw = compute_gradient(w)
 w -= learning_rate * dw

Justin Johnson



SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

for t in range(num_steps):
 dw = compute_gradient(w)
 w -= learning_rate * dw

SGD+Momentum $v_{t+1} = \rho v_t + \nabla f(x_t)$ $x_{t+1} = x_t - \alpha v_{t+1}$ v = 0

for t in range(num_steps):
 dw = compute_gradient(w)
 v = rho * v + dw
 w -= learning_rate * v

- Build up "velocity" as a running mean of gradients
 - Rho gives "friction"; typically rho=0.9 or 0.99

Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013

Justin Johnson



SGD+Momentum

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t)$$
$$x_{t+1} = x_t + v_{t+1}$$

v = 0

```
for t in range(num_steps):
    dw = compute_gradient(w)
    v = rho * v - learning_rate * dw
    w += v
```

SGD+Momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$

$$x_{t+1} = x_t - \alpha v_{t+1}$$

v = 0
for t in range(num_steps):
 dw = compute_gradient(w)
 v = rho * v + dw
 w -= learning_rate * v

You may see SGD+Momentum formulated different ways, but they are equivalent - give same sequence of x

Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013



Local Minima Saddle points **Poor Conditioning**

Gradient Noise



SGD

Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013

Justin Johnson

Lecture 4 - 227



SGD+Momentum

Momentum update:



Combine gradient at current point with velocity to get step used to update weights

Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k^2)", 1983 Nesterov, "Introductory lectures on convex optimization: a basic course", 2004 Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013

Justin Johnson



Momentum update:



Gradient

Nesterov Momentum



Combine gradient at current point with velocity to get step used to update weights

Nesterov, "A method of solving a convex programming problem with convergence rate O(1/k^2)", 1983 Nesterov, "Introductory lectures on convex optimization: a basic course", 2004 Sutskever et al, "On the importance of initialization and momentum in deep learning", ICML 2013 "Look ahead" to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Justin Johnson

Lecture 4 - 229

Fall 2019

$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$
$$x_{t+1} = x_t + v_{t+1}$$



"Look ahead" to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Justin Johnson



$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$
$$x_{t+1} = x_t + v_{t+1}$$

Annoying, usually we want update in terms of $x_t,
abla f(x_t)$



"Look ahead" to the point where updating using velocity would take us; compute gradient there and mix it with velocity to get actual update direction

Justin Johnson



$$v_{t+1} = \rho v_t - \alpha \nabla f(x_t + \rho v_t)$$
$$x_{t+1} = x_t + v_{t+1}$$

Annoying, usually we want update in terms of $x_t, \nabla f(x_t)$

Change of variables $\tilde{x}_t = x_t + \rho v_t$ and rearrange:

$$v_{t+1} = \rho v_t - \alpha \nabla f(\tilde{x}_t)$$

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

$$= \tilde{x}_t + v_{t+1} + \rho(v_{t+1} - v_t)$$

v = 0
for t in range(num_steps):
 dw = compute_gradient(w)
 old_v = v
 v = rho * v - learning_rate * dw
 w -= rho * old_v - (1 + rho) * v

Justin Johnson



Justin Johnson





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

"Per-parameter learning rates" or "adaptive learning rates"

Duchi et al, "Adaptive subgradient methods for online learning and stochastic optimization", JMLR 2011







Duchi et al, "Adaptive subgradient methods for online learning and stochastic optimization", JMLR 2011

Justin Johnson







Q: What happens with AdaGrad?

Justin Johnson







Q: What happens with AdaGrad?

Progress along "steep" directions is damped; progress along "flat" directions is accelerated

Justin Johnson



RMSProp: "Leaky Adagrad"



Tieleman and Hinton, 2012

Justin Johnson



RMSProp



Justin Johnson



```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
```



moment1 = 0							
moment2 = 0							
<pre>for t in range(num_steps):</pre>							
_dw = compute_gradient(w)							
moment1 = beta1 * moment1 + (1 - beta1) * dw							
momen <u>t2 = beta2 * moment2 + (1 –</u> beta2) * dw * dw							
<pre>w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)</pre>							

Momentum

SGD+Momentum

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



moment1 = 0	
moment2 = 0	Adam
<pre>for t in range(num_steps):</pre>	
<u>dw = compute_gradient(w)</u>	Momentum
<pre>moment1 = beta1 * moment1 + (1 - beta1) * dw</pre>	
moment2 = beta2 * moment2 + (1 - beta2) * dw * dw	AdaGrad / RIVISProp
<pre>w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)</pre>	



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

Justin Johnson



```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    w -= learning_rate * moment1 / (moment2.sqrt() + 1e-7)
Bias correction
```

Q: What happens at t=0? (Assume beta2 = 0.999)







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

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

Justin Johnson



```
moment1 = 0
moment2 = 0
for t in range(num_steps):
    dw = compute_gradient(w)
    moment1 = beta1 * moment1 + (1 - beta1) * dw
    moment2 = beta2 * moment2 + (1 - beta2) * dw * dw
    moment1_unbias = moment1 / (1 - beta1 ** t)
    moment2_unbias = moment2 / (1 - beta2 ** t)
    w -= learning_rate * moment1_unbias / (moment2_unbias.sqrt() + 1e-7)
```

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

Adam with beta1 = 0.9, beta2 = 0.999, and learning_rate = 1e-3, 5e-4, 1e-4 is a great starting point for many models!

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

Justin Johnson



Adam: Very Common in Practice!

for input to the CNN; each colored pixel in the image yields a 7D one-hot vector. Following common practice, the network is trained end-to-end using stochastic gradient descent with the Adam optimizer [22]. We anneal the learning rate to 0 using a half cosine schedule without restarts [28].

Bakhtin, van der Maaten, Johnson, Gustafson, and Girshick, NeurIPS 2019

We train all models using Adam [23] with learning rate 10^{-4} and batch size 32 for 1 million iterations; training takes about 3 days on a single Tesla P100. For each minibatch we first update f, then update D_{img} and D_{obj} .

Johnson, Gupta, and Fei-Fei, CVPR 2018

ganized into three residual blocks. We train for 25 epochs using Adam [27] with learning rate 10^{-4} and 32 images per batch on 8 Tesla V100 GPUs. We set the cubify thresh-

Gkioxari, Malik, and Johnson, ICCV 2019

sampled with each bit drawn uniformly at random. For gradient descent, we use Adam [29] with a learning rate of 10^{-3} and default hyperparameters. All models are trained with batch size 12. Models are trained for 200 epochs, or 400 epochs if being trained on multiple noise layers.

Zhu, Kaplan, Johnson, and Fei-Fei, ECCV 2018

16 dimensional vectors. We iteratively train the Generator and Discriminator with a batch size of 64 for 200 epochs using Adam [22] with an initial learning rate of 0.001.

Gupta, Johnson, et al, CVPR 2018

Adam with beta1 = 0.9, beta2 = 0.999, and learning_rate = 1e-3, 5e-4, 1e-4 is a great starting point for many models!

Justin Johnson



Adam



Justin Johnson



Optimization Algorithm Comparison

Algorithm	Tracks first moments (Momentum)	Tracks second moments (Adaptive learning rates)	Leaky second moments	Bias correction for moment estimates
SGD	X	X	X	X
SGD+Momentum	\checkmark	X	X	X
Nesterov	\checkmark	X	X	X
AdaGrad	X	\checkmark	X	X
RMSProp	X	\checkmark	\checkmark	X
Adam	\checkmark	\checkmark	\checkmark	\checkmark



_

- Adam is a good default choice in many cases SGD+Momentum can outperform Adam but may require more tuning
- If you can afford to do full batch updates then try out
 L-BFGS (and don't forget to disable all sources of noise)



Deep neural networks

- Loss functions
- Regularization
- Weights and layers
- Activation functions
- Backpropagation
- Training and optimization

