

CS 343: Artificial Intelligence

Deep Learning

Profs. Peter Stone and Yuke Zhu — The University of Texas at Austin

[These slides based on those of Dan Klein, Pieter Abbeel, Anca Dragan for CS188 Intro to AI at UC Berkeley. All CS188 materials are available at http://ai.berkeley.edu.]

Good morning colleagues!

- Past due:
 - HW1-7: Search, CSPs, Games, MDP, RL, Bayes Nets, Particle Filters/VPI
 - 11 reading responses: AI100 report; 10 Textbook readings
 - P0,1,2,3: tutorial, Search, Multiagent RL
 - Midterm
- Upcoming EdX Homeworks
 - HW8: Naive Bayes and Perceptrons due Monday 4/19 at 11:59 pm
 - HW9: Neural Networks due Monday 4/26 at 11:59 pm
- Upcoming programming projects
 - P4: Bayes Nets due Wednesday 4/14 at 11:59pm
 - P5: Particle Filters due Wednesday 4/21 at 11:59pm
 - Contest (Capture the flag): Qualification due 4/28; Finals 5/3 (extra credit)
- Readings: SVMs, Kernels, and Clustering Due Monday 4/19 at 9:30am

Good morning colleagues!

Midterm grades

- Some context:
 - Deep learning = neural networks
 - Al <> Deep Learning
 - But...it's definitely an important area to know about these days
 - Applications other than vision and natural language processing?
 - Robotics
 - Fraud detection
 - Game playing (e.g. go, Starcraft)
 - Election predictions

Perceptron



Two-Layer Perceptron Network



N-Layer Perceptron Network



N-Layer Neural Network



- What's the relationship between a weight and a gradient?
 - How does gradient ascent work?
 - [Policy Grad RL slides]
 - [Gradient ascent problem]
 - Can gradient ascent be changed to find a global maximum? (Michael Labarca)
- What's the purpose of the activation function?
- What exactly is backpropagation?
 - [Gradient computation problem]
- Are there other ways to train NNs? (Tyler Miller)
 - NEAT

Test Your Understanding

- Data sufficiency problem
- Practice problem in breakout rooms
- Talk about each subproblem individually

- What's the purpose of the activation function?
- Why are leaky RELU units better than sigmoids or tanh? (Cyrus Mahdavi)
- [Representation capacity problem]

What Can be Done with Non-Linear (e.g., Threshold) Units?



separating hyperplane



convex polygon region





composition of polygons: non convex regions



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -16

20 Jan 2016

- What are hyperparameters and how do you tune them?
 - Activation units
 - Learning rate
 - Momentum parameter
 - Dropout parameters
 - Normalization scheme
 - Number of layers and units (architecture)
- How do you design architectures?
 - Neural architecture search [NEAT+Q slides]
 - Can hyperparameters be made learned parameters? (Conrad Li)
 - What's the difference between adding layers vs. widening a layer? (Michael Labarca)
 - Does batch normalization correct for bad initialization? (Yuhan Zheng)

- Differences between brain's NN and artificial NNs design? (Pranooha Veeramachaneni)
 - Are NNs designed based on intuition from brain structures? (Rudraksh Garg)
- What are the limitations of NNs? (Jack Si)
 - How far can we go with deeper and larger networks? (Trong Lv)
- When shouldn't you use NNs? (Ethan Houston)

- Computation requirements of NNs vs. traditional vision classification (Vijay Vuyyuru)
 - Training vs. testing
 - Depends on hardware
 - Why does deep learning work so much better on GPUs? (Cameron Doggett)
- Why manual features favored in past, and only recently NNs favored? (Nalin Mahajan)
- If image recognition is so good, why do some websites still require you to identify images to check if you're a robot? (Jessica Ma)

- Are there good ways of introducing human knowledge into NNs? Or is that missing the point? (Tyler Miller)
 - Neurosymbolic systems

Review: Linear Classifiers



Feature Vectors





Some (Simplified) Biology

Very loose inspiration: human neurons



Linear Classifiers

- Inputs are feature values
- Each feature has a weight
- Sum is the activation



activation_w(x) =
$$\sum_{i} w_i \cdot f_i(x) = w \cdot f(x)$$

If the activation is:

Positive, output +1



Non-Linearity



Non-Linear Separators

Data that is linearly separable works out great for linear decision rules:



But what are we going to do if the dataset is just too hard?



This and next slide adapted from Ray Mooney, UT

Non-Linear Separators

General idea: the original feature space can always be mapped to some higherdimensional feature space where the training set is separable:



Computer Vision



Object Detection



Manual Feature Design







Features and Generalization



[Dalal and Triggs, 2005]

Features and Generalization





Image



Manual Feature Design Deep Learning



- Manual feature design requires:
 - Domain-specific expertise
 - Domain-specific effort

• What if we could learn the features, too?

Perceptron



Two-Layer Perceptron Network



N-Layer Perceptron Network



Performance

ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

Performance

ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai
Performance

ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

Performance

ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

Performance

ImageNet Error Rate 2010-2014



graph credit Matt Zeiler, Clarifai

Speech Recognition



N-Layer Perceptron Network



Local Search

Simple, general idea:

- Start wherever
- Repeat: move to the best neighboring state
- If no neighbors better than current, quit
- Neighbors = small perturbations of w

Properties

Plateaus and local optima



How to escape plateaus and find a good local optimum? How to deal with very large parameter vectors? E.g., $w \in \mathbb{R}^{1billion}$

Perceptron



Objective: Classification Accuracy

$$l^{\text{acc}}(w) = \frac{1}{m} \sum_{i=1}^{m} \left(\text{sign}(w^{\top} f(x^{(i)})) = = y^{(i)} \right)$$

■ Issue: many plateaus _ how to measure incremental progress toward a correct label?

Soft-Max

Score for y=1:
$$w^{\top} f(x)$$
 for y=-1: $-w^{\top} f(x)$

m

Probability of label:

$$p(y = 1|f(x); w) = \frac{e^{w^{\top}f(x^{(i)})}}{e^{w^{\top}f(x)} + e^{-w^{\top}f(x)}}$$
$$p(y = -1|f(x); w) = \frac{e^{-w^{\top}f(x)}}{e^{w^{\top}f(x)} + e^{-w^{\top}f(x)}}$$

Objective:

Log:

$$l(w) = \prod_{i=1}^{m} p(y = y^{(i)} | f(x^{(i)}); w)$$
$$ll(w) = \sum_{i=1}^{m} \log p(y = y^{(i)} | f(x^{(i)}); w)$$

Two-Layer Neural Network



N-Layer Neural Network



Our Status

• Our objective ll(w)

- Changes smoothly with changes in w
- Doesn't suffer from the same plateaus as the perceptron network

Challenge: how to find a good w?

$$\max_{w} ll(w)$$

Equivalently:

$$\min_{w} -ll(w)$$

1-d optimization



Then step in best direction

• Or, evaluate derivative: $\frac{\partial g(w_0)}{\partial w} = \lim_{h \to 0} \frac{g(w_0 + h) - g(w_0 - h)}{2h}$

Tells which direction to step in

2-D Optimization



Source: Thomas Jungblut's Blog

Steepest Descent

Idea:

- Start somewhere
- Repeat: Take a step in the steepest descent direction



Figure source: Mathworks

What is the Steepest Descent Direction?

What is the Steepest Descent Direction?

Steepest Direction = direction of the gradient

$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \cdots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$

Optimization Procedure 1: Gradient Descent

Init:
$$w$$

For i = 1, 2, ...
 $w \leftarrow w - \alpha * \nabla g(w)$

- How? Try multiple choices
 - \blacksquare Crude rule of thumb: update changes $\,w\,$ about 0.1 1 %



Q: What is the trajectory along which we converge towards the minimum with Gradient Descent?

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -



Q: What is the trajectory along which we converge towards the minimum with Gradient Descent?

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -



Q: What is the trajectory along which we converge towards the minimum with Gradient Descent? very slow progress along flat direction, jitter along steep one

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -

Optimization Procedure 2: Momentum

• Gradient Descent• Momentum• Init: w• Init: w• For i = 1, 2, ...• For i = 1, 2, ... $w \leftarrow w - \alpha * \nabla g(w)$ $v \leftarrow \mu * v - \alpha * \nabla g(w)$ $w \leftarrow w + v$

- Physical interpretation as ball rolling down the loss function + friction (mu coefficient).

- mu = usually ~0.5, 0.9, or 0.99 (Sometimes annealed over time, e.g. from 0.5 -> 0.99)



Q: What is the trajectory along which we converge towards the minimum with Momentum?

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -

How do we actually compute gradient w.r.t. weights?

Backpropagation!

Backpropagation Learning

15-486/782: Artificial Neural Networks David S. Touretzky

Fall 2006

LMS / Widrow-Hoff Rule



Works fine for a single layer of trainable weights. What about multi-layer networks? With Linear Units, Multiple Layers Don't Add Anything



Linear operators are closed under composition. Equivalent to a single layer of weights $W=U\times V$

But with non-linear units, extra layers add computational power.

What Can be Done with Non-Linear (e.g., Threshold) Units?



separating hyperplane



convex polygon region





composition of polygons: non convex regions

How Do We Train A Multi-Layer Network?



Can't use perceptron training algorithm because we don't know the 'correct' outputs for hidden units.

How Do We Train A Multi-Layer Network?

Define sum-squared error:

$$E = \frac{1}{2} \sum_{p} (d^{p} - y^{p})^{2}$$

Use gradient descent error minimization:

$$\Delta w_{ij} = -\eta \frac{\partial E}{\partial w_{ij}}$$

Works if the nonlinear transfer function is differentiable.

Deriving the LMS or "Delta" Rule As Gradient Descent Learning



How do we extend this to two layers?

Switch to Smooth <u>Nonlinear</u> Units

$$\mathbf{net}_j = \sum_i w_{ij} y_i$$

 $y_j = g(net_j)$ g must be differentiable

Common choices for g:

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

$$g'(x) = g(x) \cdot (1 - g(x))$$

$$g(x) = \tanh(x)$$

$$g'(x) = \frac{1}{\cosh^2(x)}$$



Gradient Descent with Nonlinear Units

$$x_{i} \xrightarrow{\mathbf{w}_{i}} \tanh(\Sigma w_{i} x_{i}) \xrightarrow{\mathbf{v}} y$$

$$y = g(net) = \tanh\left(\sum_{i} w_{i} x_{i}\right)$$

$$\frac{dE}{dy} = (y - d), \qquad \frac{dy}{dnet} = 1/\cosh^2(net), \qquad \frac{\partial net}{\partial w_i} = x_i$$

$$\frac{\partial E}{\partial w_i} = \frac{d E}{d y} \cdot \frac{d y}{d net} \cdot \frac{\partial net}{\partial w_i}$$
$$= (y-d)/\cosh^2 \left(\sum_i w_i x_i\right) \cdot x_i$$

Now We Can Use The Chain Rule



$$\frac{\partial E}{\partial y_k} = (y_k - d_k)$$

$$\delta_k = \frac{\partial E}{\partial net_k} = (y_k - d_k) \cdot g'(net_k)$$

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial w_{jk}} = \delta_k \cdot y_j$$

$$\frac{\partial E}{\partial y_j} = \sum_k \left(\frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial y_j} \right)$$

$$\delta_j = \frac{\partial E}{\partial net_j} = \frac{\partial E}{\partial y_j} \cdot g'(net_j)$$

$$\frac{\partial E}{\partial w_{ij}} = \delta_j \cdot y_i$$

Weight Updates

$$\frac{\partial E}{\partial w_{jk}} = \frac{\partial E}{\partial net_k} \cdot \frac{\partial net_k}{\partial w_{jk}} = \delta_k \cdot y_j$$
$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial net_j} \cdot \frac{\partial net_j}{\partial w_{ij}} = \delta_j \cdot y_i$$

$$\Delta w_{jk} = -\eta \cdot \frac{\partial E}{\partial w_{jk}} \qquad \Delta w_{ij} = -\eta \cdot \frac{\partial E}{\partial w_{ij}}$$
Classification

Retrieval



[Krizhevsky 2012]

25 Jan 2016

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -

Segmentation

Detection



[Faster R-CNN: Ren, He, Girshick, Sun 2015]

[Farabet et al., 2012]

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -



self-driving cars

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -



[Toshev, Szegedy 2014]



[Mnih 2013]

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -



着 M 只适纸店 BJ

[Ciresan et al. 2013]



25 Jan 2016

[Sermanet et al. 2011] [Ciresan et al.]

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -

Describes without errors



A person riding a motorcycle on a dirt road.



Describes with minor errors

Two dogs play in the grass.





A skateboarder does a trick on a ramp.



A dog is jumping to catch a frisbee.



A refrigerator filled with lots of food and drinks.



A yellow school bus parked in a parking lot.

Image Captioning



A group of young people playing a game of frisbee.



Two hockey players are fighting over the puck.



A little girl in a pink hat is blowing bubbles.



A herd of elephants walking across a dry grass field.



A close up of a cat laying on a couch.









reddit.com/r/deepdream

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 6 -

Remaining Pieces

- Optimizing machine learning objectives:
 - Stochastic Descent
 - Mini-batches
- Improving generalization
 - Drop-out
- Activation functions

- Initialization and batch normalization
- Computing the gradient $\nabla g(w)$
 - Backprop
 - Gradient checking

Mini-batches and Stochastic Gradient Descent

Typical objective:

$$ll(w) = \frac{1}{m} \sum_{i=1}^{m} \log p(y = y^{(i)} | f(x^{(i)}); w)$$

= average log-likelihood of label given input

$$\approx \frac{1}{k} \sum_{i=1}^{k} \log p(y = y^{(i)} | f(x^{(i)}); w)$$

= estimate based on mini-batch 1...k

- Mini-batch gradient descent: compute gradient on mini-batch (+ cycle over mini-batches: 1..k, k+1...2k, ...; make sure to randomize permutation of data!)
- Stochastic gradient descent: k = 1

Remaining Pieces

- Optimizing machine learning objectives:
 - Stochastic Descent
 - Mini-batches
- Improving generalization
 - Drop-out
- Activation functions

- Initialization and batch normalization
- Computing the gradient $\nabla g(w)$
 - Gradient checking
 - Backprop

Regularization: **Dropout**

"randomly set some neurons to zero in the forward pass"



(a) Standard Neural Net



(b) After applying dropout.

[Srivastava et al., 2014]

Waaaait a second...

How could this possibly be a good idea?



Waaaait a second...

How could this possibly be a good idea?



Forces the network to have a redundant representation.



Waaaait a second...

How could 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, gets trained on only ~one datapoint.



Ideally:

want to integrate out all the noise

Sampling-based approximation: do many forward passes with different dropout masks, average all predictions

Can in fact do this with a single forward pass! (approximately) Leave all input neurons turned on (no dropout).

Q: Suppose that with all inputs present at test time the output of this neuron is x.

What would its output be during training time, in expectation? (e.g. if p = 0.5)



Can in fact do this with a single forward pass! (approximately) Leave all input neurons turned on (no dropout).



during test: a = w0*x + w1*yduring train: $E[a] = \frac{1}{4} * (w0*0 + w1*0)$ w0*0 + w1*yw0*x + w1*0w0*x + w1*y) $= \frac{1}{4} * (2 \text{ w0*x} + 2 \text{ w1*y})$ $= \frac{1}{2} * (w0*x + w1*y)$

Can in fact do this with a single forward pass! (approximately) Leave all input neurons turned on (no dropout).

а w1 w0 Х

during test: a = w0*x + w1*yWith p=0.5, using all inputs in the forward pass would during train: inflate the activations by 2x from what the network was $E[a] = \frac{1}{4} * (w0*0 + w1*0)$ "used to" during training! => Have to compensate by w0*0 + w1*yscaling the activations back $WO^*X + W1^*d_{0}$ wn by $\frac{1}{2}$ w0*x + w1*y) $= \frac{1}{4} * (2 w0*x + 2 w1*y)$ = $\frac{1}{2} * (w0*x + w1*y)$

Remaining Pieces

- Optimizing machine learning objectives:
 - Stochastic Descent
 - Mini-batches
- Improving generalization
 - Drop-out
- Activation functions

- Initialization and batch normalization
- Computing the gradient $\nabla g(w)$
 - Gradient checking
 - Backprop



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -103

Remaining Pieces

- Optimizing machine learning objectives:
 - Stochastic Descent
 - Mini-batches
- Improving generalization
 - Drop-out
- Activation functions

- Initialization and batch normalization
- Computing the gradient $\nabla g(w)$
 - Gradient checking
 - Backprop

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -105

- First idea: Small random numbers

(gaussian with zero mean and 1e-2 standard deviation)

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -106

- First idea: Small random numbers

(gaussian with zero mean and 1e-2 standard deviation)

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

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -107

Lets look at some activation statistics

E.g. 10-layer net with 500 neurons on each layer, using tanh non-linearities, and initializing as described in last slide.

```
# assume some unit gaussian 10-D input data
D = np.random.randn(1000, 500)
hidden_layer_sizes = [500]*10
nonlinearities = ['tanh']*len(hidden_layer_sizes)
```

```
act = {'relu':lambda x:np.maximum(0,x), 'tanh':lambda x:np.tanh(x)}
Hs = {}
for i in xrange(len(hidden_layer_sizes)):
    X = D if i == 0 else Hs[i-1] # input at this layer
    fan_in = X.shape[1]
    fan_out = hidden_layer_sizes[i]
    W = np.random.randn(fan_in, fan_out) * 0.01 # layer initialization
```

```
H = np.dot(X, W) # matrix multiply
H = act[nonlinearities[i]](H) # nonlinearity
Hs[i] = H # cache result on this layer
```

```
# look at distributions at each layer
print 'input layer had mean %f and std %f' % (np.mean(D), np.std(D))
layer_means = [np.mean(H) for i,H in Hs.iteritems()]
layer_stds = [np.std(H) for i,H in Hs.iteritems()]
for i,H in Hs.iteritems():
    print 'hidden layer %d had mean %f and std %f' % (i+1, layer means[i], layer stds[i])
```

```
# plot the means and standard deviations
plt.figure()
plt.subplot(121)
plt.plot(Hs.keys(), layer_means, 'ob-')
plt.title('layer mean')
plt.subplot(122)
plt.plot(Hs.keys(), layer_stds, 'or-')
plt.title('layer std')
```

```
# plot the raw distributions
plt.figure()
for i,H in Hs.iteritems():
    plt.subplot(1,len(Hs),i+1)
    plt.hist(H.ravel(), 30, range=(-1,1))
```

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -108

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -109

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: What do the gradients look like?

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -10

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.

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -

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 0.000139 and std 0.228008

0.0012

0.0010

0.0008

0.0006

0.0004

layer mean

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

layer std

"Xavier initialization" [Glorot et al., 2010]

Reasonable initialization.

(Mathematical derivation assumes linear activations)



0.60

0.55

0.50

0.45

input layer had mean 0.000501 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 0.018408 and std 0.038583 hidden layer 10 had mean 0.018408 and std 0.026076

but when using the ReLU nonlinearity it breaks.



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -13

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 0.889348 hidden layer 9 had mean 0.557498 and std 0.845357 hidden layer 10 had mean 0.552531 and std 0.844523

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -14

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 0.889348 hidden layer 9 had mean 0.567498 and std 0.845357 hidden layer 10 had mean 0.552531 and std 0.844523

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



20 Jan 2016

\$ 00 51 01 5 2 02 53 0 \$ 00 51 01 5 2 00 51 0 \$ 00 51 01 5 2 00 51 0\$ 00 51 01 5 2 00 51 0\$ 00 51 01 5 2 00 51 0\$ 00 51 01 5 2 00 51 0\$ 0

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -115

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 -116

Batch Normalization

"you want unit gaussian activations? just make them so."

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

$$\widehat{x}^{(k)} = \frac{x^{(k)} - \mathbf{E}[x^{(k)}]}{\sqrt{\mathrm{Var}[x^{(k)}]}}$$

this is a vanilla differentiable function...

20 Jan 2016

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -117

20 Jan 2016

Batch Normalization

"you want unit gaussian activations? just make them so."

D

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -118

[loffe and Szegedy, 2015]

Batch Normalization

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?

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

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -19
Batch Normalization

Normalize:

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

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

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

[loffe and Szegedy, 2015]

Note, the network can learn:

$$\gamma^{(k)} = \sqrt{\operatorname{Var}[x^{(k)}]}$$
$$\beta^{(k)} = \operatorname{E}[x^{(k)}]$$

to recover the identity mapping.

20 Jan 2016

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -120

Batch Normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1...m}\};$ Parameters to be learned: γ , β **Output:** $\{y_i = BN_{\gamma,\beta}(x_i)\}$ $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i$ // mini-batch mean $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$ // mini-batch variance $\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}}$ // normalize $y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \mathbf{BN}_{\gamma,\beta}(x_i)$ // 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

20 Jan 2016

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -121

Batch Normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1m}\};$ Parameters to be learned: γ, β		Note: at test time BatchNorm layer functions differently:
Output: $\{y_i = BN_{\gamma,\beta}(x_i)\}$ $\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i$ $\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2$	// mini-batch mean // mini-batch variance	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)
$\widehat{x}_{i} \leftarrow \frac{x_{i} - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^{2} + \epsilon}}$ $y_{i} \leftarrow \gamma \widehat{x}_{i} + \beta \equiv BN_{\gamma,\beta}(x_{i})$	// normalize // scale and shift	

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 5 -122

Remaining Pieces

- Optimizing machine learning objectives:
 - Stochastic Descent
 - Mini-batches
- Improving generalization
 - Drop-out
- Activation functions

- Initialization and batch normalization
- Computing the gradient $\nabla g(w)$
 - Gradient checking
 - Backprop

Gradient Descent

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

Numerical gradient: slow :(, approximate :(, easy to write :) Analytic gradient: fast :), exact :), error-prone :(

In practice: Derive analytic gradient, check your implementation with numerical gradient

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -124

Computational Graph



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -125



Lecture 4 -126



Lecture 4 -127



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

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



Lecture 4 -129

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

e.g. x = -2, y = 5, z = -4
$$q = x + y \quad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$
$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$
$$\frac{\partial f}{\partial q} = \frac{\partial f}{\partial q} = \frac{\partial f}{\partial q} = \frac{\partial f}{\partial z} = \frac{\partial f}{\partial z} = \frac{\partial f}{\partial z}$$

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

Lecture 4 -130

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

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

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

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

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

Lecture 4 -131

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

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

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

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

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

Lecture 4 -132

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

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

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

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

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

Lecture 4 -133

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

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

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

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

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

Lecture 4 -134

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

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

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

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

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

Lecture 4 -135

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

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

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

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

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

Lecture 4 -136

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

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

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

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

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

Lecture 4 -137

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

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

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

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

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

$$\overset{x -2}{y = 5}, z = -4$$

$$x -2 \qquad y = 5, z = -4$$

$$y = \frac{y = 5}{-4} \qquad (x + y) = \frac{1}{2}$$

$$\frac{y = -4}{3}$$

Chain rule: $\frac{\partial f}{\partial y}$

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

Lecture 4 -138

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

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

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

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

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

Lecture 4 -139

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

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

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

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

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

$$\overset{x -2}{-4}$$

$$\frac{z -4}{3}$$

Chain rule:

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

Lecture 4 - 40



Lecture 4 -141



Lecture 4 -142



Lecture 4 -143



Lecture 4 -144



Lecture 4 -145



Lecture 4 -146

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -147

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -148

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -149

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -150

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -151

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -152

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -153

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -154

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -155
Another example:

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -156

Another example:

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -157

Another example:

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



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -158



Lecture 4 -159



Lecture 4 -160

$$\begin{aligned} f(w,x) &= \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}} & \sigma(x) &= \frac{1}{1 + e^{-x}} & \text{sigmoid function} \\ \frac{d\sigma(x)}{dx} &= \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x)) \sigma(x) \end{aligned}$$



Lecture 4 -161

$$\begin{aligned} f(w,x) &= \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}} & \sigma(x) &= \frac{1}{1 + e^{-x}} & \text{sigmoid function} \\ \frac{d\sigma(x)}{dx} &= \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x)) \sigma(x) \end{aligned}$$



Lecture 4 -162

Patterns in backward flow

add gate: gradient distributor
max gate: gradient router
mul gate: gradient...?



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -163

Gradients add at branches



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -164

Implementation: forward/backward API



Graph (or Net) object. (Rough psuedo code)

lass Co	<pre>omputationalGraph(object):</pre>
#	
def	<pre>forward(inputs):</pre>
	# 1. [pass inputs to input gates]
	# 2. forward the computational graph:
	<pre>for gate in self.graph.nodes_topologically_sorted():</pre>
	gate.forward()
	<pre>return loss # the final gate in the graph outputs the loss</pre>
def	backward():
	<pre>for gate in reversed(self.graph.nodes_topologically_sorted()):</pre>
	<pre>gate.backward() # little piece of backprop (chain rule applied)</pre>
	return inputs gradients

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -165

Implementation: forward/backward API



Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -166

Implementation: forward/backward API



(x,y,z are scalars)

CSSB0

13 Jan 2016

Fei-Fei Li & Andrej Karpathy & Justin Johnson

Lecture 4 -167

Deep Learning Frameworks

- TensorFlow (in your Project 6!)
- Theano
- Torch
- CAFFE

Computation Graph Toolkit (CGT)