# Lecture 7: Training Neural Networks, Part 2

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 1 April 24, 2018

#### Administrative

- Assignment 1 is being graded, stay tuned
- Project proposals due tomorrow by 11:59pm on Gradescope
- Assignment 2 is out, due Wednesday 5/2 11:59pm

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 2 April 24, 2018

# Last time: Activation Functions



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



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



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 3 April 24, 2018

# Last time: Activation Functions



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



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



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 4 April 24, 2018

# Last time: Weight Initialization



#### Initialization too small:

Activations go to zero, gradients also zero, No learning

Initialization too big: Activations saturate (for tanh), Gradients zero, no learning

#### Initialization just right:

Nice distribution of activations at all layers, Learning proceeds nicely

#### Fei-Fei Li & Justin Johnson & Serena Yeung

#### Lecture 7 - 5 April 24, 2018

# Last time: Data Preprocessing



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 6 April 24, 2018

# Last time: Data Preprocessing

**Before normalization**: classification loss very sensitive to changes in weight matrix; hard to optimize After normalization: less sensitive to small changes in weights; easier to optimize



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 7 April 24, 2018

## Last time: Babysitting Learning



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 8 April 24, 2018

# Last time: Hyperparameter Search



#### Coarse to fine search

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

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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

#### Lecture 7 - 9 April 24, 2018

# Today

- More normalization
- Fancier optimization
- Regularization
- Transfer Learning

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 10 April 24, 2018

#### Last time: Batch Normalization

Input: 
$$x: N \times D$$

### Learnable params:

 $\gamma, \beta: D$ 

# Intermediates: $\frac{\mu, \sigma : D}{\hat{x} \cdot N \times D}$

**Output**:  $y: N \times D$ 



Fei-Fei Li & Justin Johnson & Serena Yeung

April 24, 2018 Lecture 7 - 11

# Last time: Batch Normalization

Estimate mean and variance from minibatch; Can't do this at test-time

Input:  $x : N \times D$ 

## Learnable params:

 $\gamma,eta:D$ 

Intermediates:

 $egin{array}{ll} \mu, \sigma : D \ \hat{x} : N imes D \end{array}$ 

**Output**:  $y : N \times D$ 

$$\mu_j = \frac{1}{N} \sum_{i=1}^N x_{i,j}$$
$$\sigma_j^2 = \frac{1}{N} \sum_{i=1}^N (x_{i,j} - \mu_j)^2$$
$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$

$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 12 April 24, 2018

#### **Batch Normalization: Test Time**

Input: 
$$x : N \times D$$

#### Learnable params:

 $\gamma,eta:D$ 

 $\mu, \sigma : D$  $\hat{x} : N \times D$ 

**Output**:  $y : N \times D$ 

 $\mu_j = \mathop{\rm (Running)}_{\rm seen \ during \ training} {\rm for all us}$ 

 $\sigma_j^2 = \mathop{\rm (Running)}_{\rm seen \ during \ training} {\rm draining}$ 

$$\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$$
$$y_{i,j} = \gamma_j \hat{x}_{i,j} + \beta_j$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 13 April 24, 2018

# **Batch Normalization for ConvNets**

Batch Normalization for **fully-connected** networks

Batch Normalization for **convolutional** networks (Spatial Batchnorm, BatchNorm2D)

x: N × Dx: N×C×H×WNormalize $\checkmark$ Normalize $\mu, \sigma$ : 1 × D $\mu, \sigma$ : 1×C×1×1 $\gamma, \beta$ : 1 × D $\gamma, \beta$ : 1×C×1×1 $\gamma = \gamma(x-\mu)/\sigma+\beta$  $\gamma = \gamma(x-\mu)/\sigma+\beta$ 

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 14 April 24, 2018

# Layer Normalization

Batch Normalization for fully-connected networks

x: N × Dx: N × DNormalize $\checkmark$  $\mu, \sigma: 1 \times D$ Normalize $\gamma, \beta: 1 \times D$  $\gamma, \beta: 1 \times D$  $\gamma = \gamma(x-\mu)/\sigma+\beta$  $\gamma = \gamma(x-\mu)/\sigma+\beta$ 

Ba, Kiros, and Hinton, "Layer Normalization", arXiv 2016

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 15 April 24, 2018

Layer Normalization for

fully-connected networks

Same behavior at train and test!

Can be used in recurrent networks

# **Instance Normalization**

Batch Normalization for convolutional networks



Ulyanov et al, Improved Texture Networks: Maximizing Quality and Diversity in Feed-forward Stylization and Texture Synthesis, CVPR 2017

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 16 April 24, 2018

Instance Normalization for

Same behavior at train / test!

convolutional networks

# **Comparison of Normalization Layers**



Wu and He, "Group Normalization", arXiv 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 17 April 24, 2018

# **Group Normalization**



Wu and He, "Group Normalization", arXiv 2018 (Appeared 3/22/2018)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 18 April 24, 2018

# **Decorrelated Batch Normalization**



#### **Batch Normalization**

 $\hat{x}_{i,j} = \frac{x_{i,j} - \mu_j}{\sqrt{\sigma_j^2 + \varepsilon}}$ 

BatchNorm normalizes the data, but cannot correct for correlations among the input features

#### **Decorrelated Batch Normalization**



DBN **whitens** the data using the full covariance matrix of the minibatch; this corrects for correlations

April 24, 2018

Huang et al, "Decorrelated Batch Normalization", arXiv 2018 (Appeared 4/23/2018)

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 -

### Optimization

# Vanilla Gradient Descent

while True: weights\_grad = evaluate\_gradient(loss\_fun, data, weights) weights += - step\_size \* weights\_grad # perform parameter update



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 20 April 24, 2018

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 21 April 24, 2018

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 22 April 24, 2018

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



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 23 April 24, 2018

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

Zero gradient, gradient descent gets stuck

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 24 April 24, 2018

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

Saddle points much more common in high dimension

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 25 April 24, 2018

Our gradients come from minibatches so they can be noisy!

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

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



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 26 April 24, 2018

## SGD + Momentum

SGD

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

while True: dx = compute\_gradient(x) x -= learning\_rate \* dx

#### SGD+Momentum

 $\begin{aligned} v_{t+1} &= \rho v_t + \nabla f(x_t) \\ x_{t+1} &= x_t - \alpha v_{t+1} \end{aligned}$ vx = 0 while True: dx = compute\_gradient(x) vx = rho \* vx + dx x -= learning\_rate \* vx

- 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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 27 April 24, 2018

# SGD + Momentum

#### SGD+Momentum

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

vx = 0
while True:
 dx = compute\_gradient(x)
 vx = rho \* vx - learning\_rate \* dx
 x += vx

#### SGD+Momentum

```
\begin{aligned} v_{t+1} &= \rho v_t + \nabla f(x_t) \\ x_{t+1} &= x_t - \alpha v_{t+1} \end{aligned}
vx = 0
while True:
dx = compute_gradient(x)
vx = rho * vx + dx
x -= learning_rate * vx
```

Lecture 7 - 28

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

Fei-Fei Li & Justin Johnson & Serena Yeung

April 24, 2<u>018</u>

# SGD + Momentum Grad Local Minima Saddle points

#### **Gradient Noise**



Fei-Fei Li & Justin Johnson & Serena Yeung

Poor Conditioning

Lecture 7 - 29 April 24, 2018

∎SGD

SGD+Momentum

# 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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 30 April 24, 2018

#### Momentum update:



#### 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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

#### Lecture 7 - 31 April 24, 2018

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



"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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

#### Lecture 7 - 32 April 24, 2018

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

$$x_{t+1} = x_t + v_{t+1}$$

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



"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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 33 April 24, 2018

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

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

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

rearrange:

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



"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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

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

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

#### Lecture 7 - 34 April 24, 2018

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

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:

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

dx = compute\_gradient(x)
old\_v = v
v = rho \* v - learning\_rate \* dx
x += -rho \* old\_v + (1 + rho) \* v

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 35 April 24, 2018



#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 36 April 24, 2018
grad\_squared = 0
while True:
 dx = compute\_gradient(x)
 grad\_squared += dx \* dx
 x -= learning\_rate \* dx / (np.sqrt(grad\_squared) + 1e-7)

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 37 April 24, 2018





#### Q: What happens with AdaGrad?

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 38 April 24, 2018



Q: What happens with AdaGrad? Progress along "steep" directions is damped; progress along "flat" directions is accelerated

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 39 April 24, 2018





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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 40 April 24, 2018



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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 41 April 24, 2018

# RMSProp



Tieleman and Hinton, 2012

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 42 April 24, 2018

# **RMSProp**



#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 43 April 24, 2018

# Adam (almost)

```
first_moment = 0
second_moment = 0
while True:
    dx = compute_gradient(x)
    first_moment = beta1 * first_moment + (1 - beta1) * dx
    second_moment = beta2 * second_moment + (1 - beta2) * dx * dx
    x -= learning_rate * first_moment / (np.sqrt(second_moment) + 1e-7))
```

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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 44 April 24, 2018

# Adam (almost)



Sort of like RMSProp with momentum

#### Q: What happens at first timestep?

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 45 April 24, 2018

# Adam (full form)



Lecture 7 - 46

April 24, 2018

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

Fei-Fei Li & Justin Johnson & Serena Yeung

# Adam (full form)



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 or 5e-4 is a great starting point for many models!

Lecture 7 - 47

April 24, 2018

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

Fei-Fei Li & Justin Johnson & Serena Yeung

# Adam



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 48 April 24, 2018



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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 49 April 24, 2018



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

#### step decay:

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

exponential decay:

$$lpha=lpha_0 e^{-kt}$$

1/t decay: $lpha=lpha_0/(1+kt)$ 

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 50 April 24, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 51 April 24, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 52 April 24, 2018

# **First-Order Optimization**



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 53 April 2<u>4, 2018</u>

# **First-Order Optimization**



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 54 April 2<u>4, 2018</u>



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 55 April 24, 2018

second-order Taylor expansion:

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

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

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

#### Q: What is nice about this update?

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 56 April 24, 2018

second-order Taylor expansion:

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

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

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

No hyperparameters! No learning rate! (Though you might use one in practice)

#### Q: What is nice about this update?

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 57 April 24, 2018

second-order Taylor expansion:

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

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

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

#### Q2: Why is this bad for deep learning?

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 58 April 24, 2018

second-order Taylor expansion:

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

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

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

Hessian has O(N<sup>2</sup>) elements Inverting takes O(N<sup>3</sup>) N = (Tens or Hundreds of) Millions

#### Q2: Why is this bad for deep learning?

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 59 April 24, 2018

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

- Quasi-Newton methods (BGFS most popular): instead of inverting the Hessian (O(n^3)), approximate inverse Hessian with rank 1 updates over time (O(n^2) each).
- **L-BFGS** (Limited memory BFGS): Does not form/store the full inverse Hessian.

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 60 April 24, 2018

## L-BFGS

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 61 April 24, 2018

Le et al, "On optimization methods for deep learning, ICML 2011" Ba et al, "Distributed second-order optimization using Kronecker-factored approximations", ICLR 2017

# In practice:

- Adam is a good default choice in many cases
- **SGD+Momentum** with learning rate decay often outperforms Adam by a bit, but requires 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)

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 62 April 24, 2018

# **Beyond Training Error**



help reduce training loss

But we really care about error on new data - how to reduce the gap?

Fei-Fei Li & Justin Johnson & Serena Yeung

April 24, 2018 Lecture 7 - 63

# **Early Stopping**



Stop training the model when accuracy on the validation set decreases Or train for a long time, but always keep track of the model snapshot that worked best on val

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 64 April 24, 2018

# Model Ensembles

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

(Take average of predicted probability distributions, then choose argmax)

## Enjoy 2% extra performance

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 65 April 24, 2018

# Model Ensembles: Tips and Tricks

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



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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 66 April 24, 2018

# Model Ensembles: Tips and Tricks

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



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



Cyclic learning rate schedules can make this work even better!

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 67 April 24, 2018

# Model Ensembles: Tips and Tricks

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



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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 68

April 24, 2018

### How to improve single-model performance?



#### Regularization

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 69 April 24, 2018

# Regularization: Add term to loss

$$L = rac{1}{N} \sum_{i=1}^{N} \sum_{j 
eq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \lambda R(W)$$

# In common use:L2 regularization $R(W) = \sum_k \sum_l W_{k,l}^2$ (Weight decay)L1 regularization $R(W) = \sum_k \sum_l |W_{k,l}|$ Elastic net (L1 + L2) $R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 70 April 24, 2018

# **Regularization:** Dropout

In each forward pass, randomly set some neurons to zero Probability of dropping is a hyperparameter; 0.5 is common





Srivastava et al, "Dropout: A simple way to prevent neural networks from overfitting", JMLR 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 71 April 24, 2018

# **Regularization: Dropout**

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

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

```
# forward pass for example 3-layer neural network
H1 = np.maximum(0, np.dot(W1, X) + b1)
U1 = np.random.rand(*H1.shape)
```

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

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Example forward pass with a 3-layer network using dropout



Lecture 7 - 72 April 24, 2018
# Regularization: Dropout

How can this possibly be a good idea?



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



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 73 April 24, 2018

# **Regularization:** Dropout

How can this possibly be a good idea?



Another interpretation:

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

Each binary mask is one model

An FC layer with 4096 units has  $2^{4096} \sim 10^{1233}$  possible masks! Only ~  $10^{82}$  atoms in the universe...

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 74 April 24, 2018

Dropout makes our output random!



Want to "average out" the randomness at test-time

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

But this integral seems hard ...

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 75 April 24, 2018

Want to approximate the integral

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

Consider a single neuron.



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 76 April 24, 2018

Want to approximate the integral

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

Consider a single neuron.



At test time we have: 
$$E\left[a
ight]=w_{1}x+w_{2}y_{2}$$

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 77 April 24, 2018

Want to approximate the integral

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

Consider a single neuron.



At test time we have:  $E[a] = w_1 x + w_2 y$ During training we have:  $E[a] = \frac{1}{4}(w_1 x + w_2 y) + \frac{1}{4}(w_1 x + 0y) + \frac{1}{4}(0x + w_2 y) + \frac{1}{4}(0x + w_2 y) = \frac{1}{2}(w_1 x + w_2 y)$ 

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 78 April 24, 2018

Want to approximate the integral

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

Consider a single neuron.



At test time we have:  $E[a] = w_1 x + w_2 y$ During training we have:  $E[a] = \frac{1}{4}(w_1 x + w_2 y) + \frac{1}{4}(w_1 x + 0y)$ At test time, **multiply** by dropout probability  $= \frac{1}{2}(w_1 x + w_2 y)$ 

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 79 April 24, 2018

#### def predict(X):

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

### At test time all neurons are active always => We must scale the activations so that for each neuron: <u>output at test time</u> = <u>expected output at training time</u>

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 80

<u>April 24, 2018</u>

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

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



#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 81 April 24, 2018

**Dropout Summary** 

## More common: "Inverted dropout"

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



#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 82 April 24, 2018

## Regularization: A common pattern

**Training**: Add some kind of randomness

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

**Testing:** Average out randomness (sometimes approximate)

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 83 April 24, 2018

## Regularization: A common pattern

# **Training**: Add some kind of randomness

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

**Testing:** Average out randomness (sometimes approximate)

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

**Example**: Batch Normalization

### **Training**: Normalize using stats from random minibatches

**Testing**: Use fixed stats to normalize

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 84 April 24, 2018

## **Regularization: Data Augmentation**



#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 85 April 24, 2018

## **Regularization: Data Augmentation**



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 86 April 24, 2018

## Data Augmentation Horizontal Flips





Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 87 April 24, 2018

## Data Augmentation Random crops and scales

**Training**: sample random crops / scales ResNet:

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



#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 88 April 24, 2018

## Data Augmentation Random crops and scales

**Training**: sample random crops / scales ResNet:

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



# **Testing**: average a fixed set of crops ResNet:

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 89 April 24, 2018

## Data Augmentation Color Jitter

Simple: Randomize contrast and brightness



#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 90 April 24, 2018

## Data Augmentation Color Jitter

Simple: Randomize contrast and brightness



### More Complex:

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

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 91 April 24, 2018

## Data Augmentation

Get creative for your problem!

## Random mix/combinations of :

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 92 April 24, 2018

## Regularization: A common pattern

**Training**: Add random noise **Testing**: Marginalize over the noise

### Examples:

- Dropout
- **Batch Normalization**
- Data Augmentation

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 93 April 24, 2018

## Regularization: A common pattern Training: Add random noise

Testing: Marginalize over the noise

### Examples:

Dropout Batch Normalization Data Augmentation DropConnect





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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 94 April 24, 2018

# Regularization: A common pattern

**Training**: Add random noise **Testing**: Marginalize over the noise

## Examples:

Dropout

Batch Normalization

Data Augmentation DropConnect Fractional Max Pooling



Graham, "Fractional Max Pooling", arXiv 2014

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 95 April 24, 2018

## Regularization: A common pattern

**Training**: Add random noise **Testing**: Marginalize over the noise

### Examples:

Dropout

**Batch Normalization** 

Data Augmentation DropConnect

Fractional Max Pooling Stochastic Depth

April 24, 2018

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

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - <u>96</u>

# **Transfer Learning**

# "You need a lot of a data if you want to train/use CNNs"

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 97 April 24, 2018



Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 98 April 24, 2018

### Transfer Learning with CNNs

1. Train on Imagenet

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
0.0000 050
CONV-256
Conv-256
MaxPool
Conv-128
Conv-128
MaxPool
Conv-64
Conv-64
Image

Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014 Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

#### Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 - 99 April 24, 2018

### Transfer Learning with CNNs

1. Train on Imagenet

FC-1000
FC-4096
FC-4096
MaxPool
Conv-512
Conv-512
MaxPool
Conv-512
Conv-512
MaxPool
MaxPool Conv-256
MaxPool Conv-256 Conv-256
MaxPool Conv-256 Conv-256 MaxPool
MaxPool Conv-256 Conv-256 MaxPool Conv-128
MaxPool Conv-256 Conv-256 MaxPool Conv-128 Conv-128
MaxPool Conv-256 Conv-256 MaxPool Conv-128 Conv-128 MaxPool
MaxPool Conv-256 Conv-256 MaxPool Conv-128 Conv-128 MaxPool Conv-64

Image

2. Small Dataset (C classes)



Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014 Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014

April 24, 2018

10

 $\left( \right)$ 

Lecture 7 -

### Transfer Learning with CNNs

1. Train on Imagenet

FC-1000	
FC-4096	
FC-4096	
MaxPool	
Conv-512	
Conv-512	
MaxPool	
Conv-512	
Conv-512	
MaxPool	
MaxPool Conv-256	
MaxPool Conv-256 Conv-256	
MaxPool Conv-256 Conv-256 MaxPool	
MaxPool Conv-256 Conv-256 MaxPool Conv-128	
MaxPool Conv-256 Conv-256 MaxPool Conv-128 Conv-128	
MaxPool Conv-256 Conv-256 MaxPool Conv-128 Conv-128 MaxPool	
MaxPool Conv-256 Conv-256 MaxPool Conv-128 MaxPool Conv-64	

Image

2. Small Dataset (C classes)



Donahue et al, "DeCAF: A Deep Convolutional Activation Feature for Generic Visual Recognition", ICML 2014 Razavian et al, "CNN Features Off-the-Shelf: An Astounding Baseline for Recognition", CVPR Workshops 2014



April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung

Lecture 7 -

FC-1000 FC-4096 FC-4096 MaxPool		very similar dataset	very different dataset
Conv-512 MaxPool Conv-512 MaxPool Conv-512 MaxPool Conv-256	very little data	?	?
MaxPool Conv-128			
Conv-128 MaxPool Conv-64 Conv-64 Image	quite a lot of data	?	?

Lecture 7 -  $\frac{10}{2}$ 

April 24, 2018

FC-1000 FC-4096 FC-4096 MaxPool		very similar dataset	very different dataset
Conv-512MaxPoolMore specificConv-512More specificMaxPoolMore genericMaxPoolMore genericMaxPoolMore generic	very little data	Use Linear Classifier on top layer	?
Conv-128 Conv-128 MaxPool Conv-64 Conv-64 Image	quite a lot of data	Finetune a few layers	?

Lecture 7 -  $\frac{10}{3}$ 

April 24, 2018

FC-1000 FC-4096 FC-4096 MaxPool		very similar dataset	very different dataset
Conv-512MaxPoolMore specificConv-512MaxPoolMaxPoolConv-256Conv-256More genericMaxPoolMore generic	very little data	Use Linear Classifier on top layer	You're in trouble Try linear classifier from different stages
Conv-128 Conv-128 MaxPool Conv-64 Conv-64 Image	quite a lot of data	Finetune a few layers	Finetune a larger number of layers

Lecture 7 -  $\frac{10}{4}$ 

April 24, 2018

# Transfer learning with CNNs is pervasive... (it's the norm, not an exception)



#### Image Captioning: CNN + RNN



10

5

Lecture 7 -

#### Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015 Figure copyright IEEE, 2015. Reproduced for educational purposes.

April 24, 2018

Girshick, "Fast R-CNN", ICCV 2015 Figure copyright Ross Girshick, 2015. Reproduced with permission.

# Transfer learning with CNNs is pervasive... (it's the norm, not an exception)



Karpathy and Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions", CVPR 2015 Figure copyright IEEE, 2015. Reproduced for educational purposes.

April 24, 2018

10

6

Lecture 7 -

Girshick, "Fast R-CNN", ICCV 2015 Figure copyright Ross Girshick, 2015. Reproduced with permission.

# Transfer learning with CNNs is pervasive... (it's the norm, not an exception)



## Takeaway for your projects and beyond:

Have some dataset of interest but it has < ~1M images?

Lecture 7 -

April 24, 2018

- 1. Find a very large dataset that has similar data, train a big ConvNet there
- 2. Transfer learn to your dataset

Deep learning frameworks provide a "Model Zoo" of pretrained models so you don't need to train your own

Caffe: <u>https://github.com/BVLC/caffe/wiki/Model-Zoo</u> TensorFlow: <u>https://github.com/tensorflow/models</u> PyTorch: <u>https://github.com/pytorch/vision</u>
## Summary

- Lots of Batch Normalization variants
- Optimization
  - Momentum, RMSProp, Adam, etc

Lecture 7 -

April 24, 2018

- Regularization
  - Dropout, etc
- Transfer learning
  - Use this for your projects!

Fei-Fei Li & Justin Johnson & Serena Yeung

## Next time: Deep Learning Software!

Lecture 7 -

April 24, 2018

Fei-Fei Li & Justin Johnson & Serena Yeung