Lecture 10: Training Neural Networks (Part 1)

Justin Johnson

Lecture 1 - 1

Reminder: A3

• Due Friday, October 9

Midterm Exam

We are still working out details! Will share more on Wednesday

- Will (most likely) be online via <u>https://crabster.org/</u>
- Material up to Lecture 13 is fair game
- Mostly conceptual questions, no coding
- Some combination of:
 - True / False
 - Multiple choice
 - Short answer (math on paper)
- If you need accommodations, send your SSD letter to me

Last Time: Hardware and Software

CPU

GPU

TPU



Static Graphs vs **Dynamic Graphs**

PyTorch vs TensorFlow

Justin Johnson

Lecture 10 - 4

Overview

1.One time setup

Activation functions, data preprocessing, weight initialization, regularization

2. Training dynamics

Learning rate schedules; large-batch training;

hyperparameter optimization

3.After training

Model ensembles, transfer learning



1.One time setup

Activation functions, data preprocessing, weight initialization, regularization

2. Training dynamics

Learning rate schedules; large-batch training;

hyperparameter optimization

3.After training

Model ensembles, transfer learning

Today

Next time

Justin Johnson

Lecture 10 - 6

Activation Functions

Justin Johnson

Lecture 10 - 7

Activation Functions



Justin Johnson

Lecture 10 - 8

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

Lecture 10 - 9



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

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

Justin Johnson

Lecture 10 - 10



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

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

3 problems:

1. Saturated neurons "kill" the gradients

Justin Johnson

Lecture 10 - 11



What happens when x = -10? What happens when x = 0? What happens when x = 10?

Lecture 10 - 12



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

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

3 problems:

1. Saturated neurons "kill" the gradients

Justin Johnson

Lecture 10 - 13



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

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

3 problems:

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

$$h_i^{(\ell)} = \sum_j w_{i,j}^{(\ell)} \sigma\left(h_j^{(\ell-1)}\right) + b_i^{(\ell)}$$

 $h_i^{(\ell)}$ is the *i*th element of the hidden layer at layer ℓ (before activation) $w^{(\ell)}$, $b^{(\ell)}$ are the weights and bias of layer ℓ

What can we say about the gradients on $w^{(\ell)}$?

$$h_i^{(\ell)} = \sum_j w_{i,j}^{(\ell)} \sigma\left(h_j^{(\ell-1)}\right) + b_i^{(\ell)}$$

 $h_i^{(\ell)}$ is the *i*th element of the hidden layer at layer ℓ (before activation) $w^{(\ell)}$, $b^{(\ell)}$ are the weights and bias of layer ℓ

What can we say about the gradients on $w^{(\ell)}$?

$$h_i^{(\ell)} = \sum_j w_{i,j}^{(\ell)} \sigma\left(h_j^{(\ell-1)}\right) + b_i^{(\ell)}$$

 $h_i^{(\ell)}$ is the *i*th element of the hidden layer at layer ℓ (before activation) $w^{(\ell)}$, $b^{(\ell)}$ are the weights and bias of layer ℓ

What can we say about the gradients on $w^{(\ell)}$? Always all positive or all negative :(



Lecture 10 - 17

$$h_i^{(\ell)} = \sum_j w_{i,j}^{(\ell)} \sigma\left(h_j^{(\ell-1)}\right) + b_i^{(\ell)}$$

 $h_i^{(\ell)}$ is the *i*th element of the hidden layer at layer ℓ (before activation) $w^{(\ell)}$, $b^{(\ell)}$ are the weights and bias of layer ℓ

What can we say about the gradients on $w^{(\ell)}$? Always all positive or all negative :((For a single element! Minibatches help)



Lecture 10 - 18



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

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

3 problems:

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



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

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

3 problems:

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

Justin Johnson

Lecture 10 - 20

Activation Functions: Tanh



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

Justin Johnson

Lecture 10 - 21

Activation Functions: ReLU f(x) =



ReLU (Rectified Linear Unit) f(x) = max(0,x)

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

Justin Johnson

Lecture 10 - 22

Activation Functions: ReLU f(x) =



f(x) = max(0,x)

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

- Not zero-centered output

ReLU (Rectified Linear Unit)

Activation Functions: ReLU f(



$$f(x) = max(0,x)$$

- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

- Not zero-centered output
- An annoyance:

hint: what is the gradient when x < 0?

Justin Johnson

Lecture 10 - 24

Activation Functions: ReLU



What happens when x = -10? What happens when x = 0? What happens when x = 10?



Justin Johnson

Lecture 10 - 26



Justin Johnson

Lecture 10 - 27

Activation Functions: Leaky ReLU

-

Lecture 10 - 28



Justin Johnson

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

Activation Functions: Leaky ReLU



Justin Johnson

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

Parametric ReLU (PReLU) $f(x) = \max(\alpha x, x)$ α is learned via backprop

> He et al, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification", ICCV 2015

October 5, 2020

Lecture 10 - 29

Activation Functions: Exponential Linear Unit (ELU)



- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

Computation requires exp()

(Default alpha=1)

Justin Johnson

Lecture 10 - 30

Activation Functions: Scaled Exponential Linear Unit (SELU)



 Scaled version of ELU that works better for deep networks
"Self-Normalizing" property; can train deep SELU networks without BatchNorm

- $\alpha = 1.6732632423543772848170429916717$
- $\lambda \ = \ 1.0507009873554804934193349852946$

Klambauer et al, Self-Normalizing Neural Networks, ICLR 2017

Justin Johnson

Lecture 10 - 31

Activation Functions: Scaled Exponential Linear Unit (SELU)



Scaled version of ELU that works better for deep networks "Self-Normalizing" property; can train deep SELU networks without BatchNorm

Derivation takes 91 pages of math in appendix...

 $\alpha = 1.6732632423543772848170429916717$ $\lambda = 1.0507009873554804934193349852946$

Klambauer et al, Self-Normalizing Neural Networks, ICLR 2017

Justin Johnson

Lecture 10 - 32

Activation Functions: Gaussian Error Linear Unit (GELU)



- Idea: Multiply input by 0 or 1 at random; large values more likely to be multiplied by 1, small values more likely to be multiplied by 0 (data-dependent dropout)
- Take expectation over randomness
 - Very common in Transformers (BERT, GPT, GPT-2, GPT-3)

Hendrycks and Gimpel, Gaussian Error Linear Units (GELUs), 2016

Justin Johnson

Lecture 10 - 33

Accuracy on CIFAR10

Ramachandran et al, "Searching for activation functions", ICLR Workshop 2018

■ ReLU ■ Leaky ReLU ■ Parametric ReLU ■ Softplus ■ ELU ■ SELU ■ GELU ■ Swish



Activation Functions: Summary

- Don't think too hard. Just use ReLU
- Try out Leaky ReLU / ELU / SELU / GELU if you need to squeeze that last 0.1%
- Don't use sigmoid or tanh

Data Preprocessing

Justin Johnson

Lecture 10 - 36


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

Justin Johnson

Lecture 10 - 37

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

$$h_i^{(\ell)} = \sum_j w_{i,j}^{(\ell)} \sigma\left(h_j^{(\ell-1)}\right) + b_i^{(\ell)}$$

What can we say about the gradients on w? Always all positive or all negative :((this is also why you want zero-mean data!)



Justin Johnson

Lecture 10 - 38



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

Justin Johnson

Lecture 10 - 39

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



Justin Johnson

Lecture 10 - 40

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



Justin Johnson

Lecture 10 - 41

Data Preprocessing for Images

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

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

Not common to do PCA or whitening

Justin Johnson

Lecture 10 - 43



Q: What happens if we initialize all W=0, b=0?

Justin Johnson

Lecture 10 - 44



Q: What happens if we initialize all W=0, b=0?

A: All outputs are 0, all gradients are the same!No "symmetry breaking"

Justin Johnson

Lecture 10 - 45

Next idea: **small random numbers** (Gaussian with zero mean, std=0.01)

W = 0.01 * np.random.randn(Din, Dout)

Justin Johnson

Lecture 10 - 46

Next idea: **small random numbers** (Gaussian with zero mean, std=0.01)

W = 0.01 * np.random.randn(Din, Dout)

Works ~okay for small networks, but problems with deeper networks.

Justin Johnson

Lecture 10 - 47

```
dims = [4096] * 7 Forward pass for a 6-layer
hs = [] net with hidden size 4096
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

```
dims = [4096] * 7 Forward pass for a 6-layer
hs = [] net with hidden size 4096
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
W = 0.01 * np.random.randn(Din, Dout)
x = np.tanh(x.dot(W))
hs.append(x)
All ad
deep
```

All activations tend to zero for deeper network layers

Q: What do the gradients dL/dW look like?



Justin Johnson

Lecture 10 - 49

```
dims = [4096] * 7 Forward pass for a 6-layer
hs = [] net with hidden size 4096
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.01 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

All activations tend to zero for deeper network layers

Q: What do the gradients dL/dW look like?

A: All zero, no learning =(



Justin Johnson

Lecture 10 - 50

dims = [4096] * 7 Increase std of initial weights
hs = [] from 0.01 to 0.05
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
 W = 0.05 * np.random.randn(Din, Dout)
 x = np.tanh(x.dot(W))
 hs.append(x)





Justin Johnson

Lecture 10 - 52





Justin Johnson

Lecture 10 - 53

Glorot and Bengio, "Understanding the difficulty of training deep feedforward neural networks", AISTAT 2010

Justin Johnson

Lecture 10 - 54





Glorot and Bengio, "Understanding the difficulty of training deep feedforward neural networks", AISTAT 2010

Justin Johnson

Lecture 10 - 55





Glorot and Bengio, "Understanding the difficulty of training deep feedforward neural networks", AISTAT 2010

Justin Johnson

Lecture 10 - 56

"Xavier" initialization: std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$
 $y_i = \sum_{j=1}^{Din} x_j w_j$

Justin Johnson

Lecture 10 - 57

"Xavier" initialization: std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$
 $y_i = \sum_{j=1}^{Din} x_j w_j$

 $Var(y_i) = Din * Var(x_iw_i)$

[Assume x, w are iid]

Din

Justin Johnson

Lecture 10 - 58

"Xavier" initialization: std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

Dia

$$y = Wx$$
 $y_i = \sum_{j=1}^{Din} x_j w_j$

 $Var(y_i) = Din * Var(x_iw_i)$ $= Din * (E[x_i^2]E[w_i^2] - E[x_i]^2 E[w_i]^2)$ [Assume x, w independent]

"Xavier" initialization: std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

Dia

$$y = Wx$$
 $y_i = \sum_{j=1}^{Din} x_j w_j$

$Var(y_i) = Din * Var(x_iw_i)$ [Assume x, w are iid] = Din * (E[x_i^2]E[w_i^2] - E[x_i]^2 E[w_i]^2) [Assume x, w independent] = Din * Var(x_i) * Var(w_i) [Assume x, w are zero-mean]

"Xavier" initialization: std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$
 $y_i = \sum_{j=1}^{Din} x_j w_j$

 $Var(y_i) = Din * Var(x_iw_i)$ [Assume x, w are iid] = Din * (E[x_i^2]E[w_i^2] - E[x_i]^2 E[w_i]^2) [Assume x, w independent] = Din * Var(x_i) * Var(w_i) [Assume x, w are zero-mean]

If $Var(w_i) = 1/Din$ then $Var(y_i) = Var(x_i)$

Weight Initialization: What about ReLU?

```
dims = [4096] * 7 Change from tanh to ReLU
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

Weight Initialization: What about ReLU?





Justin Johnson

Lecture 10 - 63

Weight Initialization: Kaiming / MSRA Initialization





He et al, "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification", ICCV 2015

Justin Johnson

Lecture 10 - 64

Weight Initialization: Residual Networks



If we initialize with MSRA: then Var(F(x)) = Var(x)But then Var(F(x) + x) > Var(x) - variance growswith each block!

Residual Block

Justin Johnson

Lecture 10 - 65

Weight Initialization: Residual Networks



If we initialize with MSRA: then Var(F(x)) = Var(x) But then Var(F(x) + x) > Var(x) variance grows with each block!

Solution: Initialize first conv with MSRA, initialize second conv to zero. Then Var(x + F(x)) = Var(x)

Zhang et al, "Fixup Initialization: Residual Learning Without Normalization", ICLR 2019

Justin Johnson

Lecture 10 - 66

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

Fixup Initialization: Residual Learning Without Normalization, Zhang et al, 2019

The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks, Frankle and Carbin, 2019

Justin Johnson

Lecture 10 - 67

Now your model is training ... but it overfits!



Regularization

		•			
	1101	rin i		hnc	n
			JU		
<u> </u>			\sim \sim .		<u> </u>

Lecture 10 - 68

Regularization: Add term to the 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 L1 regularization Elastic net (L1 + L2)

 $egin{aligned} R(W) &= \sum_k \sum_l W_{k,l}^2 & ext{(Weight decay)} \ R(W) &= \sum_k \sum_l |W_{k,l}| \ R(W) &= \sum_k \sum_l eta W_{k,l}^2 + |W_{k,l}| \end{aligned}$

Justin Johnson

Lecture 10 - 69

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

Justin Johnson

Lecture 10 - 70

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)

Example forward pass with a 3-layer network using dropout



Justin Johnson

Lecture 10 - 71

Regularization: Dropout



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



Justin Johnson

Lecture 10 - 72
Regularization: Dropout



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...

Justin Johnson

Lecture 10 - 73

Dropout makes our output random!

Output Input
(label) (image)
$$\mathbf{y} = f_W(\mathbf{x}, \mathbf{z})$$
 Random
mask

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

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

But this integral seems hard ...

Justin Johnson

Lecture 10 - 74

Want to approximate the integral

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

Consider a single neuron:



At test time we have:
$$E[a] = w_1 x + w_2 y$$

Want to approximate the integral

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

Consider a single neuron:



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

Want to approximate the integral

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

 $\begin{array}{c}
 a \\
 w_1 \\
 w_2 \\
 \hline
 x \\
 y
\end{array}$

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, drop nothing and multiply by dropout probability $= \frac{1}{2}(w_1 x + w_2 y)$

Lecture 10 - 77

Consider a single neuron:

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>

Dropout Summary

Vanilla Dropout: Not recommended implementation (see notes below) """ **p** = 0.5 # probability of keeping a unit active. higher = less dropout def train step(X): """ X contains the data """ # forward pass for example 3-layer neural network H1 = np.maximum(0, np.dot(W1, X) + b1)U1 = np.random.rand(*H1.shape) H1 *= U1 # drop! H2 = np.maximum(0, np.dot(W2, H1) + b2)U2 = np.random.rand(*H2.shape) < p # second dropout mask H2 *= U2 # drop! out = np.dot(W3, H2) + b3 # backward pass: compute gradients... (not shown) # perform parameter update... (not shown) def predict(X): # ensembled forward pass H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations out = np.dot(W3, H2) + b3

drop in forward pass

scale at test time

Justin Johnson

Lecture 10 - 79

More common: "Inverted dropout"

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



Justin Johnson

Lecture 10 - 80

Dropout architectures

Recall AlexNet, VGG have most of their parameters in **fully-connected layers**; usually Dropout is applied there



Justin Johnson

Lecture 10 - 81

Dropout architectures

Recall AlexNet, VGG have most of their parameters in **fully-connected layers**; usually Dropout is applied there



Later architectures (GoogLeNet, ResNet, etc) use global average pooling instead of fully-connected layers: they don't use dropout at all!

Justin Johnson

Lecture 10 - 82

Training: Add some kind of randomness

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

Testing: Average out randomness (sometimes approximate)

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

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

Training: Add some kind of randomness

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

For ResNet and later, often L2 and Batch Normalization are the only regularizers! **Example**: Batch Normalization

Training: Normalize using stats from random minibatches

Testing: Average out randomness (sometimes approximate)

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

Testing: Use fixed stats to normalize

Data Augmentation



Justin Johnson

Lecture 10 - 86

Data Augmentation



Justin Johnson

Lecture 10 - 87

Data Augmentation: Horizontal Flips





Justin Johnson

Lecture 10 - 88

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



Justin Johnson

Lecture 10 - 89

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



Justin Johnson

Lecture 10 - 90

Data Augmentation: Color Jitter

Simple: Randomize contrast and brightness





More Complex:

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

(Used in AlexNet, ResNet, etc)

Justin Johnson

Lecture 10 - 91

Data Augmentation: Get creative for your problem!

Random mix/combinations of :

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

Training: Add some randomness **Testing**: Marginalize over randomness

Examples:

Dropout Batch Normalization Data Augmentation

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

Lecture 10 - 93

Regularization: DropConnect

Training: Drop random connections between neurons (set weight=0) **Testing**: Use all the connections

Examples:

Dropout Batch Normalization Data Augmentation DropConnect





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

Justin Johnson

Lecture 10 - 94

Regularization: Fractional Pooling

Training: Use randomized pooling regions **Testing**: Average predictions over different samples

Examples:

Dropout Batch Normalization Data Augmentation DropConnect Fractional Max Pooling



Lecture 10 - 95

Regularization: Stochastic Depth

Training: Skip some residual blocks in ResNet **Testing**: Use the whole network

Examples:

Dropout Batch Normalization Data Augmentation DropConnect Fractional Max Pooling Stochastic Depth



Justin Johnson

Regularization: Stochastic Depth

Training: Set random images regions to 0 **Testing**: Use the whole image

Examples:

Dropout Batch Normalization Data Augmentation DropConnect Fractional Max Pooling Stochastic Depth Cutout





Works very well for small datasets like CIFAR, less common for large datasets like ImageNet

Justin Johnson

Lecture 10 - 97

Regularization: Mixup

Training: Train on random blends of images **Testing**: Use original images

Examples:

Dropout Batch Normalization Data Augmentation DropConnect Fractional Max Pooling Stochastic Depth Cutout Mixup







CNN

Target label: cat: 0.4 dog: 0.6

Randomly blend the pixels of pairs of training images, e.g. 40% cat, 60% dog

Zhang et al, "mixup: Beyond Empirical Risk Minimization", ICLR 2018

Justin Johnson

Lecture 10 - 98

Regularization: Mixup

Training: Train on random blends of images **Testing**: Use original images

Examples:

Dropout Batch Normalization Data Augmentation DropConnect Fractional Max Pooling Stochastic Depth Cutout Mixup







Target label: cat: 0.4 dog: 0.6

Randomly blend the pixels of pairs of training images, e.g. 40% cat, 60% dog

Zhang et al, "mixup: Beyond Empirical Risk Minimization", ICLR 2018

Justin Johnson

Lecture 10 - 99

October 5, 2020

CNN

Regularization: Mixup

Training: Train on random blends of images **Testing**: Use original images

Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect Fractional Max Pooling Stochastic Depth

Cutout Mixup

- Consider dropout for large fullyconnected layers
 - Batch normalization and data augmentation almost always a good idea
 - Try cutout and mixup especially for small classification datasets

Justin Johnson

Lecture 10 - 100

_



1. One time setup

Activation functions, data preprocessing, weight initialization, regularization

2. Training dynamics

Learning rate schedules; large-batch training;

hyperparameter optimization

3. After training

Model ensembles, transfer learning

Today

Next time

Justin Johnson

Lecture 10 - 101

Next time: Training Neural Networks (part 2)

Justin Johnson

Lecture 10 - 102