

Machine Intelligence & Deep Learning Workshop

Raymond Ptucha, Majid Rabbani, Mark Smith

The Kate Gleason **COLLEGE OF
ENGINEERING**

Introduction to Deep Learning



Raymond Ptucha
June 27-29, 2018
Rochester Institute of Technology
www.rit.edu/kgcoe/cqas/machinelearning



© 2018 Ray Ptucha, Rochester Institute of Technology

1

Fair Use Agreement

This agreement covers the use of all slides in this document, please read carefully.

- You may freely use these slides for personal use, if:
 - My name (R. Ptucha) appears on each slide.
- You may freely use these slides externally, if:
 - You send me an email telling me the conference/venue/company name in advance, and which slides you wish to use.
 - You receive a positive confirmation email back from me.
 - My name (R. Ptucha) appears on each slide you use.

(c) Raymond Ptucha, rwpeec@rit.edu

© 2018 Ray Ptucha, Rochester Institute of Technology

2

Agenda

- **Wed, June 27**
 - 9-10:30am Regression and Classification
 - 10:30-10:45pm Break
 - 10:45-12:15pm Boosting and SVM
 - 12:15-1:30pm Lunch
 - 1:30-3:30pm Neural Networks and Dimensionality Reduction
 - 3:30-5pm Hands-on Python and Machine Learning
- **Thur, June 28**
 - 9-10:30am **Introduction to deep learning**
 - 10:30-10:45pm Break
 - 10:45-12:15pm Convolutional Neural Networks
 - 12:15-1:30pm Lunch
 - 1:30-3:30pm Region and pixel-level convolutions
 - 3:30-5pm Hands-on CNNs
- **Fri, June 29**
 - 9-10:30am Recurrent neural networks
 - 10:30-10:45pm Break
 - 10:45-12:15pm Language and Vision
 - 12:15-1:30pm Lunch
 - 1:30-3:30pm Graph convolutional neural networks; Generative adversarial networks
 - 3:30-5pm Hands-on regional CNNs, RNNs

© 2018 Ray Ptucha, Rochester Institute of Technology

3

Machine Learning



- Machine learning is giving computers the ability to analyze, generalize, think/reason/behave like humans.
- Machine learning is transforming medical research, financial markets, international security, and generally making humans more efficient and improving quality of life.
- Inspired by the mammalian brain, deep learning is machine learning on steroids- bigger, faster, better!

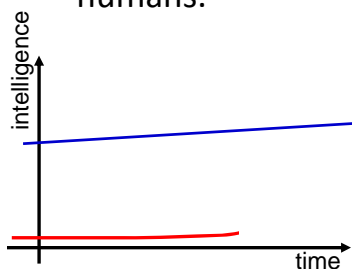


© 2018 Ray Ptucha, Rochester Institute of Technology

4

The point of Singularity

- The point of singularity is when computers become smarter than humans.



— Evolution of biology
— Advancement of technology

© 2018 Ray Ptucha, Rochester Institute of Technology

5

Unleashing of Intelligence

- Machines will slowly match, then quickly surpass human capabilities.
- Today it is exciting/scary/fun to drive next to an autonomous car.
- Tomorrow it may be considered irresponsible for a human to relinquish control from a car that has faster reaction times, doesn't drink/text/get distracted/tired, and is communicating with surrounding vehicles and objects.

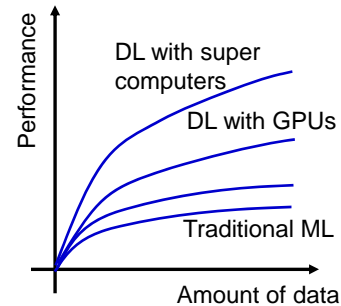


© 2018 Ray Ptucha, Rochester Institute of Technology

7

Why is AI (Deep Learning) Just Now Becoming Practical in Many Day-to-Day Situations?

- Availability of data;
- Sustained advances in hardware capabilities (including GPUs running machine learning workloads);
- Omnipresent connectivity;
- Lower cost and power consumption;
- Sustained advances in algorithmic learning techniques.



Hot trend:
High performance
architecture experts
teaming up with deep
learning experts

© 2018 Ray Ptucha, Rochester Institute of Technology

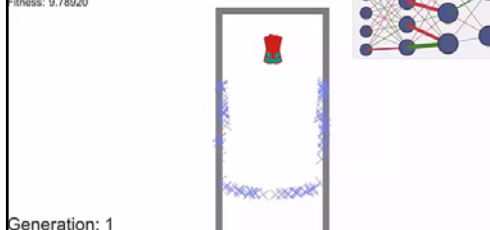
8

2017: The Year of AI: The Wall Street Journal, Forbes, and Fortune



NEC Face Recognition

Turn: -0.7932
Engine: 0.99999
Fitness: 9.78920



SONY Playstation Virtual Reality

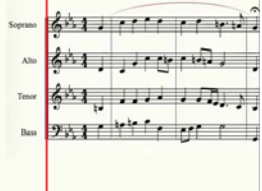
Evolutionary Reinforcement Learning

© 2018 Ray Ptucha, Rochester Institute of Technology

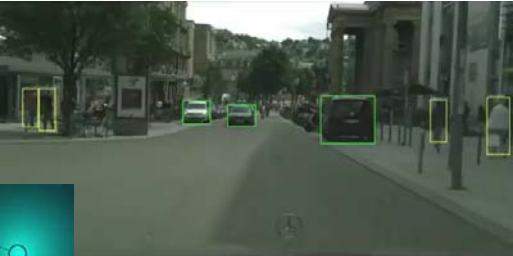
11

2017: The Year of AI:


The Wall Street Journal, Forbes, and Fortune



DeepBach



NVIDIA Autonomous Car Detection & Segmentation




<http://pjreddie.com/yolo>

YOLO v2 Object Detection

© 2018 Ray Ptucha, Rochester Institute of Technology

12

Some Things to Look for in 2018



CelebA-HQ
1024 x 1024

Progressive growing

CelebA-HQ
1024 x 1024

Latent space interpolations

http://research.nvidia.com/sites/default/files/pubs/2017-10_Progressive-Growing-of/karras2017gan-paper.pdf

© 2018 Ray Ptucha, Rochester Institute of Technology

17

Some Things to Look for in 2018

Faceshift GDC



Apple iPhone X, Animoji Yourself



© 2018 Ray Ptucha, Rochester Institute of Technology

18

Some Things to Look for in 2018

NVIDIA DRIVE
Autonomous Vehicle Platform
October 10, 2017



NVIDIA Drive

© 2018 Ray Ptucha, Rochester Institute of Technology

19

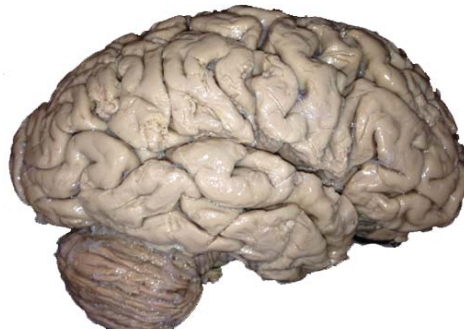
Prerequisites

- Basic probability and statistics
- Some programming experience (we will be using Python)
- Simple linear algebra
- Simple multivariate calculus
- Don't worry if you don't understand all concepts, over time things will make more sense.
- Feel free to stick around from 5-7pm in the evenings where the instructor and TAs can help you setup your own personal laptop, reinforce concepts, or answer personal applications of deep learning for your particular need.

© 2018 Ray Ptucha, Rochester Institute of Technology

20

The Human Brain



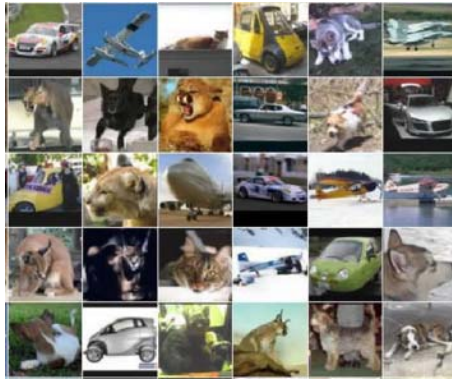
- We've learned more about the brain in the last 5 years than we have learned in the last 5000 years!
- It controls every aspect of our lives, but we still don't understand exactly how it works.

© 2018 Ray Ptucha, Rochester Institute of Technology

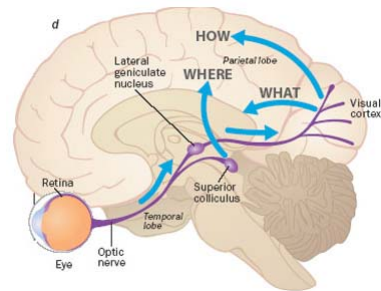
21

The Brain on Pattern Recognition

- Airplane, Cat, Car, Dog



STL-10 dataset



<http://thebraingeek.blogspot.com/2012/08/blindsight.html>

© 2018 Ray Ptucha, Rochester Institute of Technology

24

The Brain on Pattern Recognition

Despite Changes in Deformation:

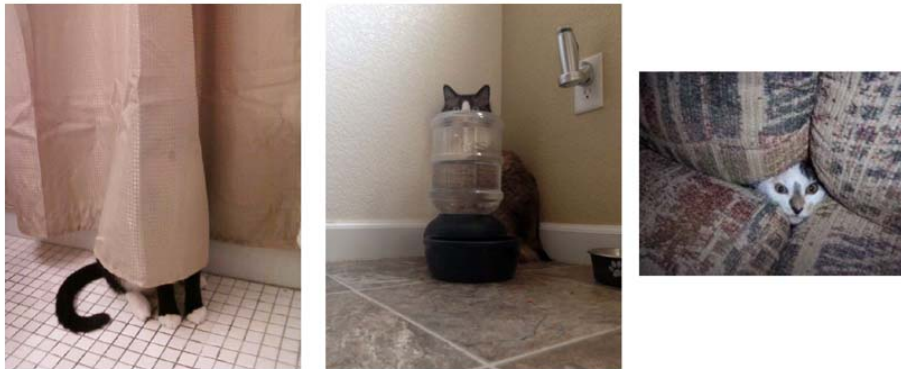


© 2018 Ray Ptucha, Rochester Institute of Technology

25

The Brain on Pattern Recognition

Despite Changes in Occlusion:

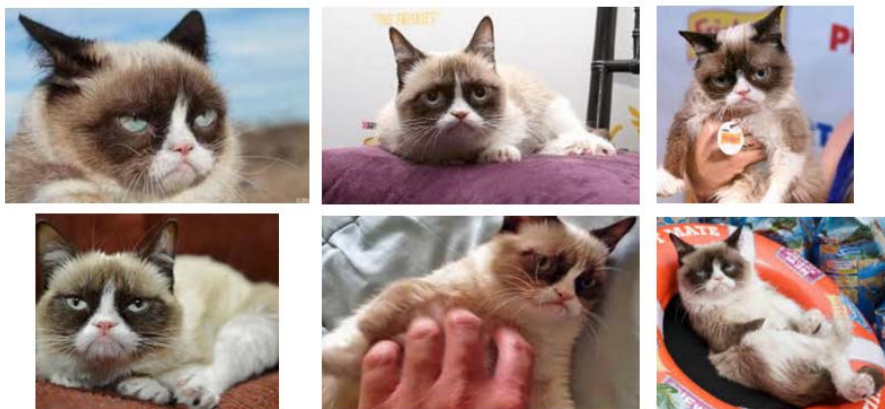


© 2018 Ray Ptucha, Rochester Institute of Technology

26

The Brain on Pattern Recognition

Despite Changes in Size, Pose, Angle:



Tardar Sauce "Grumpy Cat"

© 2018 Ray Ptucha, Rochester Institute of Technology

27

The Brain on Pattern Recognition

Despite Changes in Background Clutter:



© 2018 Ray Ptucha, Rochester Institute of Technology

28

The Brain on Pattern Recognition

Despite Changes in Class Variation...

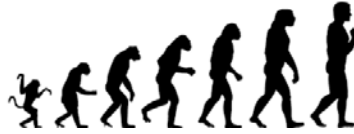


© 2018 Ray Ptucha, Rochester Institute of Technology

29

Teaching Computers to See

- It took evolution 540M years to develop the marvel of the eye-brain.
- Lets say a child collects a new image every 200msec.
- By age 3, this child has processed ~250M images.



$(5 \text{ images/sec})(60 \text{ sec/min})(60 \text{ min/hr})(12 \text{ hr/day})(365 \text{ days/yr})(3 \text{ yrs}) = 236 \text{ M}$

- Today's computers can do this in a few days...

© 2018 Ray Ptucha, Rochester Institute of Technology

30

Neural Nets on Pattern Recognition

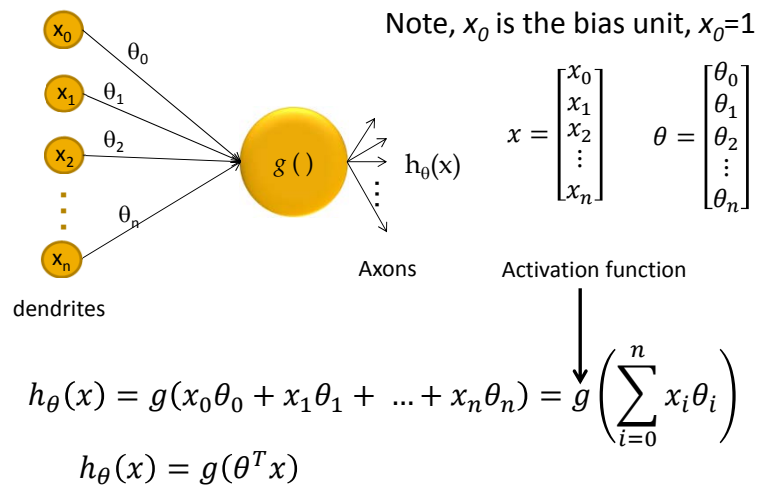
- Instead of trying to code simple intuitions/rules on what makes an airplane, car, cat, and dog...
- We feed neural networks a large number of training samples, and it will automatically learn the rules!
- We will learn the magic behind this today!



© 2018 Ray Ptucha, Rochester Institute of Technology

32

Artificial Neuron

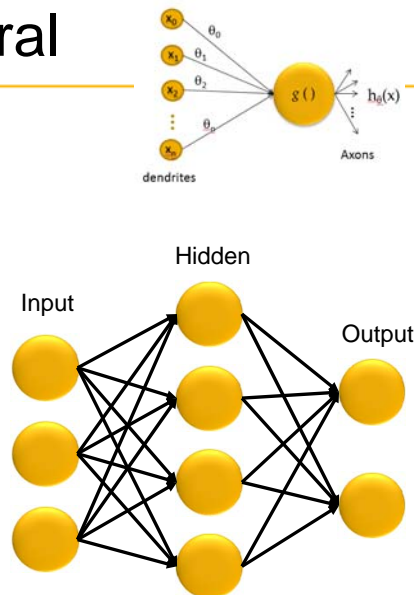


© 2018 Ray Ptucha, Rochester Institute of Technology

33

Artificial Neural Networks

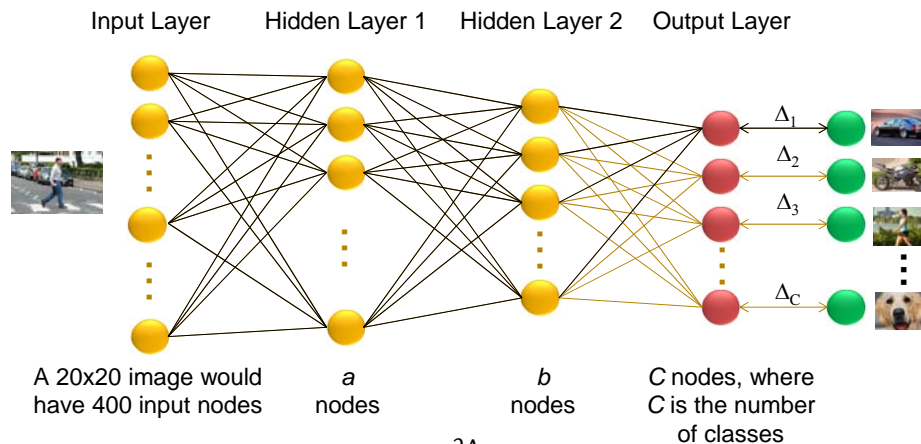
- Artificial Neural Network (ANN) –
A network of interconnected nodes that “mimic” the properties of a biological network of neurons



© 2018 Ray Ptucha, Rochester Institute of Technology

34

4-Layer ANN Fully Connected Topology

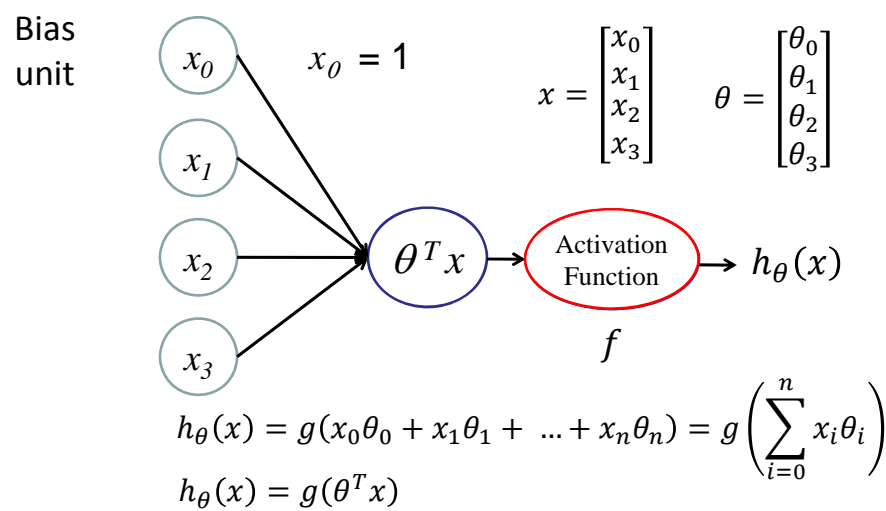


Backpropagation (~1985) uses $\frac{\partial \Delta}{\partial w}$ for learning
 Learning happens in the weights- each line is a weight.

© 2018 Ray Ptucha, Rochester Institute of Technology

35

Neuron Model



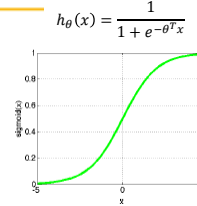
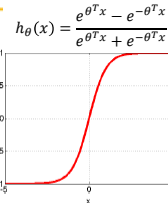
© 2018 Ray Ptucha, Rochester Institute of Technology

36

Activation Function Comparison

- Tanh
- Sigmoid

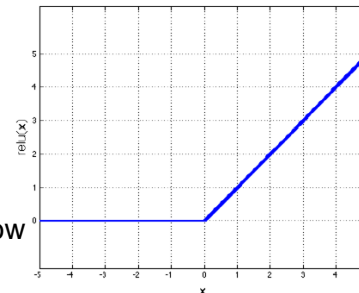
Gradient of both saturates at zero. Sigmoid also non-zero centered, so in practice tanh performs better.



- Rectified Linear Units (ReLU)

- Better for high dynamic range
- Faster learning
- Overall better result
- Neurons can “die” if allowed to grow unconstrained

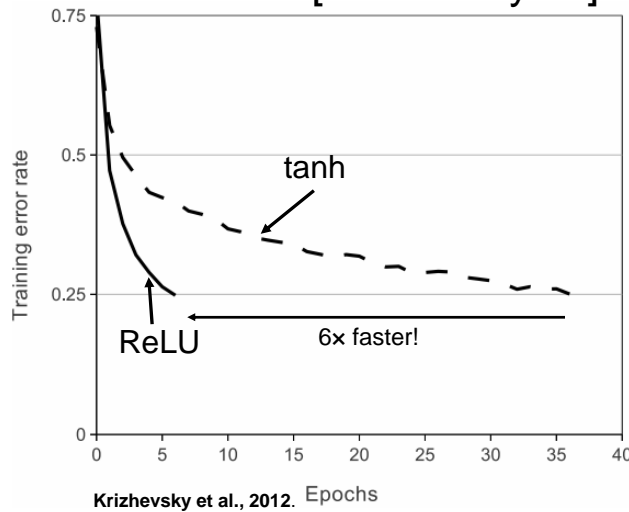
$$h_{\theta}(x) = \max(0, x)$$



© 2018 Ray Ptucha, Rochester Institute of Technology

37

Tanh vs. ReLU on CIFAR-10 dataset [Krizhevsky'12]



ReLU reaches 25% error 6x faster!
 Note: Learning rates optimized for each, no regularization, four layer CNN.

© 2018 Ray Ptucha, Rochester Institute of Technology

38

Where Do Weights Come From?

- The weights in a neural network need to be learned such that the errors are minimized.
- Just like logistic regression, we can write a cost function.
- Similar to gradient descent, we can write an iterative procedure to update weights, with each iteration decreasing our cost.
- These iterative methods may be less efficient than a direct analytical solution, but are easier to generalize.

© 2018 Ray Ptucha, Rochester Institute of Technology

39

Backpropagation

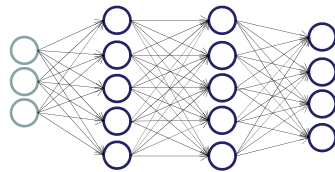
- We need to solve weights of a network so that the error is minimized.
- Weights can be refined by changing each weight by an amount proportional to the partial derivative of the error with respect to each weight.
- Partial derivatives can be calculated by iteratively changing each weight and measuring the corresponding change in error.
- Hard to do with millions of weights!
- In 1986, a technique called back-propagation was introduced (D. E. Rumelhart, G. E. Hinton, and R. J. Williams "Learning representations by back-propagating errors," *J. Nature* 323, 533-536, 1986).

© 2018 Ray Ptucha, Rochester Institute of Technology

40

Hyperparameters vs. parameters

- Hyperparameters are the tuning parameters for a nnet- say number of layers, nodes per layer, learning rate, momentum, regularization, etc.
- Parameters are the weights that are being learned. Ignoring bias terms, the below network has $3 \times 5 + 5 \times 5 + 5 \times 4 = 60$ parameters.
 - If we include bias terms, we have $4 \times 5 + 6 \times 5 + 6 \times 4 = 74$ learnable parameters.



Note: deep nets may contain 100M parameters with 100 layers!

© 2018 Ray Ptucha, Rochester Institute of Technology

41

Initialization of Weights

- We initialize weights to small \pm values centered on 0.
- If they were all initialized to 0, the network wouldn't learn anything.
 - ie: the output of each node would be equal, therefore the update would update each term identically.
- If they were all initialized to large values (values $> +1$ or < -1), the inputs to each node would be saturated.
- If they were all initialized to small values around 0, we would be in the linear portion of the activation function.
 - $W1 = 0.001 * \text{rand}(\text{length}(h1), \text{length}(h2))$

© 2018 Ray Ptucha, Rochester Institute of Technology

43

Initialization of Weights

- Although the uniform distribution is good, the more inputs to a node, the greater its variance.
- To set output distribution of all nodes equal (this empirically improves convergence), use

Transfer Learning

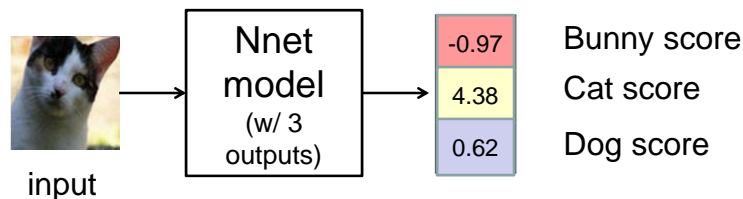
h1=#
input
nodes

- Other solutions restrict weights: $-\epsilon_{init} < \theta_{ji}^{(l)} < \epsilon_{init}$
 $\epsilon_{init} = \sqrt{6} / \sqrt{s_l + s_{l+1}}$ s_l, s_{l+1} are the No. nodes in layers around $\theta^{(l)}$
- For deep ReLU networks, He2015 showed:
 - **W1 = 0.001*rand(length(h1),length(h2))./sqrt(2*length(h1))**
- Works best and is recommended for them

© 2018 Ray Ptucha, Rochester Institute of Technology

44

Multiclass Loss Functions



- The input image scores highest against cat, but is also somewhat similar to dog.
- How do we assign a loss function?

© 2018 Ray Ptucha, Rochester Institute of Technology

45

Most Common Loss Functions

- The cost function we previously used was a direct copy from logistic regression and works great for binary classification.

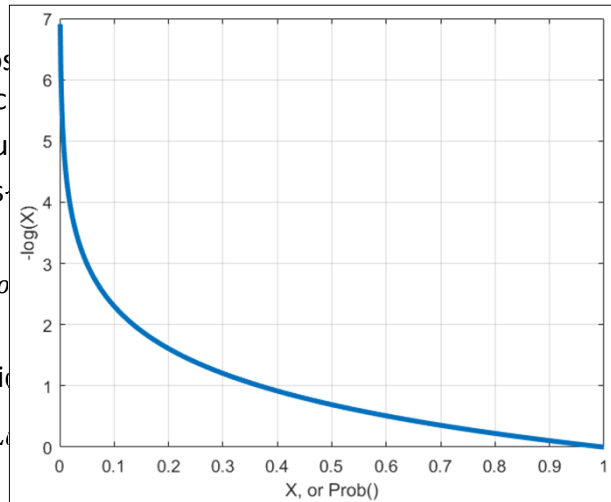
- For multi-class, there are two popular data loss methods:

1. Cross-entropy loss, which uses softmax:

$$Loss^{(i)} = -\log\left(\frac{\exp(out_{y_i}^{(i)})}{\sum_{c=1:C} \exp(out_c^{(i)})}\right)$$

2. Multiclass SVM Loss (Weston Watkins formulation):

$$Loss^{(i)} = \sum_{j \neq y_i} \max(0, out_j - out_{y_i} + \Delta)$$



© 2018 Ray Ptucha, Rochester Institute of Technology

48

Most Common Loss Functions

- The cost function we previously used was a direct copy from logistic regression and works great for binary classification.

- For multi-class, there are two popular data loss methods:

1. Cross-entropy loss, which uses softmax:

$$Loss^{(i)} = -\log\left(\frac{\exp(out_{y_i}^{(i)})}{\sum_{c=1:C} \exp(out_c^{(i)})}\right) \quad \text{Loss for sample } i = \frac{\exp(\text{output of GT node})}{\text{Sum of exp(output) of all nodes}}$$

2. Multiclass SVM Loss (Weston Watkins formulation):

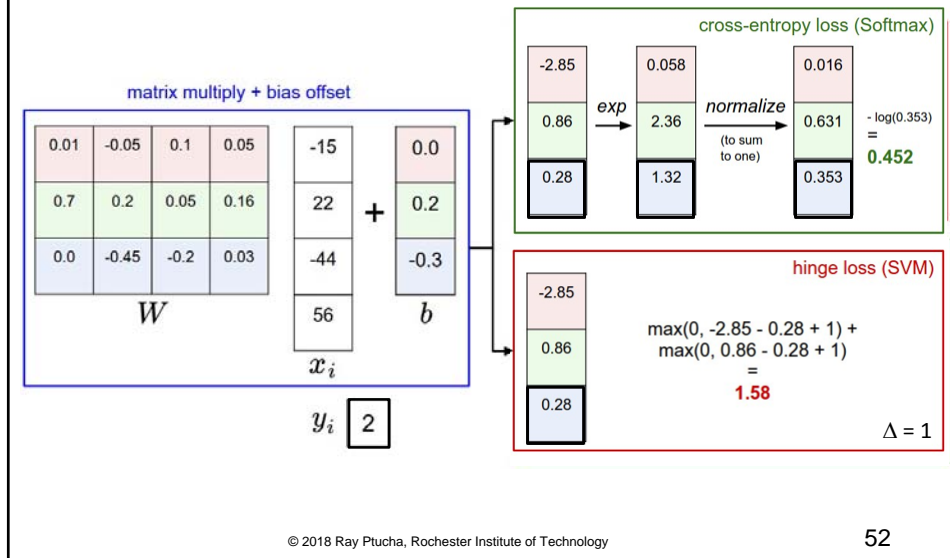
$$Loss^{(i)} = \sum_{j \neq y_i} \max(0, out_j - out_{y_i} + \Delta) \quad \text{Sum of incorrect - correct classes}$$

© 2018 Ray Ptucha, Rochester Institute of Technology

49

Loss Example

(based on cs231n, Li/Karpathy 2016)



52

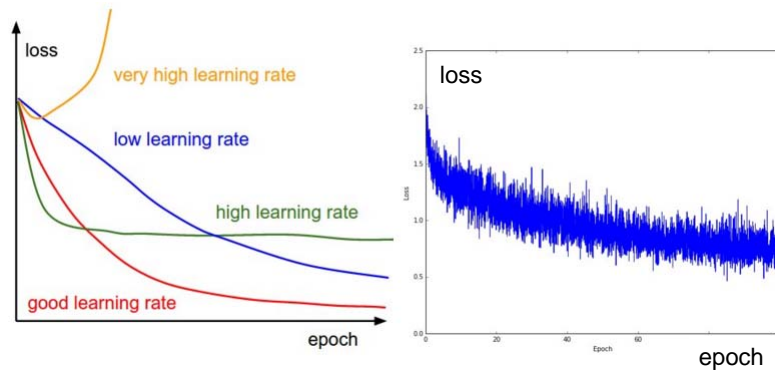
Sequential vs. Batch Training

- Back propagation can be done:
 - Sequential: one sample at a time,
 - Weights are shifted back and forth quite a bit
 - Group (minibatch): a group of samples at a time, or
 - Batch: all training samples at once
 - Weights are shifted in direction that makes most input samples better
 - Generally converges the fastest
- Recommended to use largest minibatch possible and stay within memory constraints of hardware.

© 2018 Ray Ptucha, Rochester Institute of Technology

53

Examples of Learning Rate and Batch Size

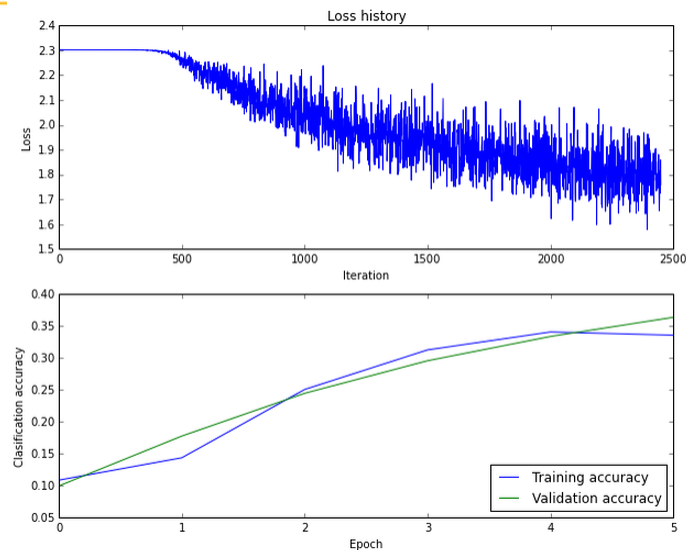


This batch size could be made a little larger to shrink the variance

© 2018 Ray Ptucha, Rochester Institute of Technology

56

Tuning Parameters



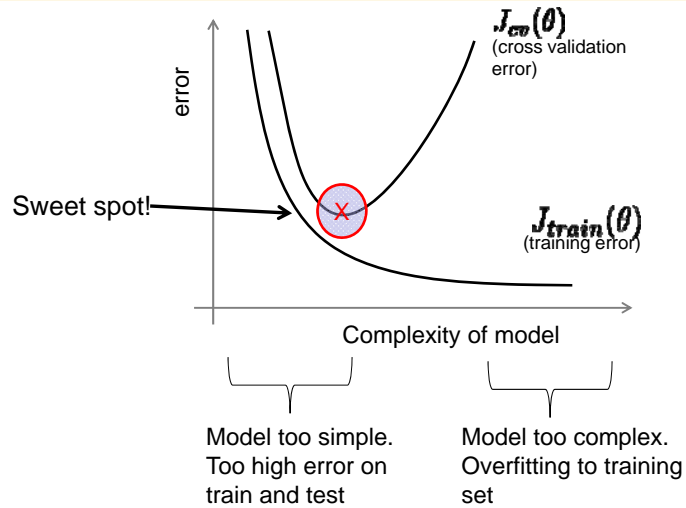
Curve too linear-increase learning rate.

Train and Test sets too similar-increase model complexity. Be careful, larger models susceptible to overfitting. Perhaps increase regularization

© 2018 Ray Ptucha, Rochester Institute of Technology

57

Bias (underfit) vs. Variance (overfit) errors

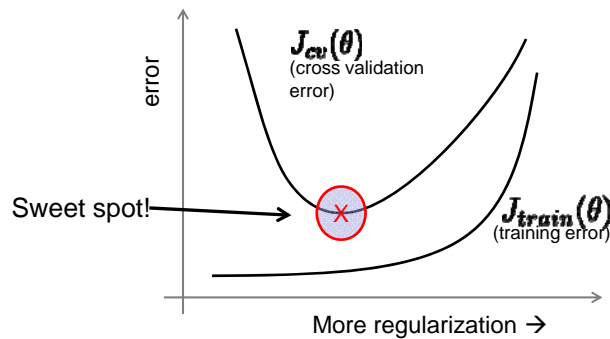


Adopted from: Andrew Ng, ML class

© 2018 Ray Ptucha, Rochester Institute of Technology

58

Regularization Tuning

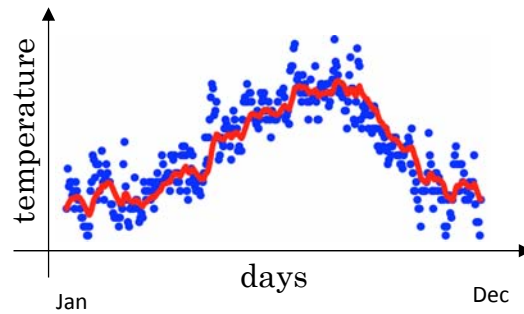


Adopted from: Andrew Ng, ML class

© 2018 Ray Ptucha, Rochester Institute of Technology

59

Exponentially Weight Average



Smoothing algorithm:

$$t_1 = 0.9t_0 + 0.1t_1$$

$$t_2 = 0.9t_1 + 0.1t_2$$

$$t_3 = 0.9t_2 + 0.1t_3$$

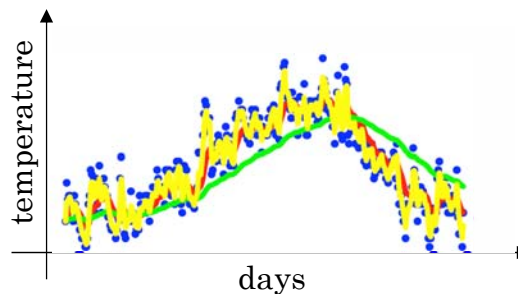
...

<https://www.coursera.org/learn/deep-neural-network/lecture/Ud7t0/understanding-exponentially-weighted-averages>

© 2018 Ray Ptucha, Rochester Institute of Technology

60

Exponentially Weight Average



Typically use v for velocity:

$$v_t = \beta v_{t-1} + (1 - \beta)v_t$$

Where β controls smoothness

$$avg \approx \frac{1}{(1 - \beta)} \text{readings}$$

$\beta = 0.9$ (~10 readings average)

$\beta = 0.98$ (~50 readings average)

$\beta = 0.5$ (~2 readings average)

Big latency, but smooth

Small latency, but overfit

- One big advantage of this method is that it takes up minimal memory

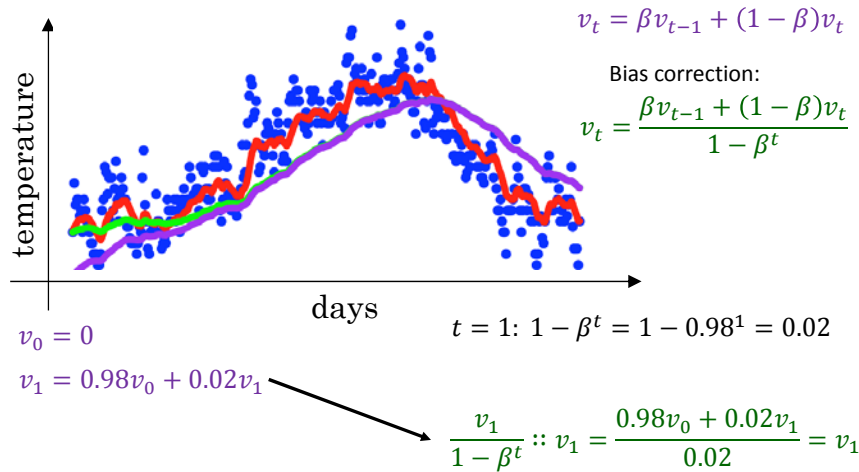
<https://www.coursera.org/learn/deep-neural-network/lecture/Ud7t0/understanding-exponentially-weighted-averages>

© 2018 Ray Ptucha, Rochester Institute of Technology

61

Exponentially Weight Average with Bias Correction

When $\beta=0.98$, we don't actually get green curve, we get purple curve if $v_0=0$

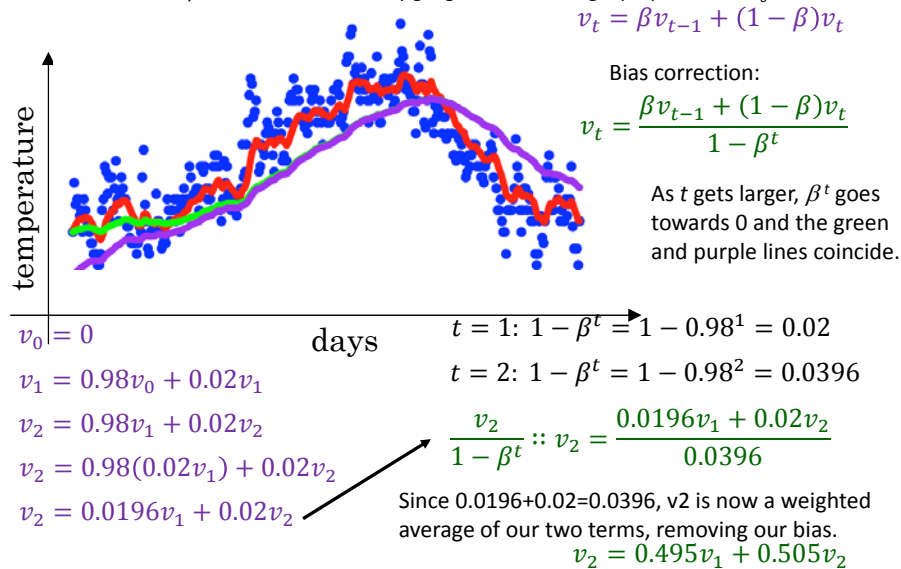


© 2018 Ray Ptucha, Rochester Institute of Technology

62

Exponentially Weight Average with Bias Correction

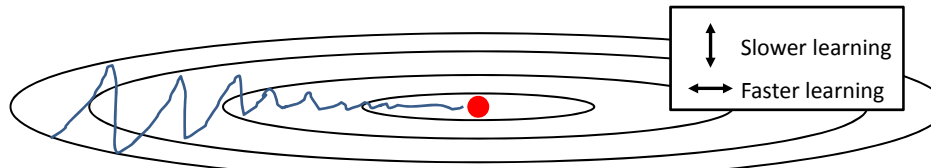
When $\beta=0.98$, we don't actually get green curve, we get purple curve if $v_0=0$



© 2018 Ray Ptucha, Rochester Institute of Technology

63

Gradient Descent with Momentum



- Typical gradient descent- oscillations slow down progress
- Further, if learning rate too large, initial step would cause significant overshooting, and perhaps divergence.

For each iteration, t :

Compute dW, db on current mini-batch

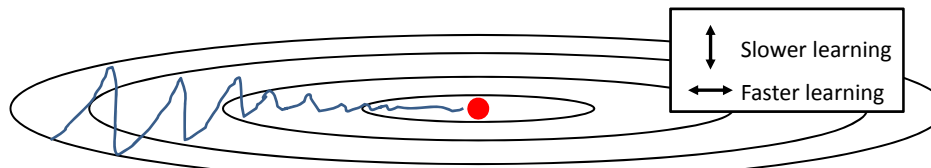
$$\left. \begin{aligned} V_{dW} &= \beta V_{dW} + (1 - \beta) dW \\ V_{db} &= \beta V_{db} + (1 - \beta) db \end{aligned} \right\} \text{Replace each partial derivative with a smoothed version}$$

$$\left. \begin{aligned} W &= W - \eta V_{dW} \\ b &= b - \eta V_{db} \end{aligned} \right\} \text{Update using smoothed derivatives}$$

© 2018 Ray Ptucha, Rochester Institute of Technology

64

Gradient Descent with Momentum



For each iteration, t :

Compute dW, db on current mini-batch

$$\left. \begin{aligned} V_{dW} &= \beta V_{dW} + (1 - \beta) dW \\ V_{db} &= \beta V_{db} + (1 - \beta) db \end{aligned} \right\} \text{Replace each partial derivative with a smoothed version}$$

$$\left. \begin{aligned} W &= W - \eta V_{dW} \\ b &= b - \eta V_{db} \end{aligned} \right\} \text{Update using smoothed derivatives}$$

We end up with two hyperparameters, η and β (learning rate and momentum)

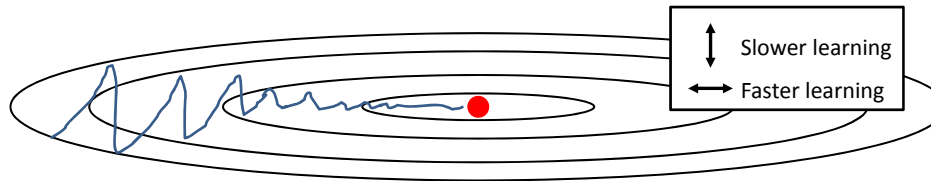
Very common to set $\beta=0.9$

Because moving average moves up quickly, bias correction not commonly used

© 2018 Ray Ptucha, Rochester Institute of Technology

65

Gradient Descent with Momentum



For each iteration, t :

Compute dW, db on current mini-batch

$$V_{dW} = \beta V_{dW} + (1 - \beta) dW \quad \text{Sometimes approximated with:} \quad V_{dW} = \beta V_{dW} + dW$$

$$V_{db} = \beta V_{db} + (1 - \beta) db \quad \text{Sometimes approximated with:} \quad V_{db} = \beta V_{db} + db$$

$$W = W - \eta V_{dW}$$

$$b = b - \eta V_{db}$$

We end up with two hyperparameters, η and β (learning rate and momentum)

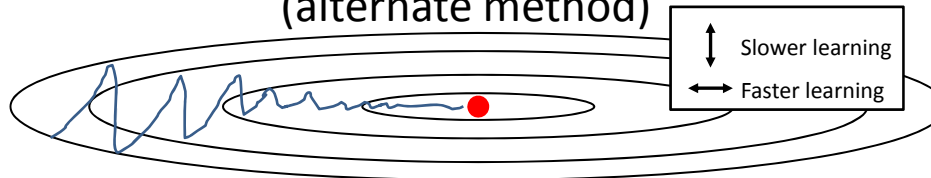
Very common to set $\beta=0.9$

Because moving average moves up quickly, bias correction not commonly used

© 2018 Ray Ptucha, Rochester Institute of Technology

66

Gradient Descent with Momentum (alternate method)



For each iteration, t :

Compute dW, db on current mini-batch

$$\left. \begin{aligned} V_{dW} &= \beta V_{dW} - \eta dW \\ V_{db} &= \beta V_{db} - \eta db \end{aligned} \right\} \begin{array}{l} \text{Smoothed version of } dW, \text{ but} \\ \text{include learning rate to increase} \\ \text{momentum more.} \end{array}$$

$$\left. \begin{aligned} W &= W + V_{dW} \\ b &= b + V_{db} \end{aligned} \right\} \begin{array}{l} \text{Since we were learning '}' dW \\ \text{above, we add, not subtract here.} \end{array}$$

We end up with two hyperparameters, η and β (learning rate and momentum)

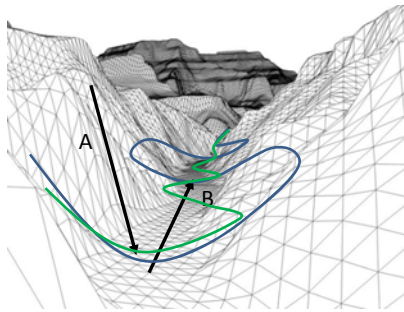
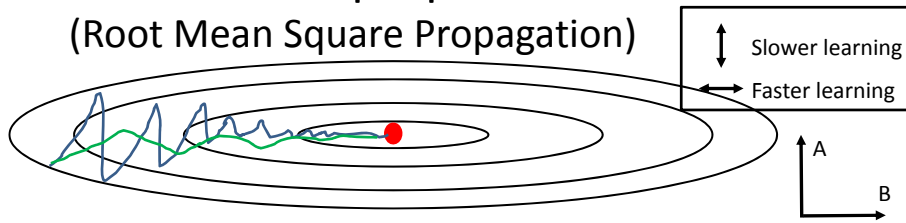
Very common to set $\beta=0.9$

Because moving average moves up quickly, bias correction not commonly used

© 2018 Ray Ptucha, Rochester Institute of Technology

67

RMSprop (Root Mean Square Propagation)

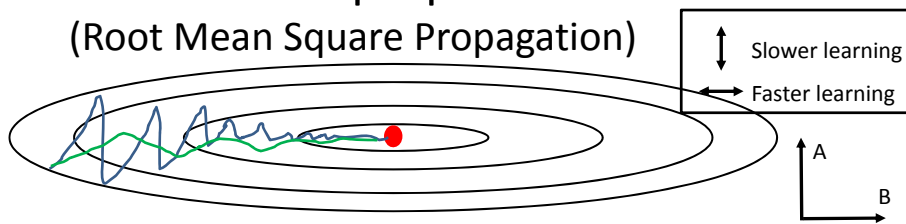


- In general the derivative in direction A will be large and B will be small.
- Weighted derivative squares in A/B will be large/small.
- When A/B large, want to dampen
- When A/B small, don't dampen as much

© 2018 Ray Ptucha, Rochester Institute of Technology

68

RMSprop (Root Mean Square Propagation)



For each iteration, t :

Compute dW , db on current mini-batch

$$\begin{aligned} S_{dW} &= \beta S_{dW} + (1 - \beta) dW * dW \\ S_{db} &= \beta S_{db} + (1 - \beta) db * db \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \begin{array}{l} \text{Keep track of exponentially} \\ \text{weighted squares of} \\ \text{derivatives} \end{array}$$

$$W = W - \eta \frac{dW}{\sqrt{S_{dW}}}$$

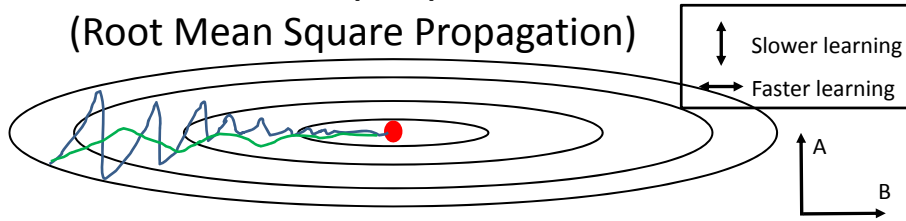
$$b = b - \eta \frac{db}{\sqrt{S_{db}}}$$

- In general the derivative in direction A will be large and B will be small.
- Weighted squares in A/B will be large/small.
- When A/B large, want to dampen
- When A/B small, don't dampen as much

© 2018 Ray Ptucha, Rochester Institute of Technology

69

RMSprop (Root Mean Square Propagation)



For each iteration, t :

Compute dW, db on current mini-batch

$$\begin{aligned} S_{dW} &= \beta_2 S_{dW} + (1 - \beta_2) dW .* dW \\ S_{db} &= \beta_2 S_{db} + (1 - \beta_2) db .* db \end{aligned} \quad \left. \begin{array}{l} \\ \end{array} \right\} \begin{array}{l} \text{Keep track of exponentially} \\ \text{weighted squares of} \\ \text{derivatives} \end{array}$$

$$W = W - \eta \frac{dW}{\sqrt{S_{dW} + \epsilon}}$$

$$b = b - \eta \frac{db}{\sqrt{S_{db} + \epsilon}}$$

- Next we will combine RMSprop with Momentum so lets rename β_2 to β_2
- Also, lets make sure we don't divide by zero

© 2018 Ray Ptucha, Rochester Institute of Technology

70

Adam Optimization

Combine momentum with RMSprop with Bias Correction

$$V_{dW}=0, S_{dW}=0, V_{db}=0, S_{db}=0$$

For each iteration, t :

Compute dW, db on current mini-batch

$$\begin{aligned} \text{momentum} \quad & \begin{cases} V_{dW} = \beta_1 V_{dW} + (1 - \beta_1) dW \\ V_{db} = \beta_1 V_{db} + (1 - \beta_1) db \end{cases} \\ \text{RMSprop} \quad & \begin{cases} S_{dW} = \beta_2 S_{dW} + (1 - \beta_2) dW .* dW \\ S_{db} = \beta_2 S_{db} + (1 - \beta_2) db .* db \end{cases} \end{aligned}$$

Hyperparameters to tune:

β_1 : 0.9
 β_2 : 0.999
 ϵ : 10^{-8}
 η : need to solve

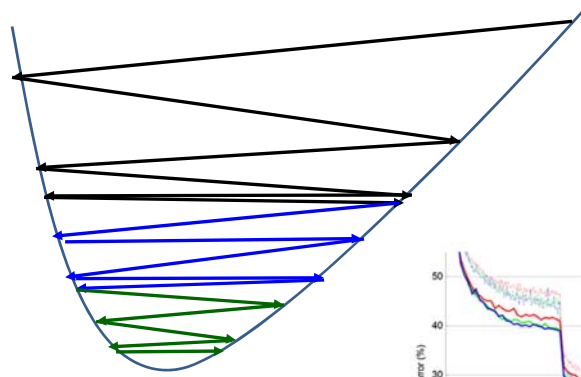
$$\begin{aligned} \text{Bias Correction} \quad & \begin{cases} V_{dW}^{\text{Corrected}} = V_{dW} / (1 - \beta_1^t); & V_{db}^{\text{Corrected}} = V_{db} / (1 - \beta_1^t) \\ S_{dW}^{\text{Corrected}} = S_{dW} / (1 - \beta_2^t); & S_{db}^{\text{Corrected}} = S_{db} / (1 - \beta_2^t) \end{cases} \end{aligned}$$

$$W = W - \eta \left(V_{dW}^{\text{Corrected}} / \sqrt{S_{dW}^{\text{Corrected}} + \epsilon} \right); \quad b = b - \eta \left(V_{db}^{\text{Corrected}} / \sqrt{S_{db}^{\text{Corrected}} + \epsilon} \right)$$

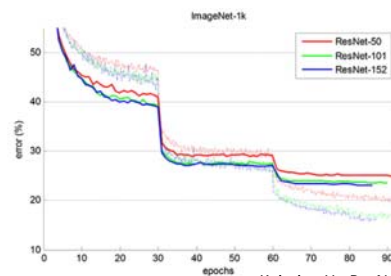
© 2018 Ray Ptucha, Rochester Institute of Technology

71

Learning Rate Decay



- If shrink step size, we can find a better solution
- If shrink step size, we can find a better solution



© 2018 Ray Ptucha, Rochester Institute of Technology

Kaiming He ResNet 72

Learning Rate Decay

$$\eta = \frac{1}{1 + \text{decayRate} * \text{epochNum}} \eta_0$$

Initial learning rate \nearrow

$$\eta_0 = 0.2$$

$$\text{decayRate} = 1$$

epochNum	η
1	0.1
2	0.067
3	0.05
4	0.04
...	...

© 2018 Ray Ptucha, Rochester Institute of Technology

73

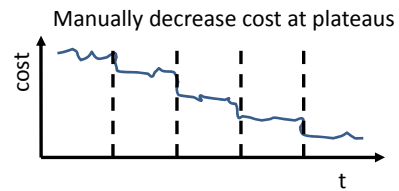
Learning Rate Decay- Other formulas

Some hyperparameter <1

$$\eta = 0.95^{\text{epochNum}} \cdot \eta_0 \quad (\text{Exponential decay})$$

$$\eta = \frac{k}{\sqrt{\text{epochNum}}} \cdot \eta_0 \quad \text{--OR--} \quad \eta = \frac{k}{\sqrt{t}} \cdot \eta_0 \quad t \text{ can be custom to problem}$$

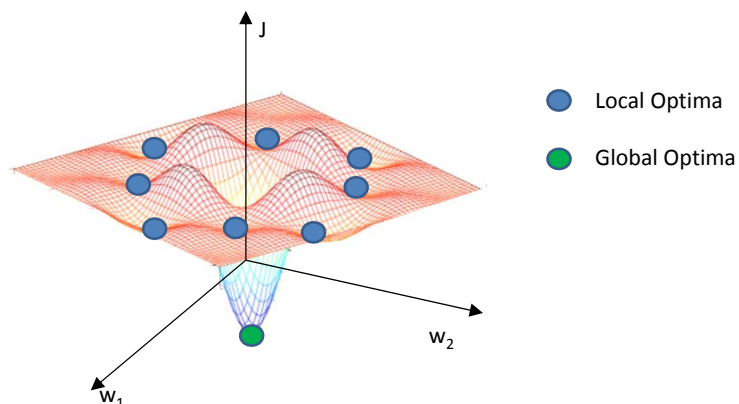
Discrete staircase:
Divide learning rate in half every
100 iterations



© 2018 Ray Ptucha, Rochester Institute of Technology

74

Local Optima in Neural Networks

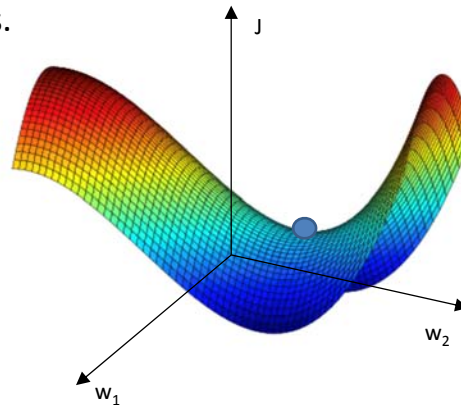


© 2018 Ray Ptucha, Rochester Institute of Technology

75

Local Optima in Neural Networks

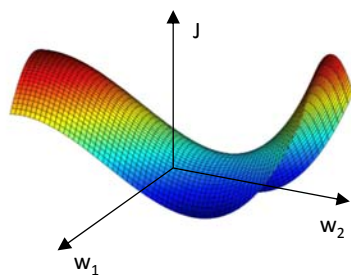
- As luck would have it, most points with gradient zero are not like previous slide, but occur on saddle points.



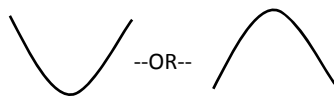
© 2018 Ray Ptucha, Rochester Institute of Technology

76

Local Optima in Neural Networks



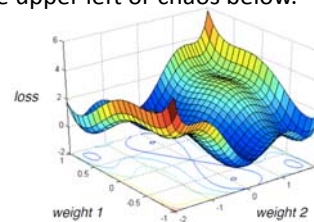
For a gradient to be zero, in each direction we either have a concave or a convex shape.



Now...in a very high dimensional space, say 10,000 dimensions, what are the odds all will be concave?

Pretty small!

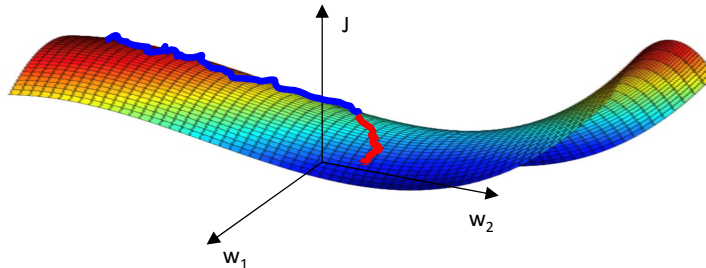
We are much more likely to have some directions run up, and some run down, like the saddle point on the upper left or chaos below.



© 2018 Ray Ptucha, Rochester Institute of Technology

77

Problem of Plateaus



- Given a surface like this, we traverse the cost function for a long time making little progress.
- Eventually, or perhaps by luck, we step to the side and finally begin to learn more aggressively again.
- But...thankfully in high dimensional spaces, unlikely to get stuck in local optima

© 2018 Ray Ptucha, Rochester Institute of Technology

78

Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift

Sergey Ioffe
Christian Szegedy

Google, 1600 Amphitheatre Pkwy, Mountain View, CA 94043

SIOFFE@GOOGLE.COM
SZEGEDY@GOOGLE.COM

Abstract

Training Deep Neural Networks is complicated by the fact that the distribution of each layer's inputs changes during training, as the parameters of the previous layers change. This slows down the training by requiring lower learning rates and careful parameter initialization, and makes it notoriously hard to train models with saturating nonlinearities. We refer to this phenomenon as *internal covariate shift*, and address the problem by normalizing layer inputs. Our method draws its strength from making normalization a part of the model architecture and performing the normalization for each training

minimize the loss

$$\Theta = \arg \min_{\Theta} \frac{1}{N} \sum_{i=1}^N \ell(x_i, \Theta)$$

where $x_{1..N}$ is the training data set. With SGD, the training proceeds in steps, at each step considering a *mini-batch* $x_{1..m}$ of size m . Using mini-batches of examples, as opposed to one example at a time, is helpful in several ways. First, the gradient of the loss over a mini-batch $\frac{1}{m} \sum_{i=1}^m \frac{\partial \ell(x_i, \Theta)}{\partial \Theta}$ is an estimate of the gradient over the training set, whose quality improves as the batch size increases. Second, computation over a mini-batch can be more efficient than m computations for individual examples on modern computing platforms.

<http://jmlr.org/proceedings/papers/v37/ioffe15.pdf>

© 2018 Ray Ptucha, Rochester Institute of Technology

79

Batch Normalization

- Makes hyperparameter search easier
- Makes learning more robust to the choice of hyperparameters
- Makes learning work with a larger range of hyperparameters
- Enables easier and often faster training of very deep networks

© 2018 Ray Ptucha, Rochester Institute of Technology

80

Normalization of Features

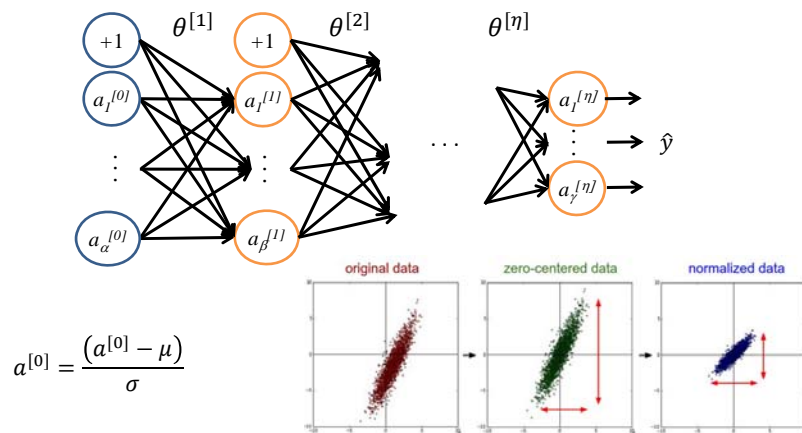
- We have seen how it is beneficial to subtract the mean and divide by the standard deviation for our input vector X (note, each dimension has its own μ and σ).
- In a deeper net, it would make the training of $\theta^{[l]}$ more efficient if the input to the layer, $a^{[l-1]}$ was also normalized.
- Note, in batch normalization, we will actually normalize z , not a . There is some disagreement on this, but most normalize z .

© 2018 Ray Ptucha, Rochester Institute of Technology

81

Implementing Batch Normalization

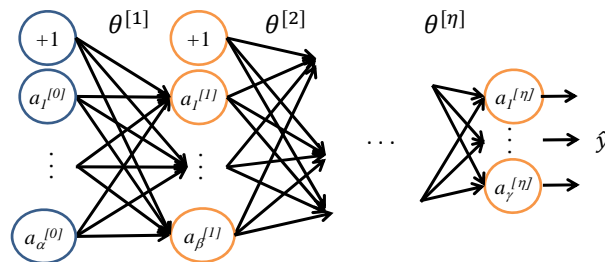
- Analogous to the way we normalized the input to $\theta^{[1]}$, normalize the input to all layers:



© 2018 Ray Ptucha, Rochester Institute of Technology

82

Implementing Batch Normalization



- As we then pass training samples through the net, replace $z^{[l]}$ with $z_{norm}^{[l]}$:

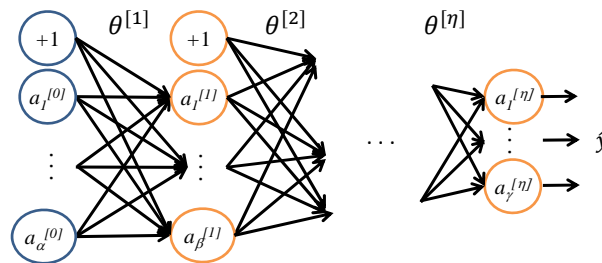
$$z_{norm}^{[l]} = \frac{(z^{[l]} - \mu^{[l]})}{\sqrt{\sigma^{[l]2} + \epsilon}}$$

- Normalize z , not a
- Calculate μ and σ^2 from mini-batch training samples:
- $\mu^{[l]} = \frac{1}{n} \sum_i z^{[l](i)}$
- $(\sigma^{[l]})^2 = \frac{1}{n} \sum_i (z^{[l](i)} - \mu^{[l]})^2$

© 2018 Ray Ptucha, Rochester Institute of Technology

83

Implementing Batch Normalization



$$z_{norm}^{[l]} = \frac{(z^{[l]} - \mu^{[l]})}{\sqrt{\sigma^{[l]2} + \epsilon}}$$

$$\hat{z}^{[l]} = \gamma z_{norm}^{[l]} + \beta$$

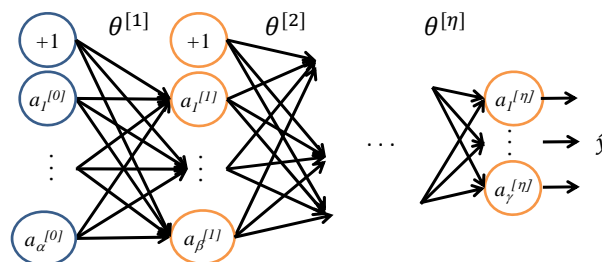
- Then go through activation function

- But, do we don't always want our distribution to have mean=0 and unit variance.
- So we learn two new parameters, γ and β , these will be learned along with our $\theta^{[l]}$'s

© 2018 Ray Ptucha, Rochester Institute of Technology

84

Implementing Batch Normalization



$$z_{norm}^{[l]} = \frac{(z^{[l]} - \mu^{[l]})}{\sqrt{\sigma^{[l]2} + \epsilon}}$$

$$\hat{z}^{[l]} = \gamma^{[l]} z_{norm}^{[l]} + \beta^{[l]}$$

- Then go through activation function

- Note $\gamma^{[l]}$ and $\beta^{[l]}$ are vectors for each layer, and get applied independently to each dimension.
- If using batch normalization with $\beta^{[l]}$, no need to use bias values as they both do the same job.

© 2018 Ray Ptucha, Rochester Institute of Technology

85

Batch Normalization Intuition

- When training set statistics differ from test set statistics, we call this covariate shift.



blog.bigml.com

© 2018 Ray Ptucha, Rochester Institute of Technology

86

Batch Normalization Intuition

- When training set statistics differ from test set statistics, we call this covariate shift.
- Typically we would have to relearn the model if we have a covariate shift, but batch norm minimizes covariate shift.
- In effect, it forces the input to each layer to have a predictable distribution, $\mu=0$, $\sigma=1$, or whatever distribution the γ and β determine works best for that layer.
- This predictable distribution makes later layers somewhat independent of earlier layers, enabling robust behavior even in deep networks.

© 2018 Ray Ptucha, Rochester Institute of Technology

87

Batch Normalization at Test Time

- Using mini-batch training samples, we calculate μ and σ^2 for each layer:
- $\mu^{[l]} = \frac{1}{n} \sum_i z^{[l](i)}$
- $(\sigma^{[l]})^2 = \frac{1}{n} \sum_i (z^{[l](i)} - \mu^{[l]})^2$
- During test, if we were testing a mini-batch at once, we could use the test mini-batch to calculate μ and σ .
- Alternately, if only have one or a few samples you need an estimate for μ and σ .

$$z_{norm}^{[l]} = \frac{(z^{[l]} - \mu^{[l]})}{\sqrt{\sigma^{[l]2} + \epsilon}}$$

$$\hat{z}^{[l]} = \gamma z_{norm}^{[l]} + \beta$$

- Can compute during training using a full epoch.
- Can use an exponentially weighted average for each:

$$v_t = \beta v_{t-1} + (1 - \beta) v_t$$

Note: this β is from exponentially weighted average and not related to batch norm β !

© 2018 Ray Ptucha, Rochester Institute of Technology

90



Andrew Ng, 2017

<https://www.deeplearning.ai/>

Deep Learning Specialization, Five courses:

1. Neural Networks and Deep Learning
2. Improving Deep Neural Networks
3. Structured Machine Learning Projects
4. Convolutional Neural Networks
5. Sequence Models

© 2018 Ray Ptucha, Rochester Institute of Technology

91



CS231n: Convolutional Neural Networks for Visual Recognition



Fei-Fei Li



Justin Johnson

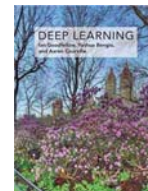


Serena Yeung

Li, Johnson, Yeung 2017

<http://cs231n.stanford.edu/>

Books on Deep Learning



- **[Grokking Deep Learning](https://www.manning.com/books/grokking-deep-learning)** by Andrew Trask. Use Udacity discount code **traskud17** for 40% off. This provides a very gentle introduction to Deep Learning and covers the intuition more than the theory.
 - <https://www.manning.com/books/grokking-deep-learning> (\$49)
- **[Neural Networks And Deep Learning](http://neuralnetworksanddeeplearning.com/)** by Michael Nielsen. This book is more rigorous than Grokking Deep Learning and includes a lot of fun, interactive visualizations to play with.
 - <http://neuralnetworksanddeeplearning.com/> (free!)
- **[The Deep Learning Textbook](http://www.deeplearningbook.org/)** from Ian Goodfellow, Yoshua Bengio, and Aaron Courville. This online book contains a lot of material and is the most rigorous of the three books suggested.
 - <http://www.deeplearningbook.org/> (html- free!)
 - <https://www.amazon.com/Deep-Learning-Adaptive-Computation-Machine/dp/0262035618/> (hard copy \$50)

Thank you!!

Ray Ptucha

rwpeec@rit.edu



<https://www.rit.edu/mil>