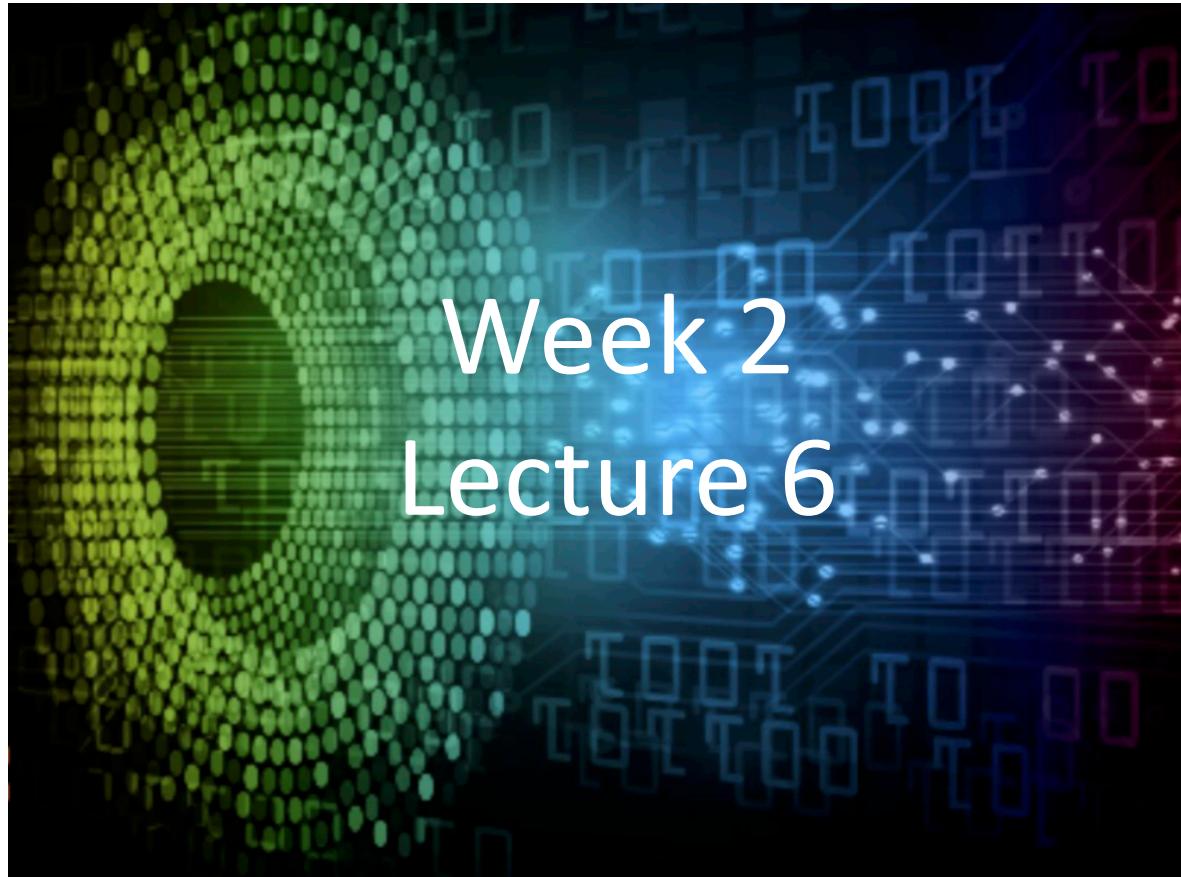


Introduction to Deep Learning Applications and Theory



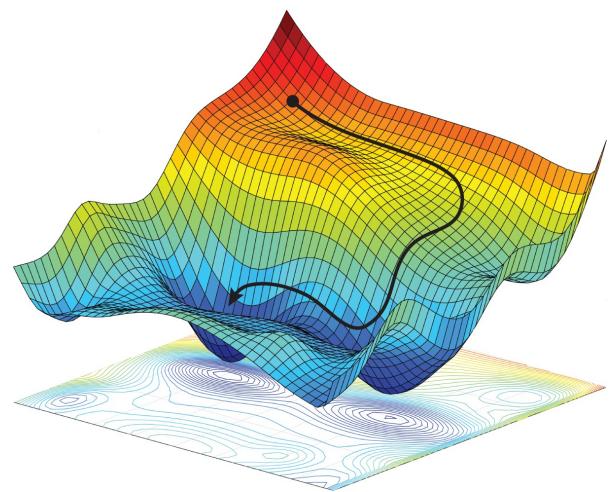
ECE 596 / AMATH 563

Previous Lecture: Optimization and Gradient Descent

1. Gradient Descent

2. Training the Logistic Regression Loss

3. Classification/ Regression with NN

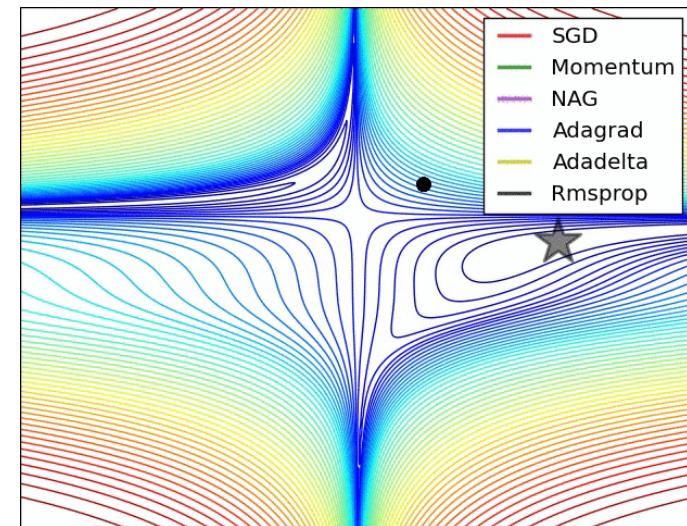


Current Lecture: Network Training/Optimization

1. Stochastic Gradient Descent

2. Extensions

3. Training NN



Gradient Descent

$$\vec{w}_{k+1} = \vec{w}_k - \alpha \nabla_{\vec{w}} J(\vec{w}_k; b)$$

learning rate



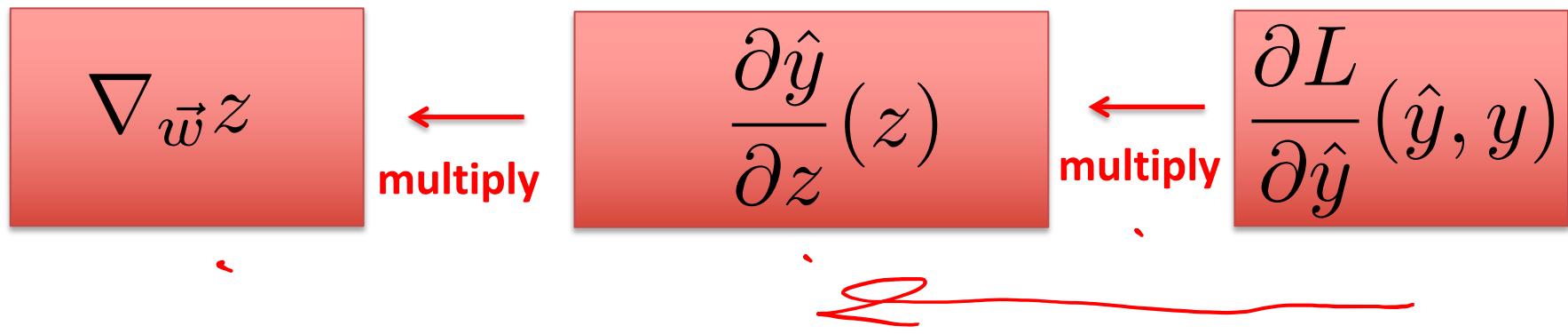
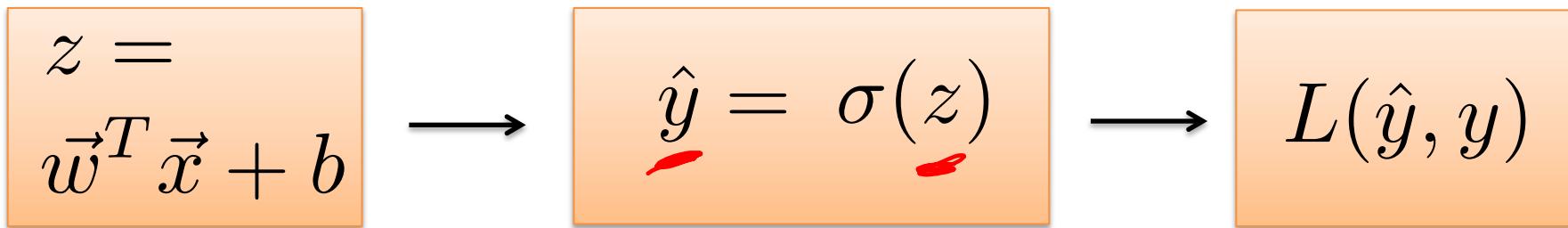
$$\vec{b}_{k+1} = \vec{b}_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

- Assume $i = 1$, and initial \vec{w}_0, b_0 *parameters*
- How do we compute \vec{w}_1, b_1 ?

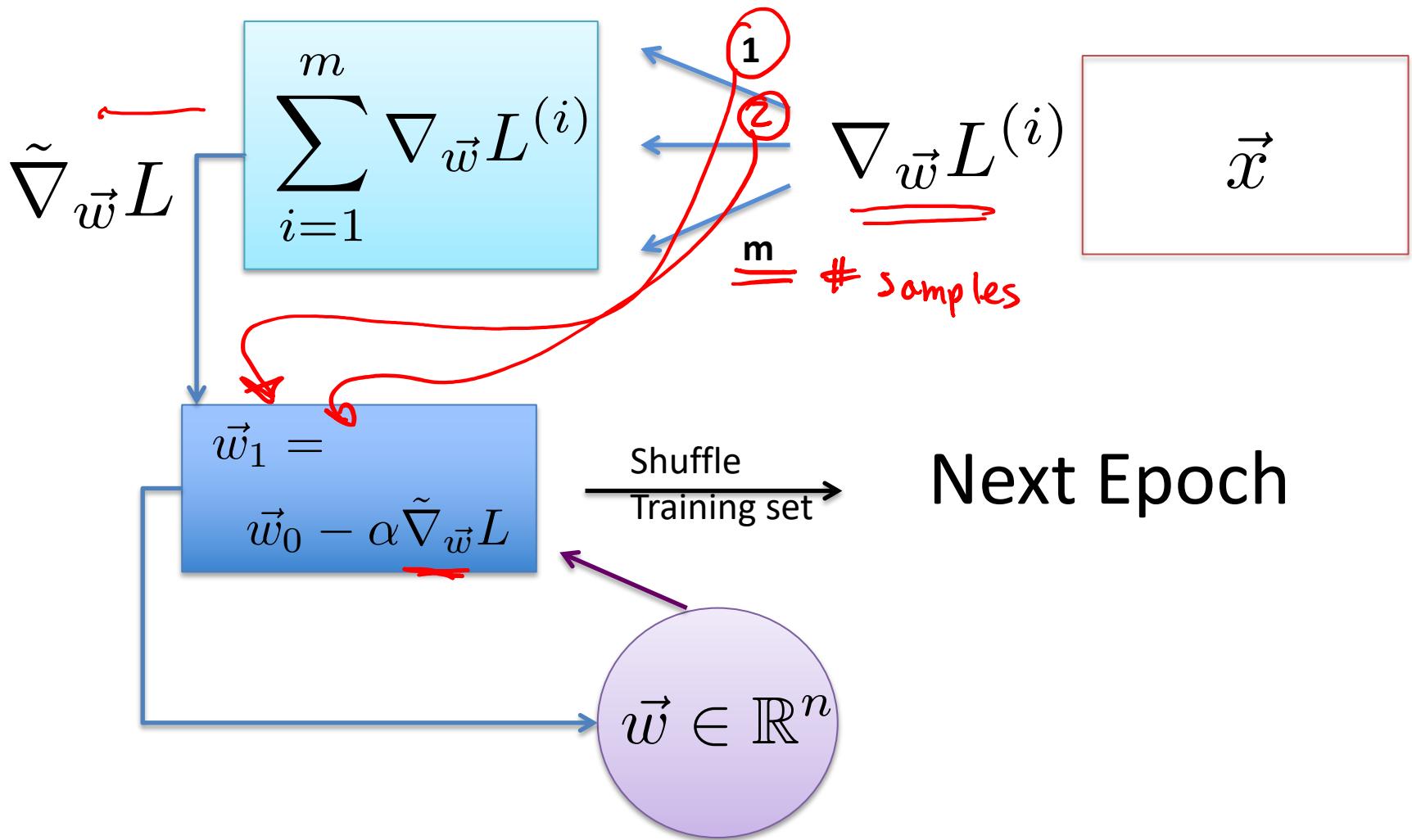
Our cost: $J = L(\hat{y}, y)$

Back Propagation – Chain Rule

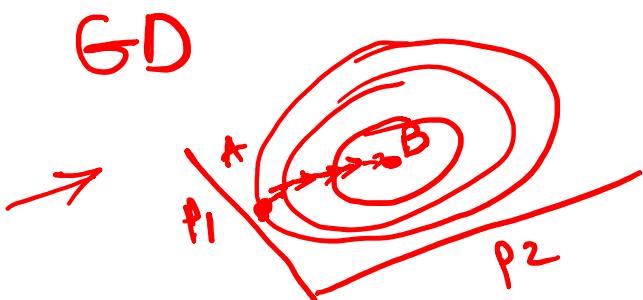
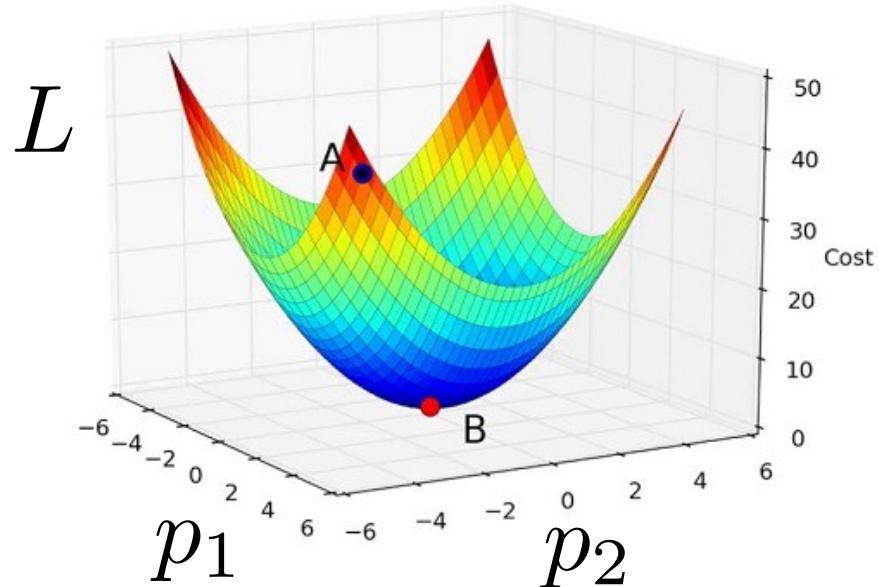
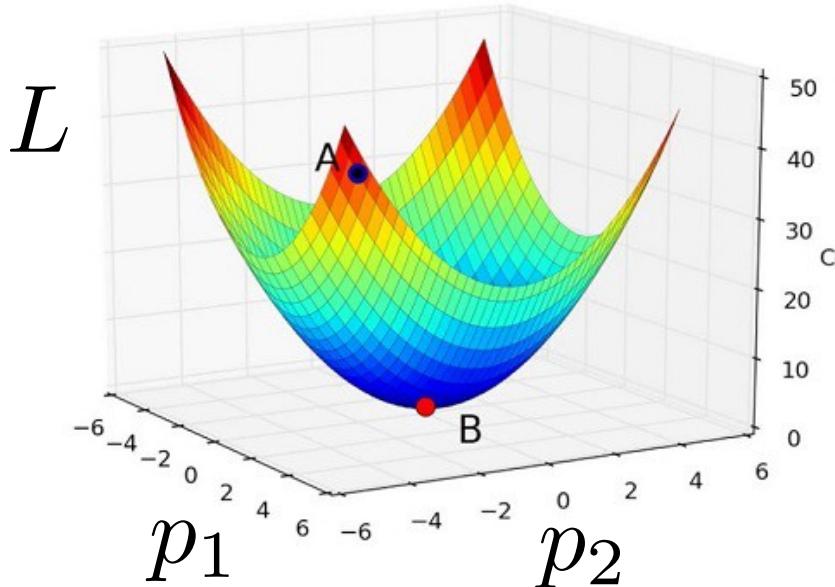
$$\boxed{\nabla_{\vec{w}} L(\hat{y}, y)} = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \nabla_{\vec{w}} z$$



Graph



SGD Convergence



Journal of Mathematical Analysis and Applications
Volume 114, Issue 2, March 1986, Pages 512-527



On convergence of the stochastic subgradient method with on-line stepsize rules

Andrzej Ruszczyński *, Wojciech Syski

Almost surely convergence

⇒ neurips 2020

Stochastic (Online) GD

$$w_{k+1} = w_k - \alpha \cdot \nabla_w J(w; x^{(i)}; y^{(i)})$$

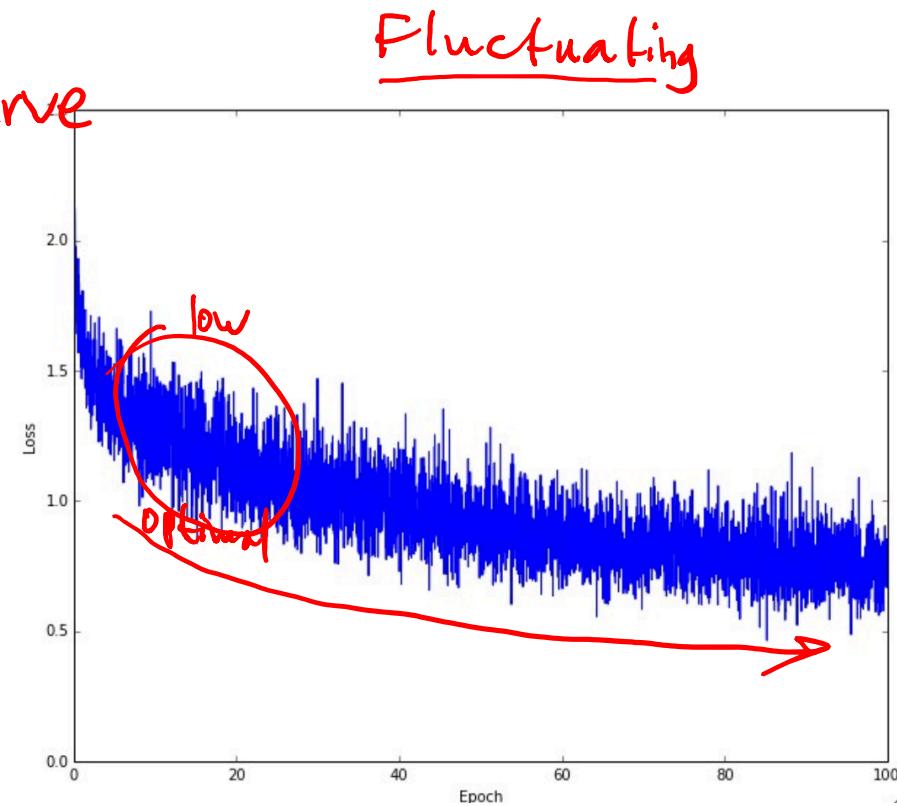
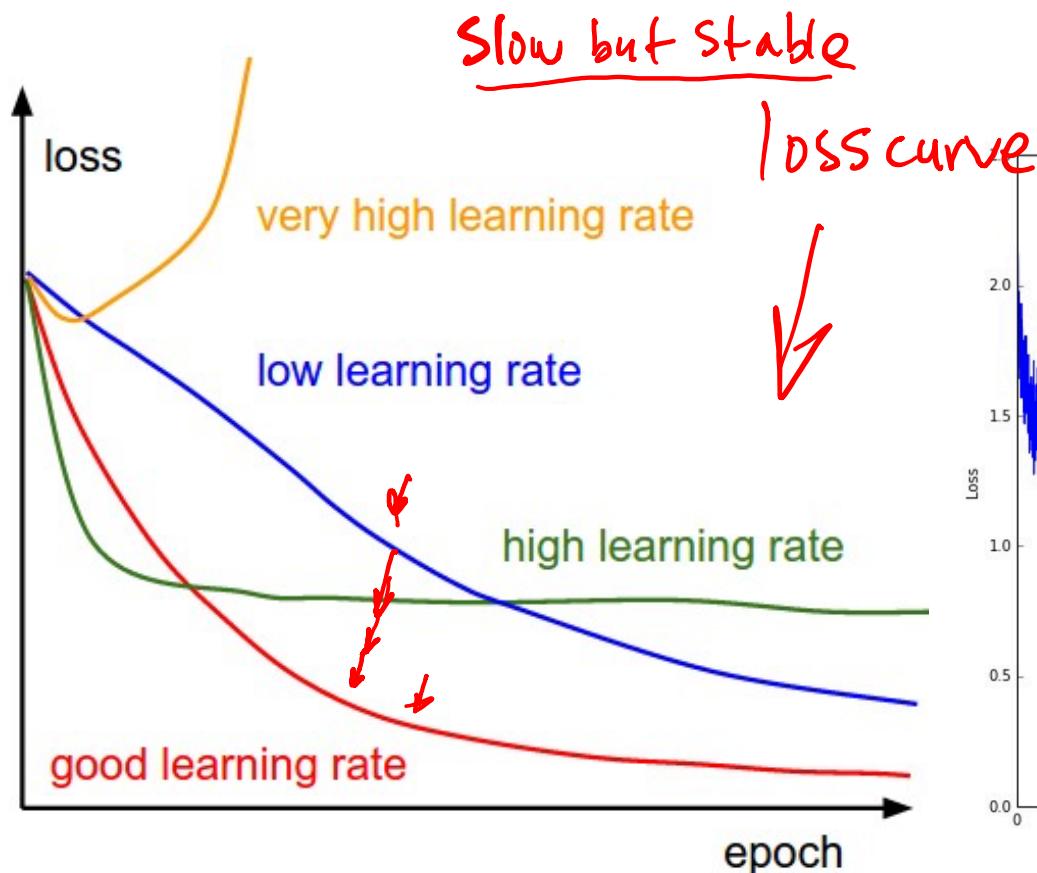
- Performs an update for **each training example $x^{(i)}$ and label $y^{(i)}$**
- The values of the loss and parameters will fluctuate.
 - (+) will discover better minimums
 - (-) convergence to chosen minimum will keep overshooting
- Learning rate plays a very important role

Need Metrics

```
training step: 4
training step: 5
training step: 6
training step: 7
training step: 8
training step: 9
training step: 10
training step: 11
training step: 12
training step: 13
training step: 14
training step: 15
training step: 16
training step: 17
training step: 18
training step: 19
training step: 20
training step: 21
training step: 22
training step: 23
training step: 24
training step: 25
training step: 26
training step: 27
training step: 28
training step: 29
training step: 30
training step: 31
training step: 32
training step: 33
training step: 34
training step: 35
training step: 36
training step: 37
training step: 38
training step: 39
training step: 40
training step: 41
training step: 42
training step: 43
training step: 44
training step: 45
training step: 46
training step: 47
training step: 48
training step: 49
~ >
```



Experimenting with GD



Hyperparameters = α

Mini-batch GD

$$w_{k+1} = w_k - \alpha \cdot \nabla_w J(w; x_{\underline{i:i+n}}^{(i:i+n)}; y_{\underline{i:i+n}}^{(i:i+n)})$$

- Mini-batch GD is a hybrid method between GD and SGD.
- Performs an update for every mini-batch of **n training examples**.

(+) reduces the variance of the parameter updates
(+) efficient in computing the gradient w.r.t. a mini-batch

- mini-batch sizes range between 50-256.

Challenges

extensions

- Choosing a proper learning rate can be difficult
- Learning rate smart schedule
 - Annealing
 - Change of J below threshold
- Variable learning for different parameters
- Suboptimal local (saddle points)



Comparing GD variants

SGD

mbsize = 1

(-) Can loose speedup from oscillations

(-) hard to parallelize

Mini-batch GD

$mbsize = m/R$

(+) The whole minibatch is evaluated in parallel
(+) Mostly consistent convergence

Batch GD

$mbsize = m$

(+) Consistent convergence

(-) Too long per iteration

Hyperparameters = mbsize (powers of 2, CPU/GPU memory)

Parameters

- **Parameters**

Model Parameters:

W, b, activation, output, cost

- **Hyper-parameters**



External Parameters:

Batch/minibatch size

Learning parameters

I. Learning rate decay

$$\tilde{\alpha} = \frac{1}{1 + \underbrace{\textit{decr}}_{\text{epnum}} \cdot \underbrace{\textit{epnum}}_{-}} \alpha_0$$

$$\underline{\alpha} = d^{\textit{epnum}} \cdot \alpha_0$$

$$\tilde{\alpha} = \frac{d}{\sqrt{\textit{epnum}}} \cdot \alpha_0$$

2. Momentum Method

$$\begin{aligned} \underline{\underline{v_{k+1}}} &= \underline{\gamma v_k} + \underline{\alpha \cdot \nabla_w J(w_k)} \\ \underline{\underline{w_{k+1}}} &= \underline{\underline{w_k}} - \underline{\underline{v_{k+1}}} \end{aligned}$$

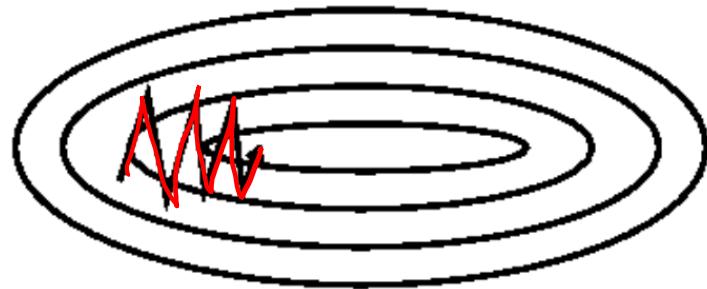


Image 2: SGD without momentum

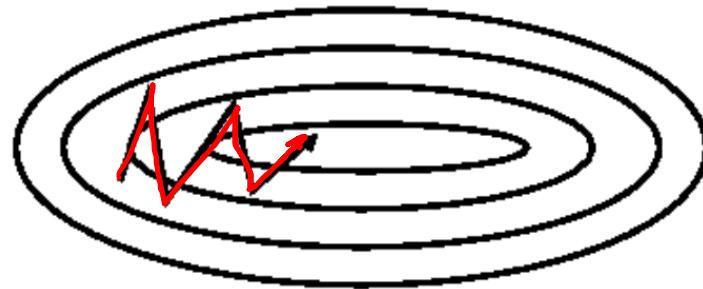
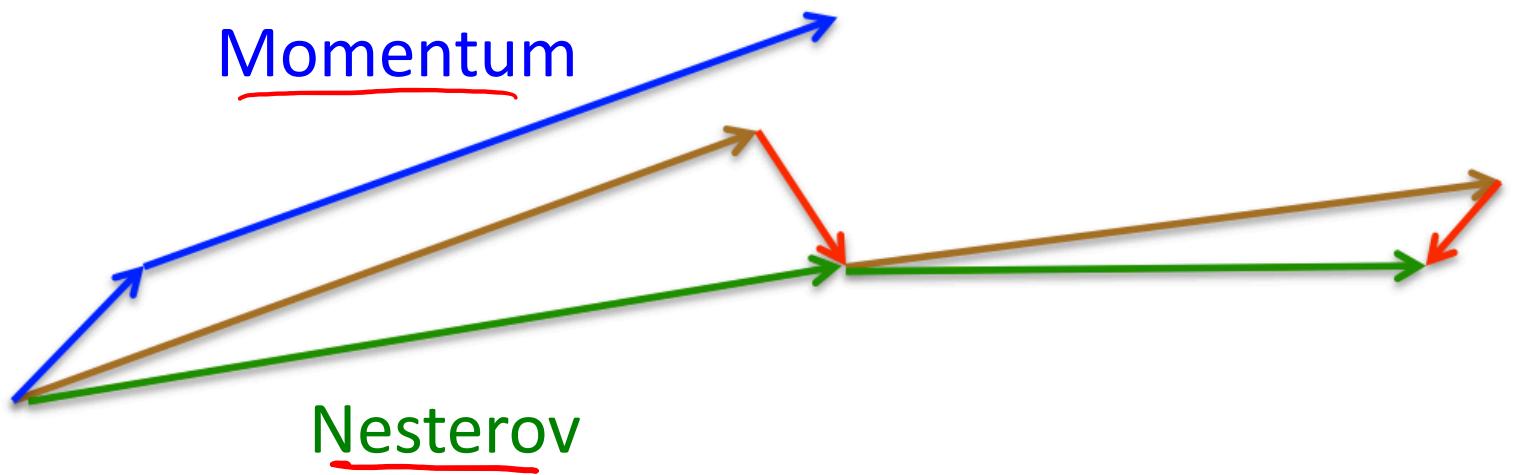


Image 3: SGD with momentum

$$\gamma \approx \underline{\underline{0.9}}$$

3. Nesterov Accelerated Gradient (NAG)

$$v_{k+1} = \underline{\gamma} v_k + \underline{\alpha} \cdot \nabla_w J(w_k - \gamma v_k)$$
$$w_{k+1} = w_k - v_{k+1}$$



4. Adagrad

$$w_{k+1,j} = w_{k,j} - \frac{\alpha}{\sqrt{G_{k,jj} + \epsilon}} \cdot \underline{\frac{g_{k,j}}{\text{gradient}}}$$

- Adagrad uses a different learning rate for every parameter w_j at every step k . G is diagonal matrix of sum squared gradient values.
- Performs **smaller updates** (i.e. low learning rates) for parameters associated with frequently occurring features, and **larger updates** (i.e. high learning rates) for parameters associated with infrequent features.

5. RMSProp

$$\underline{E[g^2]}_k = \gamma \underline{E[g^2]}_{k-1} + (1 - \gamma) g_k^2$$

$$w_{k+1} = w_k - \frac{\eta}{\sqrt{\underline{E[g^2]}_k + \epsilon}} g_k$$

- Prevents accumulation by adding regularizing term in the running average (exponentially decaying)
- Beneficial for RNNs

6. Adadelta

$$E[\Delta w^2]_k = \gamma E[\Delta w^2]_{k-1} + (1 - \gamma) \Delta w_k^2 .$$

$$\underline{RMS[\Delta w]}_k = \sqrt{E[\Delta w^2]_k + \epsilon} ;$$

$$\Delta w_k = -\frac{RMS[\Delta w]_{k-1}}{RMS[g]_k} g_k ;$$

$$w_{k+1} = w_k + \Delta w_k ;$$

- Generalizes RMSProp / Adagrad for considering RMS instead of accumulation of grad
- No learning rate parameter

7. AdaM- Adaptive Moment Estimation

$$m_k = \beta_1 m_{k-1} + (1 - \beta_1) g_k$$

$$v_k = \beta_2 v_{k-1} + (1 - \beta_2) g_k^2$$

$$w_{k+1} = w_k - \frac{\eta}{\sqrt{\hat{v}_k} + \epsilon} \hat{m}_k$$

- Keeps track of 2 moments: mean and variance
- Normalizes them to prevent biases

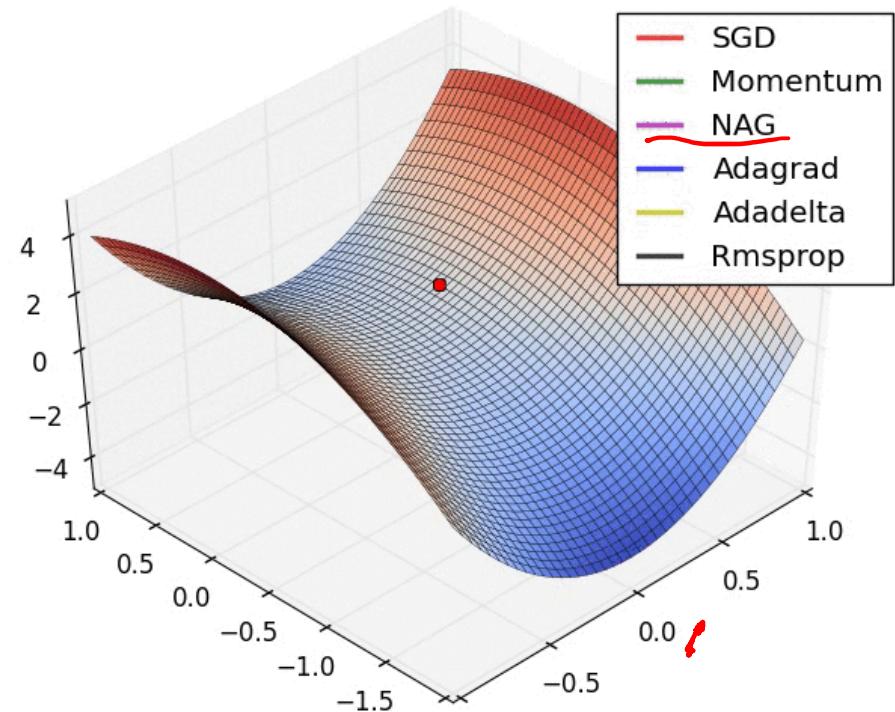
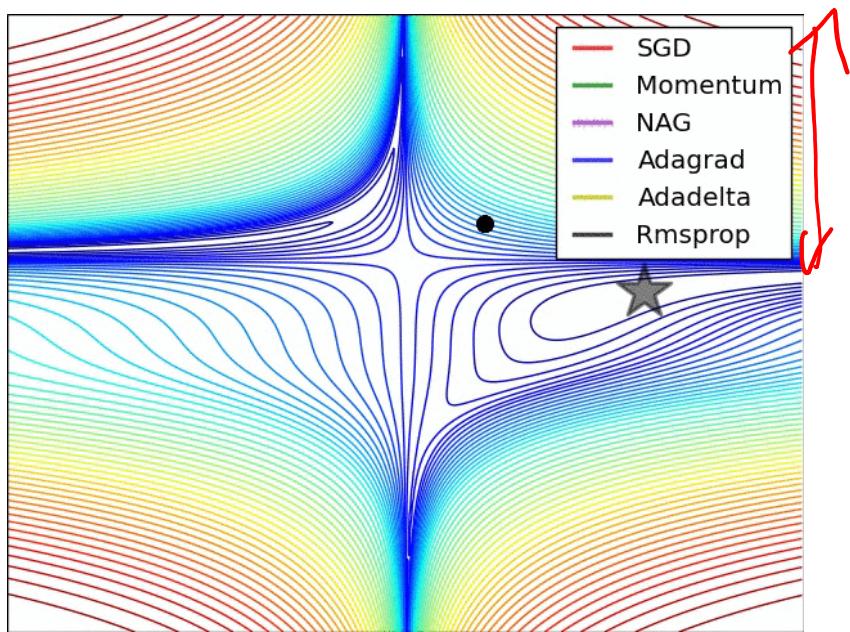
$$\hat{m}_k = \frac{m_k}{1 - \beta_1^k}$$

$$\hat{v}_k = \frac{v_k}{1 - \beta_2^k}$$

Additional

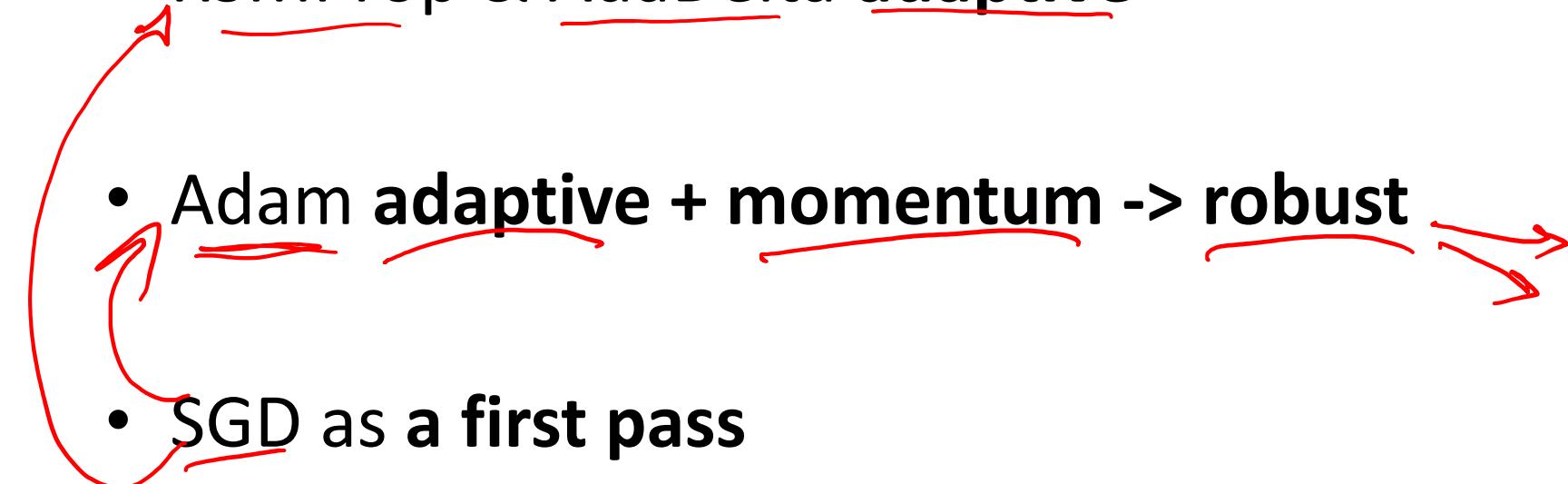
- 8 AdaMax
 - Generalization of AdaM to L-infinity norm
- 9 Nadam
 - Nesterov Adam
- 10 AMSgrad
 - Max normalization instead of exponential in Adam

Visual Comparison



Which optimizer to use?



- RMSProp & AdaDelta adaptive
 - Adam adaptive + momentum -> **robust**
 - SGD as a first pass
- 

The Right Optimization is Key



deeplearning.ai presents
Heroes of Deep Learning

Andrej Karpathy
Director of AI at Tesla



Lab This Week

The MNIST Dataset

0 0 0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3 3 3
4 4 4 4 4 4 4 4 4 4 4 4 4 4
5 5 5 5 5 5 5 5 5 5 5 5 5 5
6 6 6 6 6 6 6 6 6 6 6 6 6 6
7 7 7 7 7 7 7 7 7 7 7 7 7 7
8 8 8 8 8 8 8 8 8 8 8 8 8 8
9 9 9 9 9 9 9 9 9 9 9 9 9 9

- Adjust your setup to achieve a testing accuracy of 90% or better

Next Week - Deep Learning Practices

- Practical Optimization
- Operators
- Drop out
- Initialization
- Normalization
- Project Cycle with Deep Learning Methodology
- Introduction to CNN