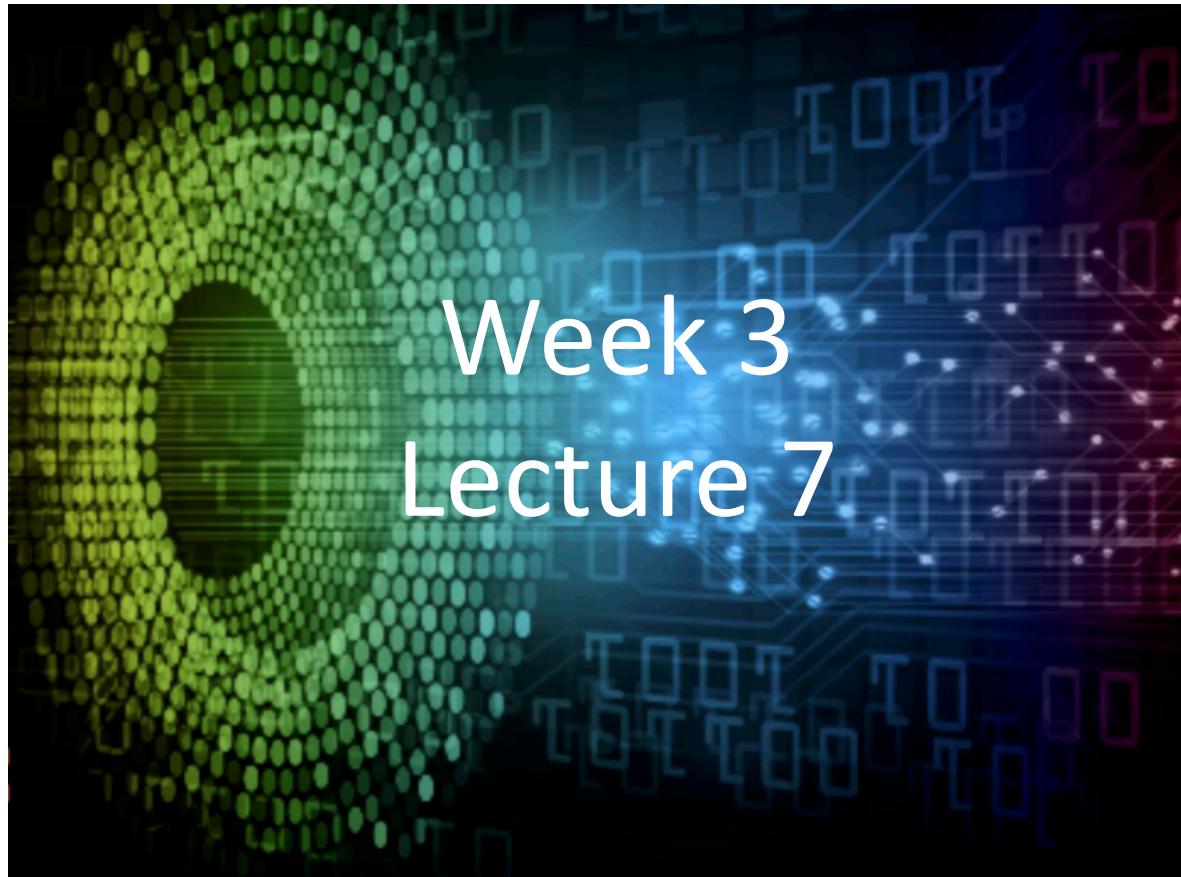


Introduction to Deep Learning Applications and Theory



ECE 596 / AMATH 563

Previous Lecture/Week: Network Training/Optimization

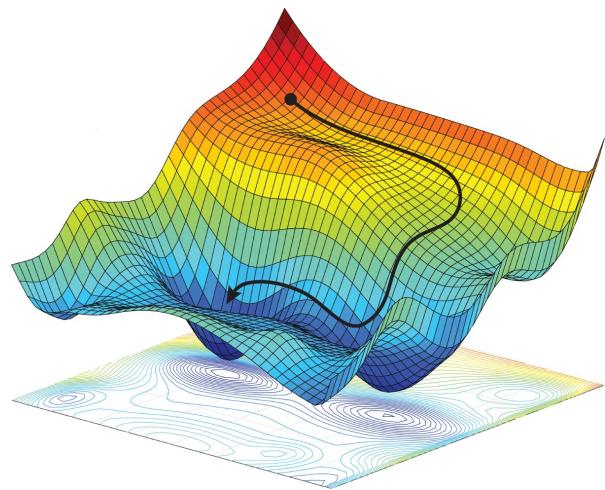
1. Model Optimization

1. Gradient Descent

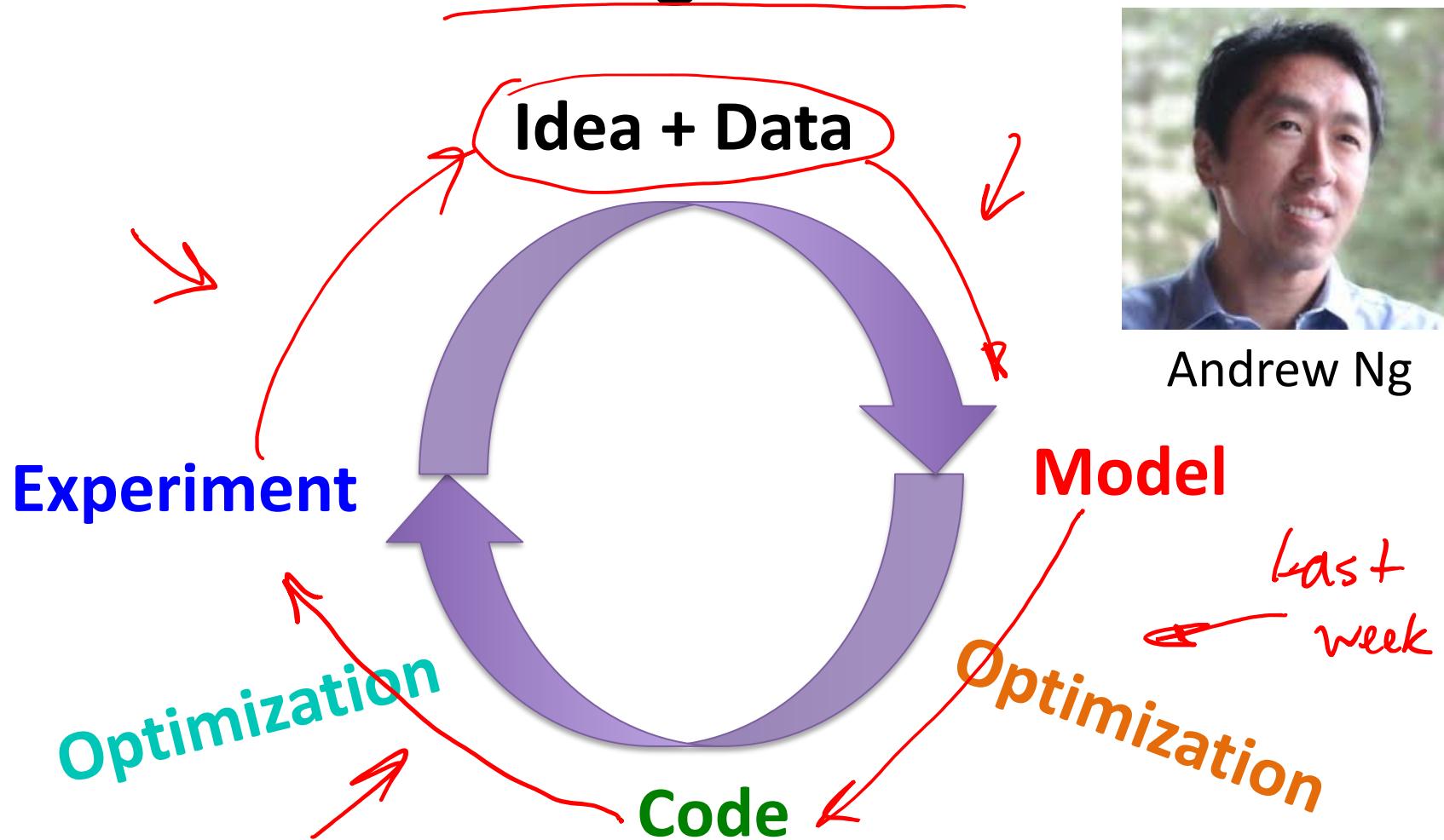
2. SGD

3. Extensions

2. Training NN



Current Week: The Data Driven Deep Learning Process



Andrew Ng

Current Week: The Data Driven Deep Learning Process

1. Hyperparameters and Optimization

1. Processing Stochastic Data
(EMA) *rederive*
2. Optimization with Stochastic
Data

Parameters

- **Parameters**

Model Parameters:

W, b, activation, output, cost

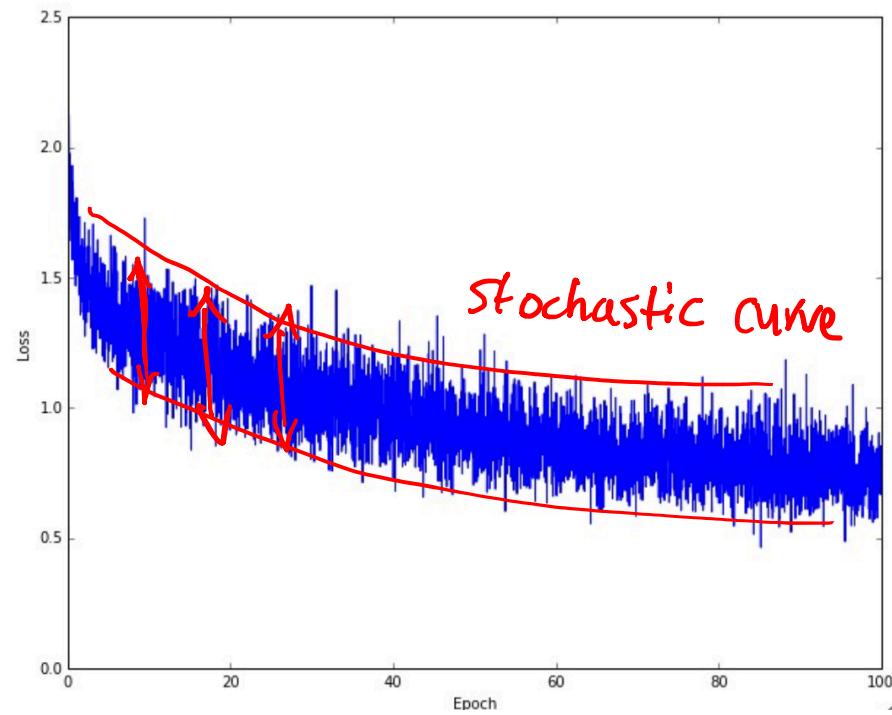
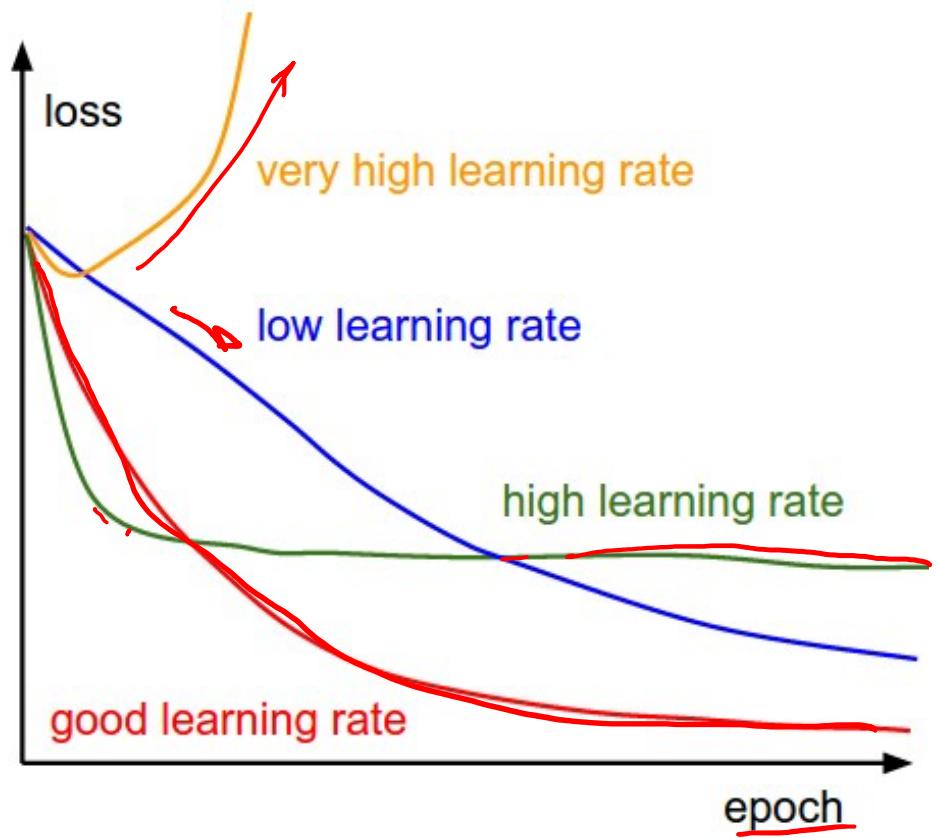
- **Hyper-parameters**

External Parameters:

Batch/minibatch size

Learning parameters

Experimenting with GD



Hyperparameters = α

Stochastic Data – TSLA Stock

TSLA \$276.0600 \$3.7500 ↑ 1.38%

Volume (M): 4,168 Open : 276.74 52 Wk Avg : 310.02



Exponentially Moving Average

(EMA)

$$v_k = \gamma v_{k-1} + (1 - \gamma) \theta_k$$

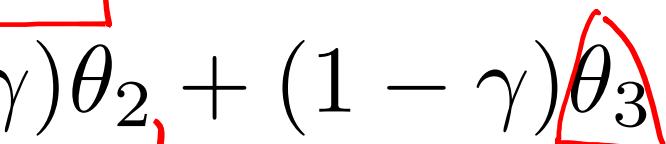
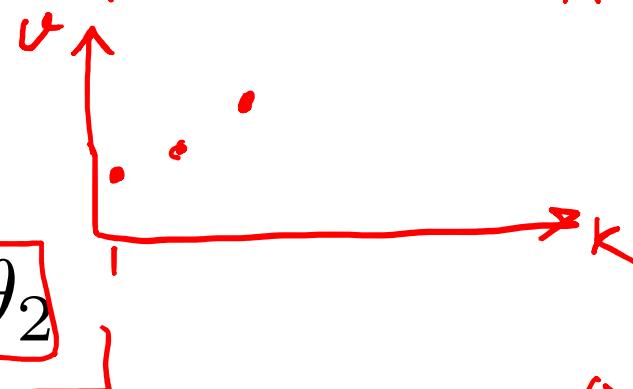
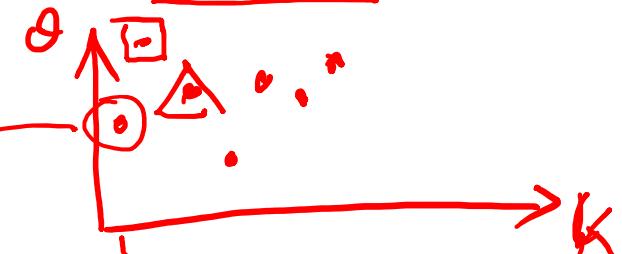
$$v_0 = 0$$

Let's iterate it:

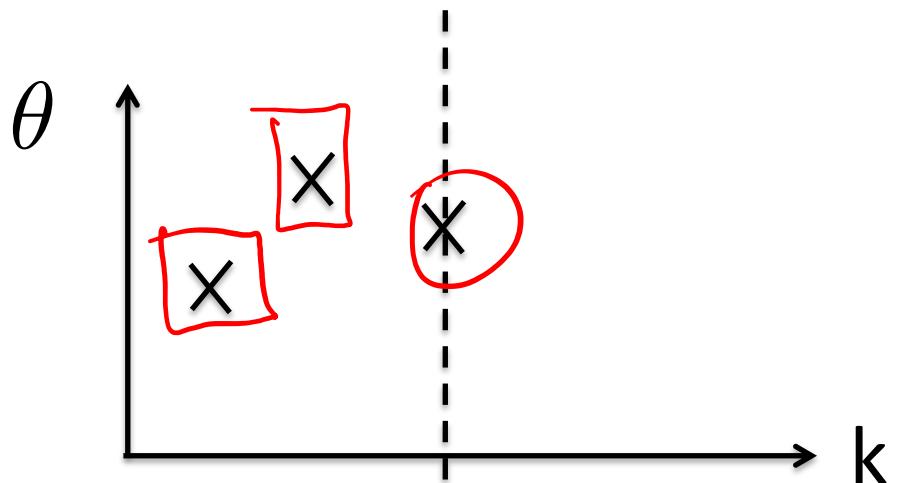
$$v_1 = (1 - \gamma) \theta_1$$

$$v_2 = \gamma(1 - \gamma) \theta_1 + (1 - \gamma) \theta_2$$

$$v_3 = \gamma^2(1 - \gamma) \theta_1 + \gamma(1 - \gamma) \theta_2 + (1 - \gamma) \theta_3$$



Exponentially Moving Average



Window width ~

$$\frac{1}{1 - \gamma}$$

$\gamma = 0.9 \approx 10$ steps window

$\gamma = 0.98 \approx 50$ steps window

$$\times (1 - \gamma)$$

Exponentially Moving Average

TSLA \$276.0600 \$3.7500 ↑ 1.38%

Volume (M): 4,168 Open : 276.74 52 Wk Avg : 310.02

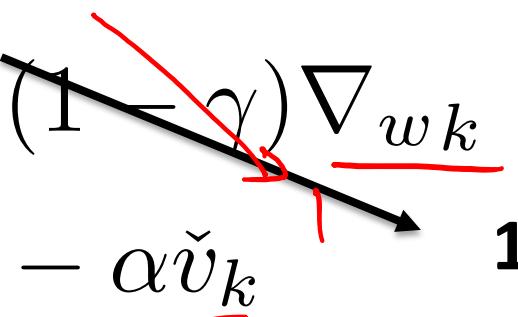
Close: \$276.06 Performance: +89.82% Open: 149.50 Low: 116.10 High: 389.61 Period Selected: Aug 13, 2013 - Apr 10, 2019

chart by amCharts.com



Momentum Method

- Let's apply it to the **gradient** values

$$\begin{aligned}\check{v}_k &= \gamma \check{v}_{k-1} + (1 - \gamma) \nabla_{w_k} \\ w_{k+1} &= w_k - \alpha \check{v}_k\end{aligned}$$


Smoothens the gradient

Hyperparameters = $\alpha, \frac{\gamma}{\beta}$

RMS Prop

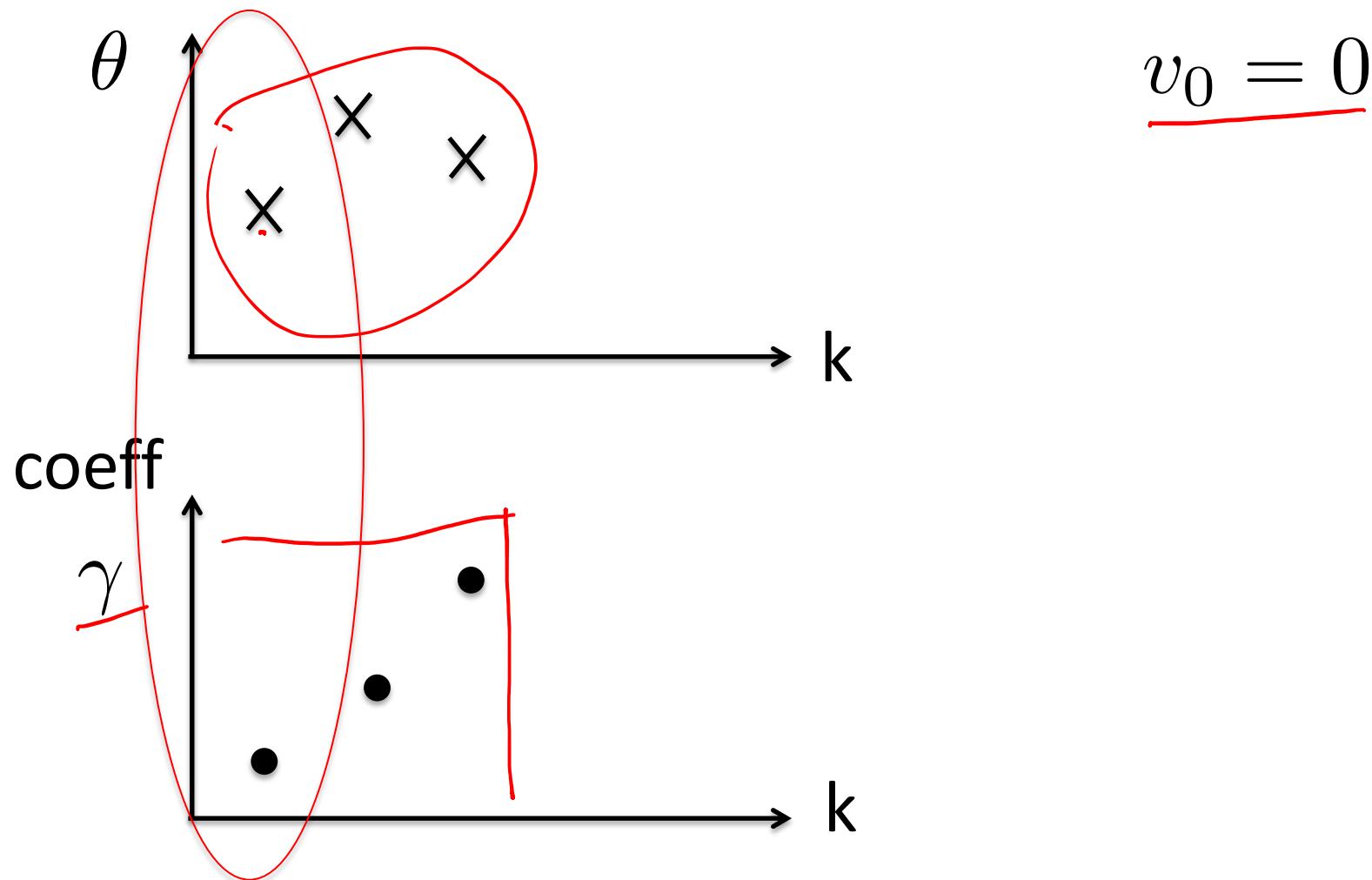
- Let's apply EMA to **square gradient values**

$$\underline{S}_k = \gamma S_{k-1} + (1 - \gamma) (\nabla_{w_k})^2$$

$$w_{k+1} = w_k - \frac{\alpha}{\sqrt{S_k} + \epsilon} \nabla_{w_k}$$

Hyperparameters = $\alpha, \gamma; \epsilon$

Biases in EMA



Exponentially Moving Average

TSLA \$276.0600 \$3.7500 ↑ 1.38%

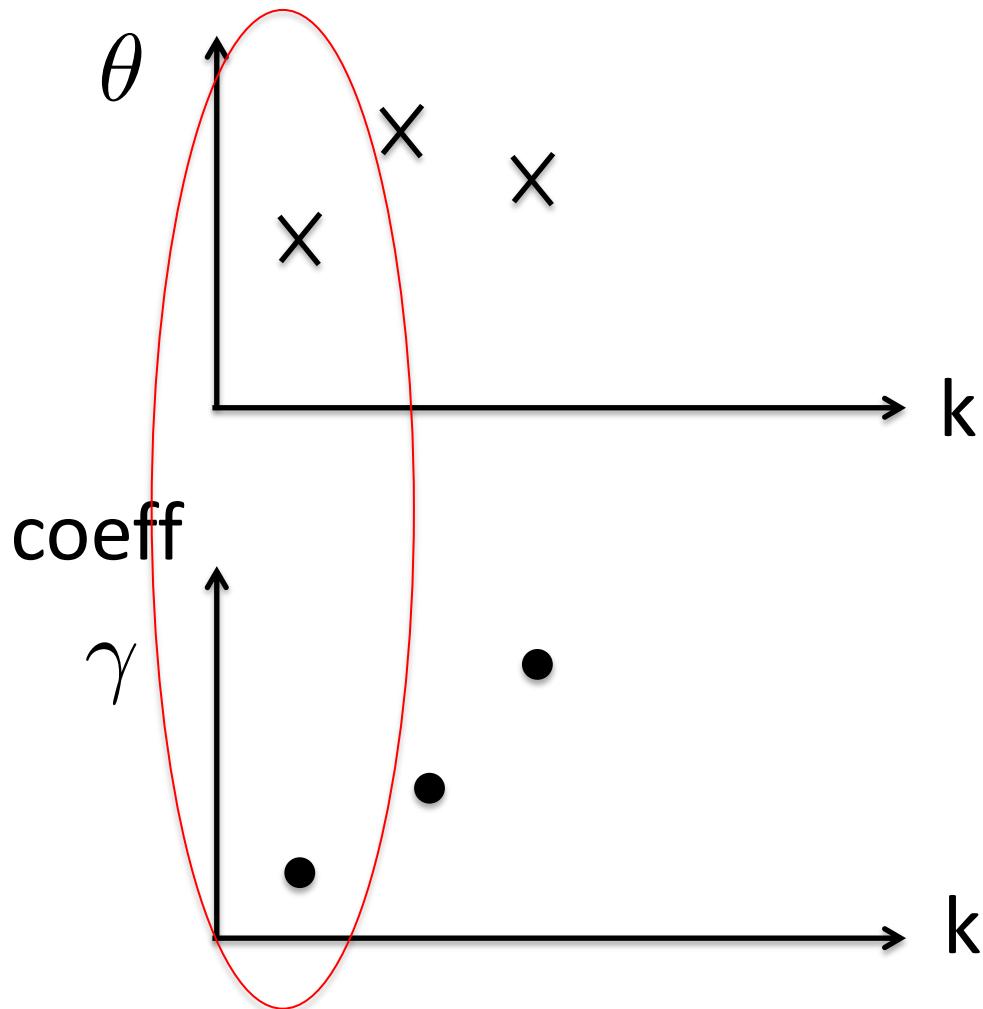
Volume (M): 4,168 Open : 276.74 52 Wk Avg : 310.02

Close: \$276.06 Performance: +89.82% Open: 149.50 Low: 116.10 High: 389.61 Period Selected: Aug 13, 2013 - Apr 10, 2019

chart by amCharts.com



Biases in EMA Correction



$$v_0 = 0$$

Correction:

$$\tilde{v}_k = \frac{v_k}{1 - \gamma^k}$$

small v_k γ small
1 v_k large γ large

Adam

- What if we combine Momentum with RMSProp

$$\underline{\check{v}_k} = \gamma_1 \check{v}_{k-1} + (1 - \gamma_1) \underline{\nabla_{w_k}}$$

$$\underline{S_k} = \gamma_2 S_{k-1} + (1 - \gamma_2) \underline{(\nabla_{w_k})^2}$$

- And do bias correction

$$\check{v}_k^c = \frac{\check{v}_k}{1 - \gamma_1^k}$$

$$S_k^c = \frac{S_k}{1 - \gamma_2^k}$$

$$w_{k+1} = w_k - \frac{\alpha}{\sqrt{S_k^c} + \epsilon} \check{v}_k^c$$

Hyperparameters = $\alpha, \gamma_1, \gamma_2; \epsilon$



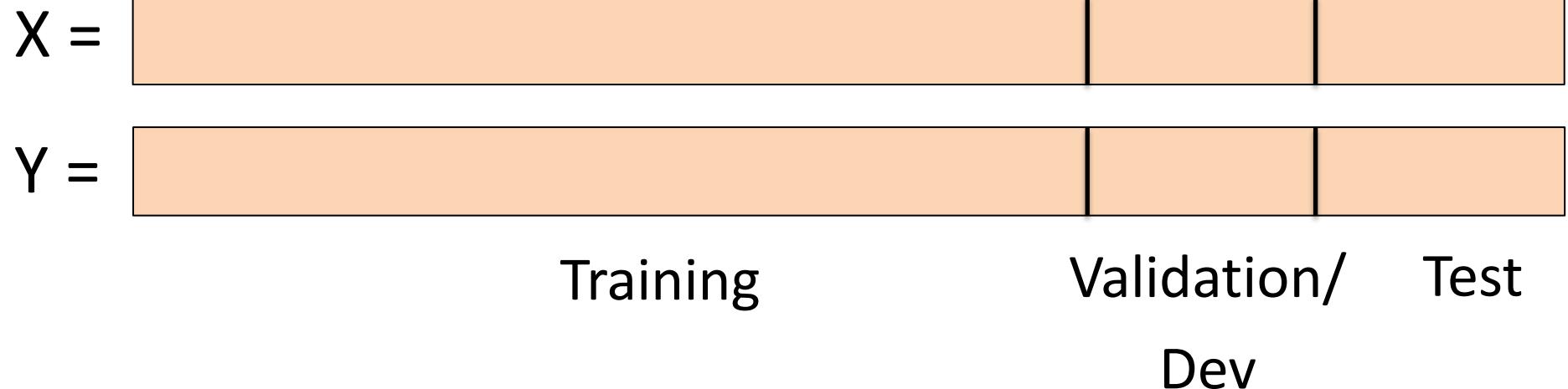
2019

2019

2021

Supervised Learning

Dataset D =



- Plausible divisions
- Distributions of each set

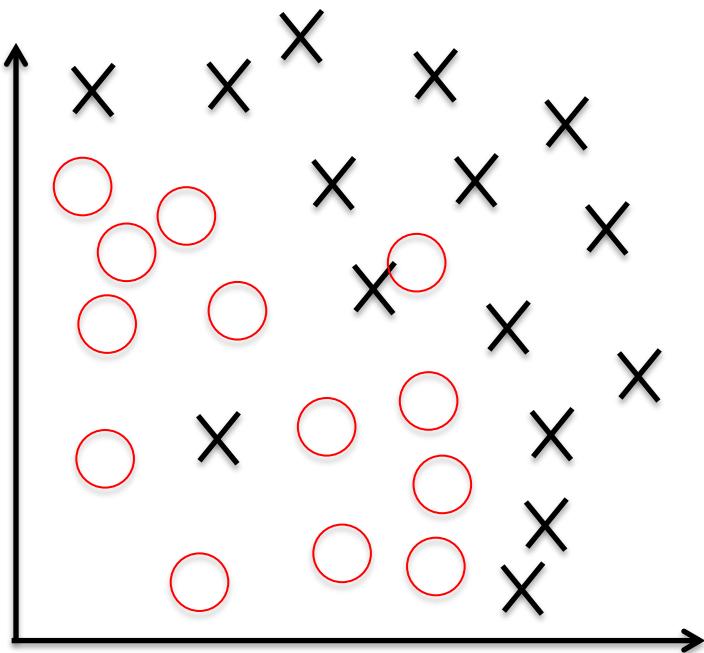
Over-fitting

- **Discrepancy** between
Training performance and
Testing performance

High **variance** – over-fitting

High **bias** – under-fitting

High Variance vs High Bias



Gradient Descent

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

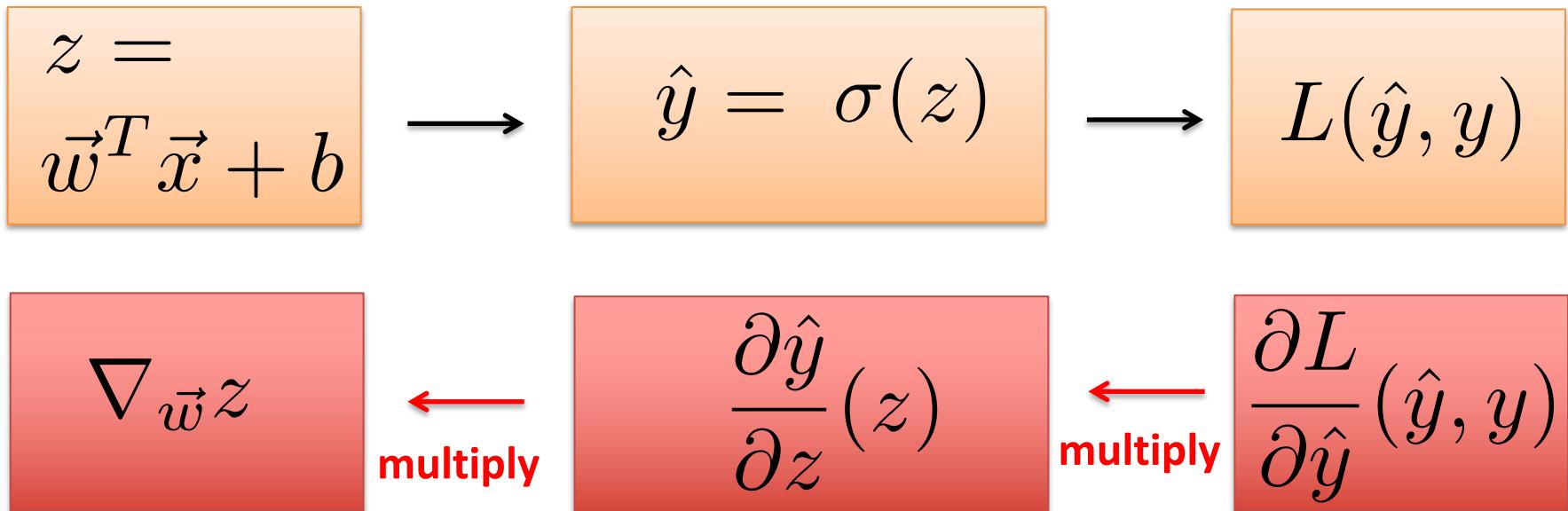
$$\vec{b}_{k+1} = \vec{b}_k - \alpha \frac{\partial}{\partial b} J(\vec{w}; b_k)$$

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

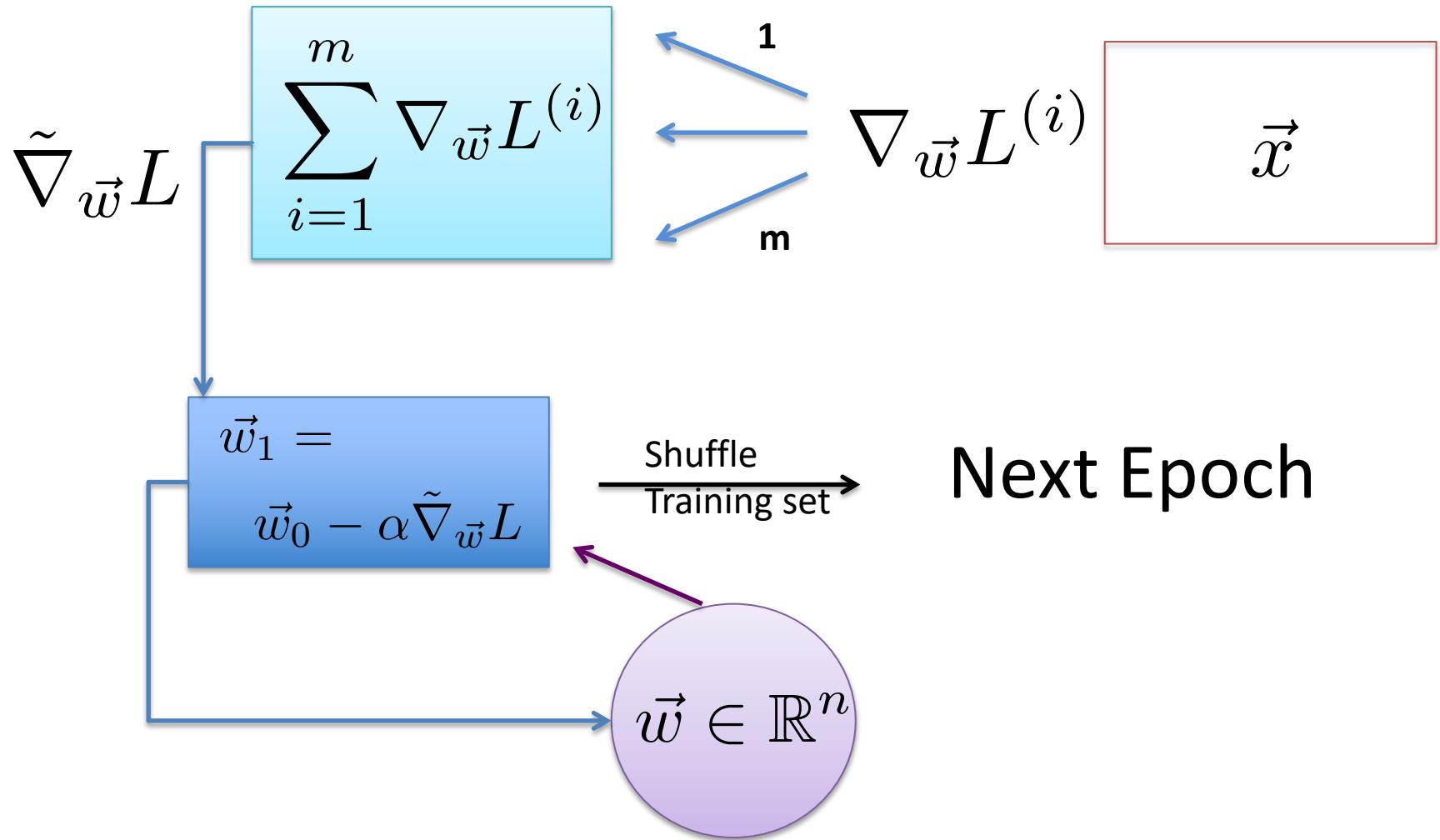
Our cost: $J = L(\hat{y}, y)$

Back Propagation – Chain Rule

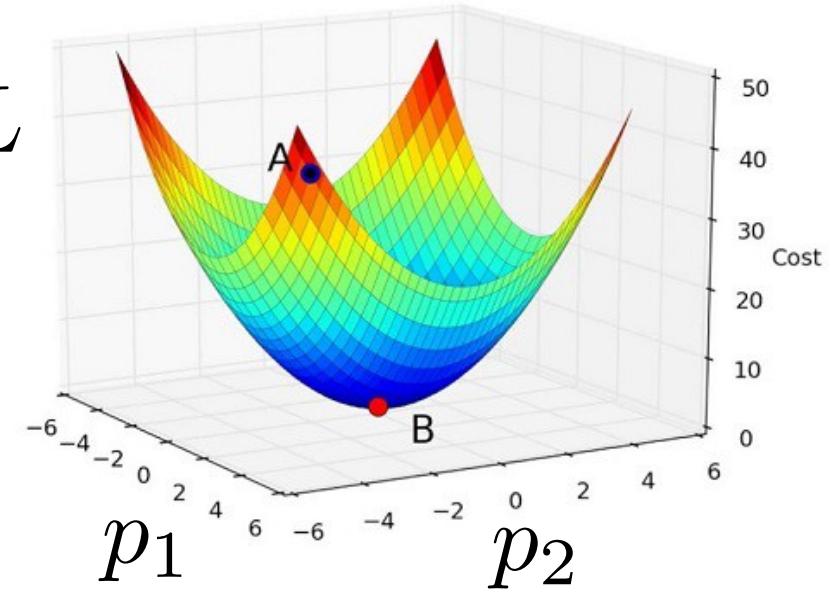
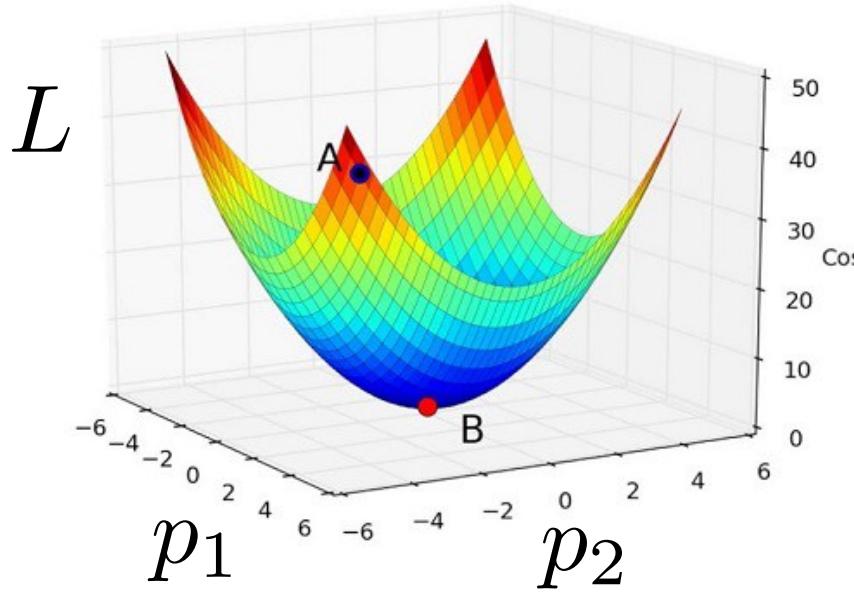
$$\nabla_{\vec{w}} L(\hat{y}, y) = \frac{\partial L}{\partial \hat{y}} \frac{\partial \hat{y}}{\partial z} \nabla_{\vec{w}} z$$



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

Stochastic (Online) GD

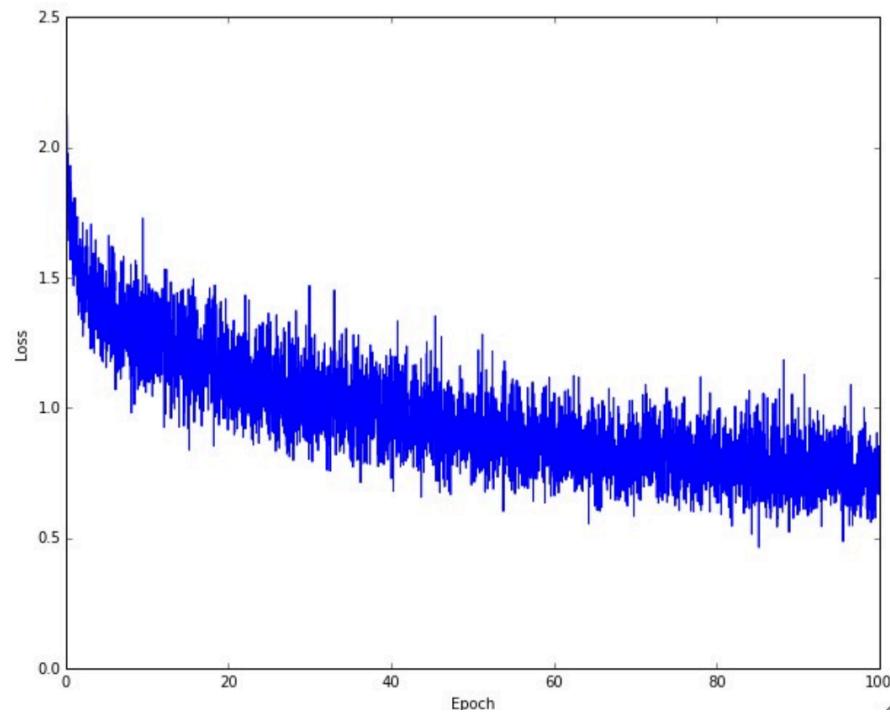
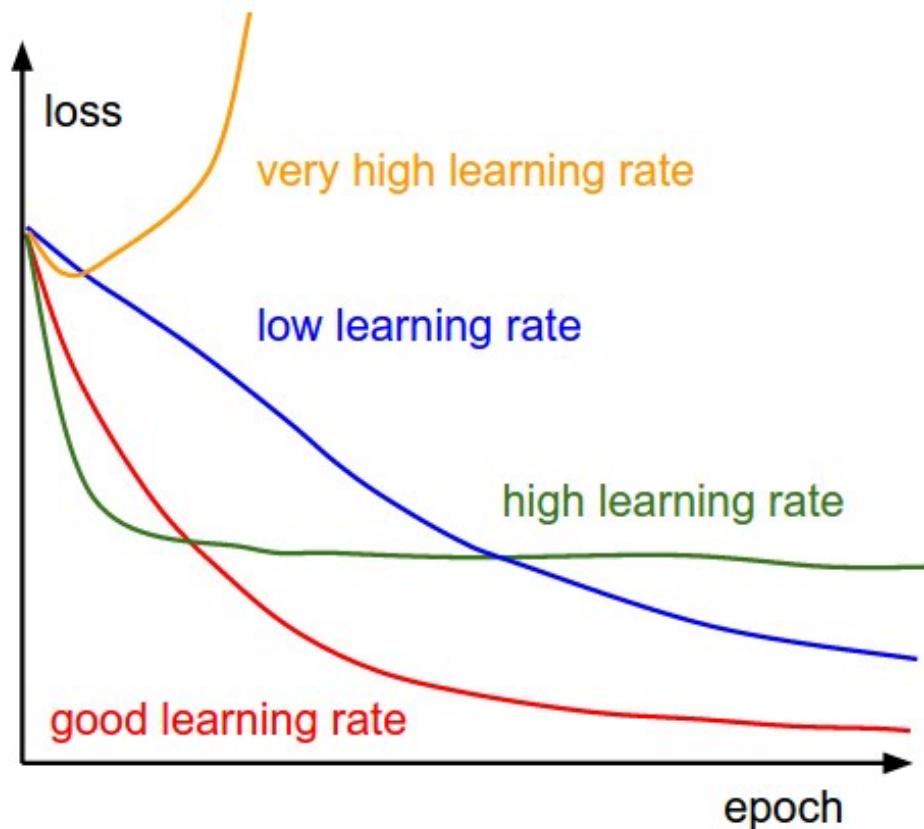
$$w_{k+1} = w_k - \alpha \cdot \nabla_w J(w; x^{(i)}; y^{(i)})$$

- 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

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

- 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

- 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

Learning rate decay

$$\alpha = \frac{1}{1 + \textit{decr} \cdot \textit{epnum}} \alpha_0$$

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

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

Momentum Method

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

$$w_{k+1} = w_k - v_{k+1}$$



Image 2: SGD without momentum

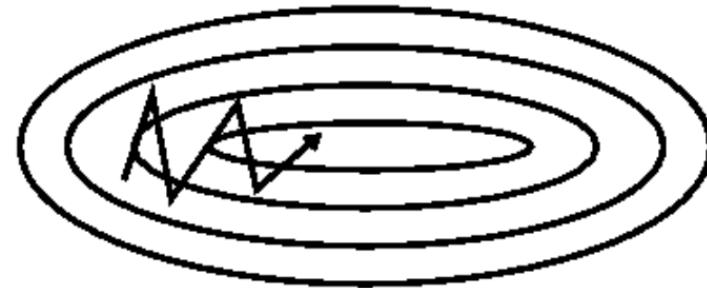


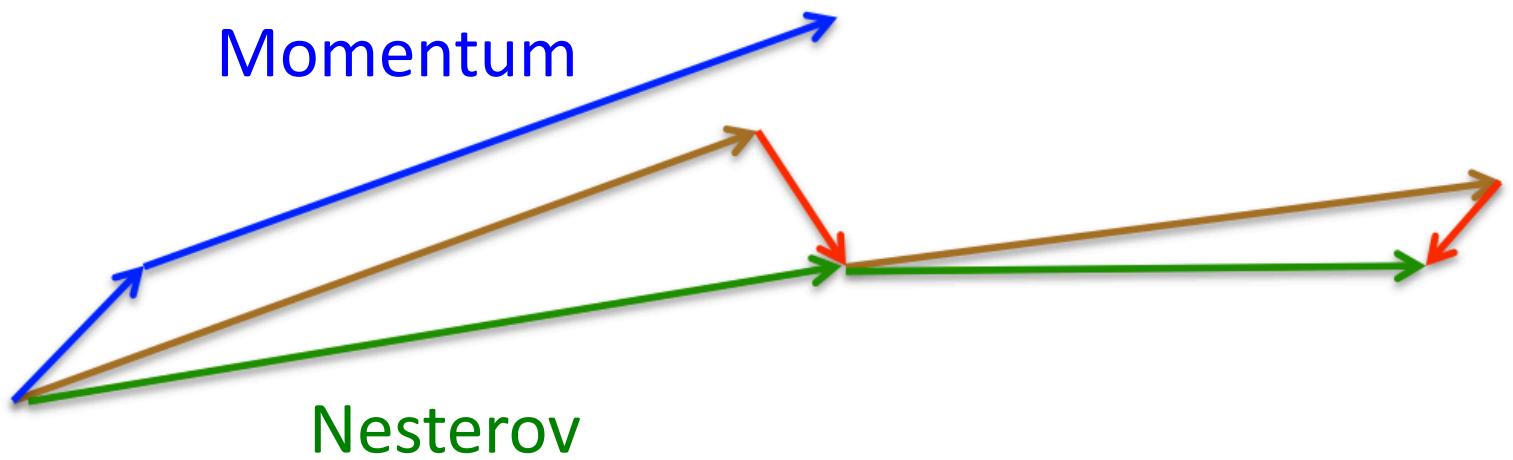
Image 3: SGD with momentum

$$\gamma \approx 0.9$$

Nesterov Accelerated Gradient

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

$$w_{k+1} = w_k - v_{k+1}$$



Adagrad

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

- 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**.

RMSProp

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

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

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

Adadelta

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

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

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

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

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

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

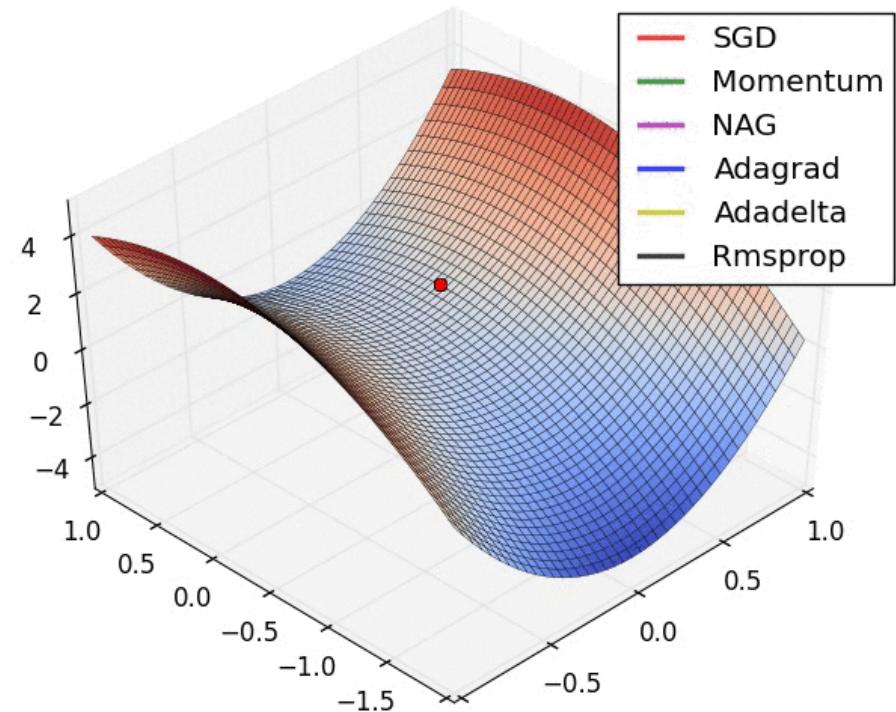
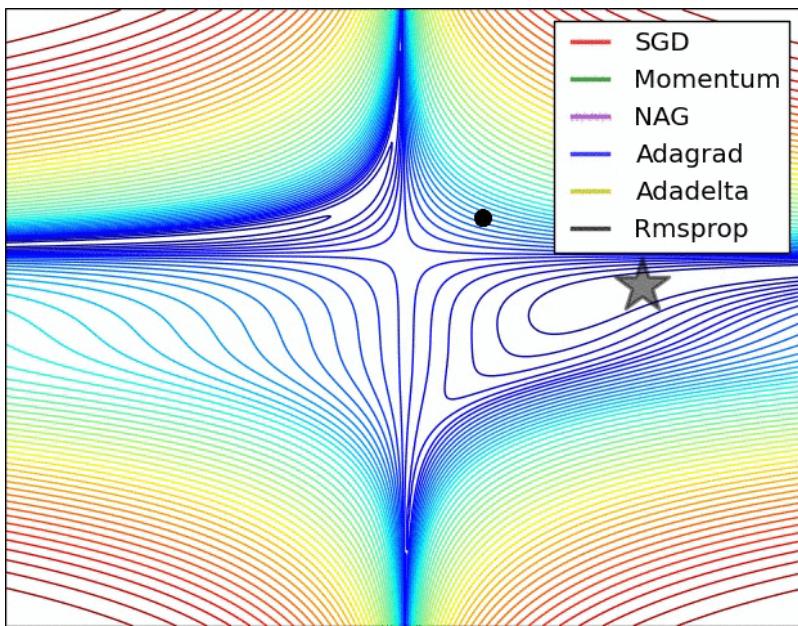
$$\hat{m}_k = \frac{m_k}{1 - \beta_1^k}$$

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

Additional

- AdaMax
 - Generalization of AdaM to L-infinity norm
- Nadam
 - Nesterov Adam
- 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

