

UVA CS 4774: Machine Learning

□

Lecture 4: More optimization for Linear Regression

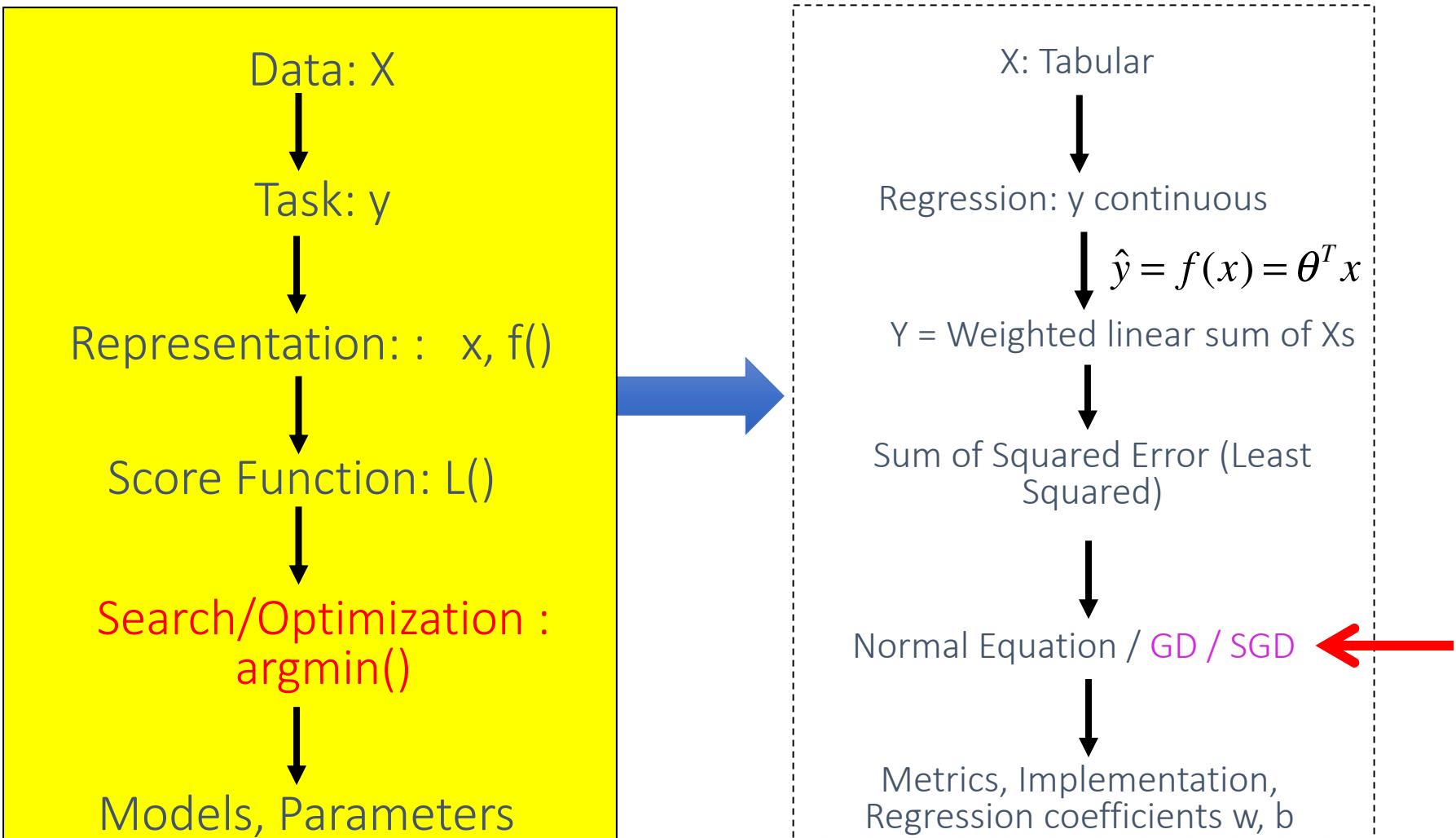
Dr. Yanjun Qi

University of Virginia
Department of Computer Science

Rough Sectioning of this Course

- 
- 1. Basic Supervised Regression + on Tabular Data
 - 2. Basic Deep Learning + on 2D Imaging Data
 - 3. Advanced Supervised learning + on Tabular Data
 - 4. Generative and Deep + on 1D Sequence Text Data
 - 5. Not Supervised + Mostly on Tabular Data

Today : GD and SGD for Multivariate Linear Regression

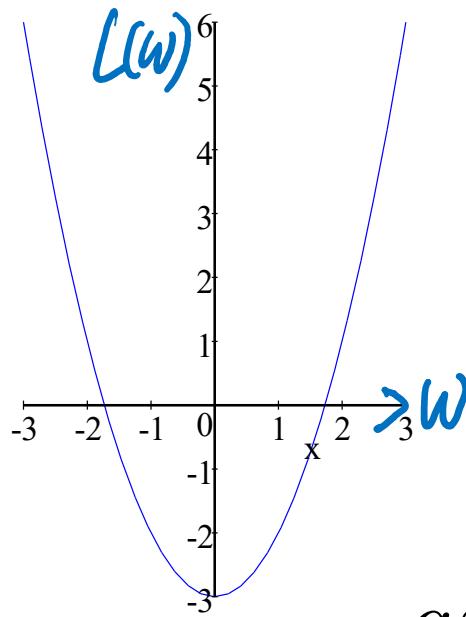


A little bit more about [Optimization]

- Objective function $F(x) \rightarrow J(\theta)$
- Variables $x \rightarrow \theta$
- Constraints $\rightarrow \theta \in \mathbb{R}^P$

To find values of the variables
that minimize or maximize the objective function
while satisfying the constraints

Method 1: Directly Optimize



Minimizing a Quadratic Function

$$L(w) = w^2 - 3$$

$$\underline{L'(w) = 2w = 0}$$

$$\operatorname{argmin}_w L(w)$$

$$= \underline{(L'(w)=0)} \quad || \quad w^*$$

This quadrative (convex) function is minimized @ the unique point whose derivative (slope) is zero.

→ When we find zeros of the derivative of this function, we also find the minima (or maxima) of that function.

Method I: normal equations to solve for LR training

- Write the cost function in matrix form:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

$$= \frac{1}{2} (\mathbf{X}\theta - \bar{\mathbf{y}})^T (\mathbf{X}\theta - \bar{\mathbf{y}})$$

$$= \frac{1}{2} (\theta^T \mathbf{X}^T \mathbf{X} \theta - \theta^T \mathbf{X}^T \bar{\mathbf{y}} - \bar{\mathbf{y}}^T \mathbf{X} \theta + \bar{\mathbf{y}}^T \bar{\mathbf{y}})$$

$$\mathbf{X}_{train} = \begin{bmatrix} \cdots & \mathbf{x}_1^T & \cdots \\ \cdots & \mathbf{x}_2^T & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \mathbf{x}_n^T & \cdots \end{bmatrix} \quad \bar{\mathbf{y}}_{train} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

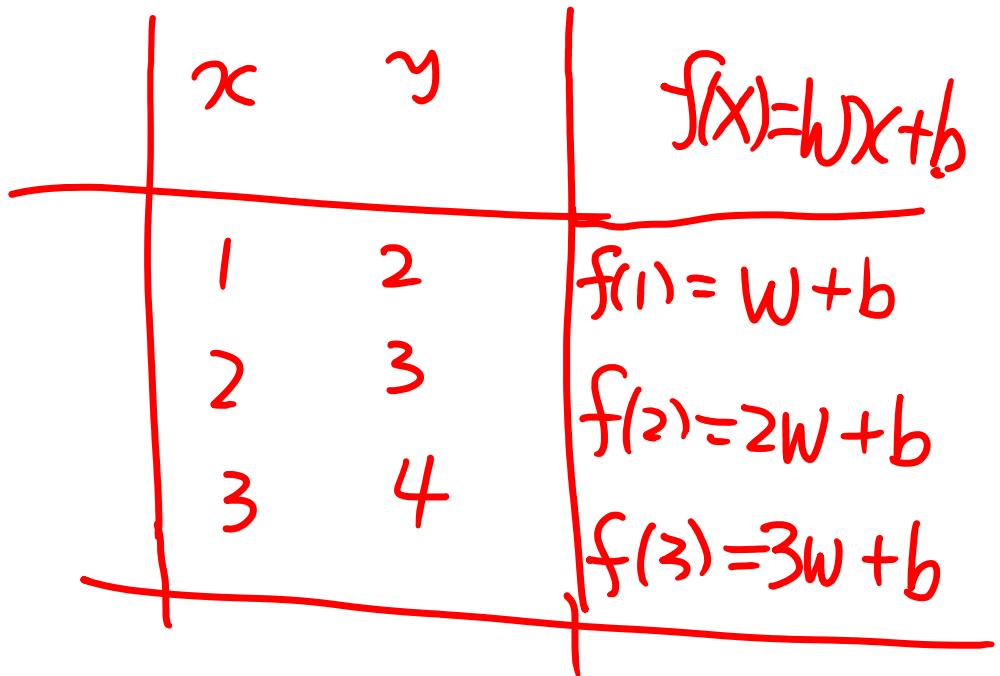
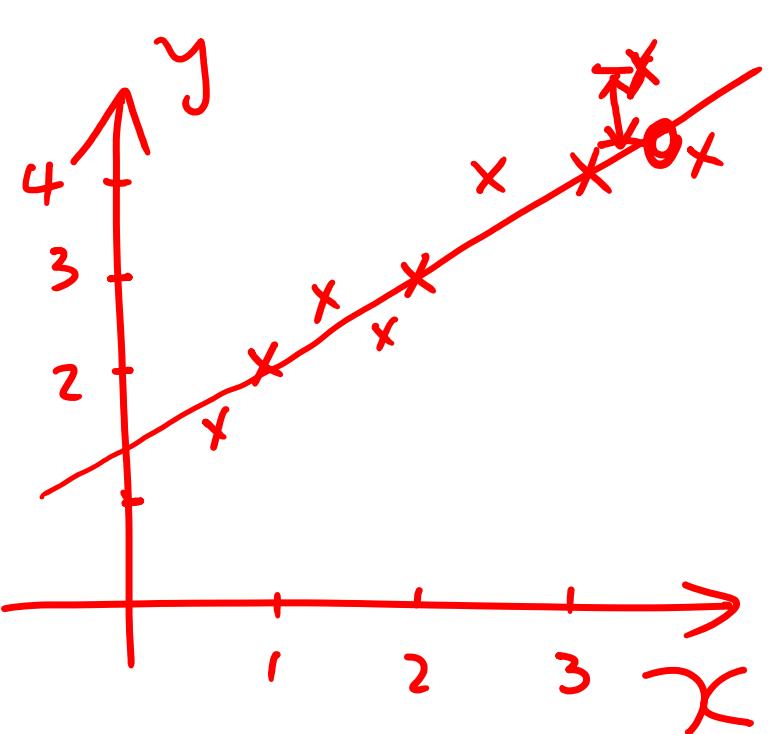
To minimize $J(\theta)$, take its gradient and set to zero:

$$\nabla_{\theta} J(\theta) = 0 \Rightarrow X^T X \theta = X^T \bar{y}$$

The normal equations

$$\theta^* = \underbrace{(X^T X)^{-1} X^T \bar{y}}$$

$X^T X$ Gram matrix
↓
PSD
↓
 $J(\theta)$ convex
↓
 $\nabla J(\theta) = 0 \Rightarrow \theta^*$



$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2 =$$

loss(w, b) =

$\frac{1}{2}$
 $J(w, b)$
 $J(\theta_1, \theta_0)$

$$\left. \begin{aligned}
 & (w+b-2)^2 + \\
 & (2w+b-3)^2 + \\
 & (3w+b-4)^2
 \end{aligned} \right\} \Rightarrow = \arg \min_{w, b} \text{loss}(\cdot)$$

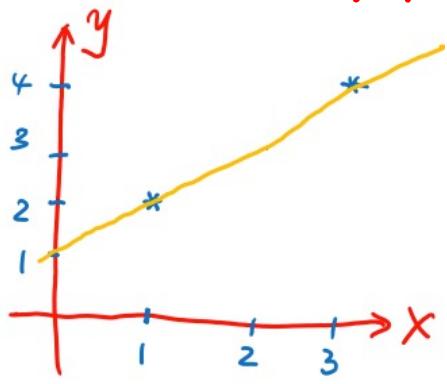
One concrete example

$$\Rightarrow \underline{J(w, b)} = \frac{1}{2}((w+b-2)^2 + (2w+b-3)^2)$$

$$\Rightarrow \begin{bmatrix} \frac{\partial J(w, b)}{\partial w} \\ \frac{\partial J(w, b)}{\partial b} \end{bmatrix} = \begin{bmatrix} (w+b-2) + (2w+b-3) \cdot 2 \\ w+b-2 + (2w+b-3) \end{bmatrix} = 0$$

2x1 Vector

$$\begin{cases} 5w + 3b - 8 = 0 \\ 3w + 2b - 5 = 0 \end{cases} \Rightarrow \begin{cases} 10w + 6b - 16 = 0 \\ 9w + 6b - 15 = 0 \end{cases}$$



$$\Rightarrow w = 1, b = \frac{1}{2}(5 - 3w) = 1$$

we solve the matrix equation via Gaussian Elimination

Method 2: Iteratively Optimize via Gradient Descent

Gradient Descent (GD): An iterative Algorithm

$$F(x)$$

- Initialize $k=0$, (randomly or by prior) choose x_0
- While $k < k_{\max}$ For the k -th epoch

$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$

Gradient Descent (GD): An iterative Algorithm

$$F(x)$$

- Initialize $k=0$, (randomly or by prior) choose x_0
- While $k < k_{\max}$

For the k -th epoch

$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$

α : learning rate / defined by users

$x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow x_3 \rightarrow \dots \rightarrow x_T$
epoch

Review: Definitions of gradient (more in Algebra-note)

- Size of gradient vector is always the same as the size of the variable vector

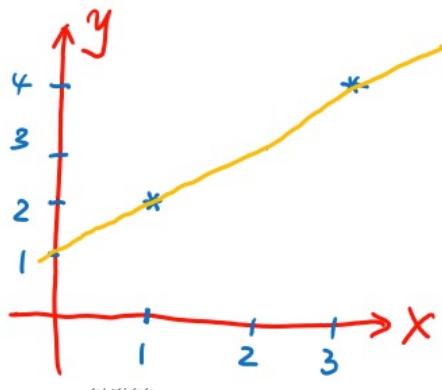
$$\nabla_{\underline{\underline{x}}} F(\underline{\underline{x}}) = \begin{bmatrix} \frac{\partial F(\underline{\underline{x}})}{\partial x_1} \\ \frac{\partial F(\underline{\underline{x}})}{\partial x_2} \\ \vdots \\ \frac{\partial F(\underline{\underline{x}})}{\partial x_p} \end{bmatrix} \in \mathbb{R}^p \quad \text{if } \underline{\underline{x}} \in \mathbb{R}^p$$

A vector whose entries, respectively, contain the p partial derivatives

Our concrete example

$$\Rightarrow J(w, b) = \frac{1}{2}((w+b-2)^2 + (2w+b-3)^2)$$

$$\Rightarrow \begin{bmatrix} \frac{\partial J(w, b)}{\partial w} \\ \frac{\partial J(w, b)}{\partial b} \end{bmatrix} = \begin{array}{l} (w+b-2) + (2w+b-3) \cdot 2 = 0 \\ w+b-2 + (2w+b-3) \cdot 2 = 0 \end{array}$$

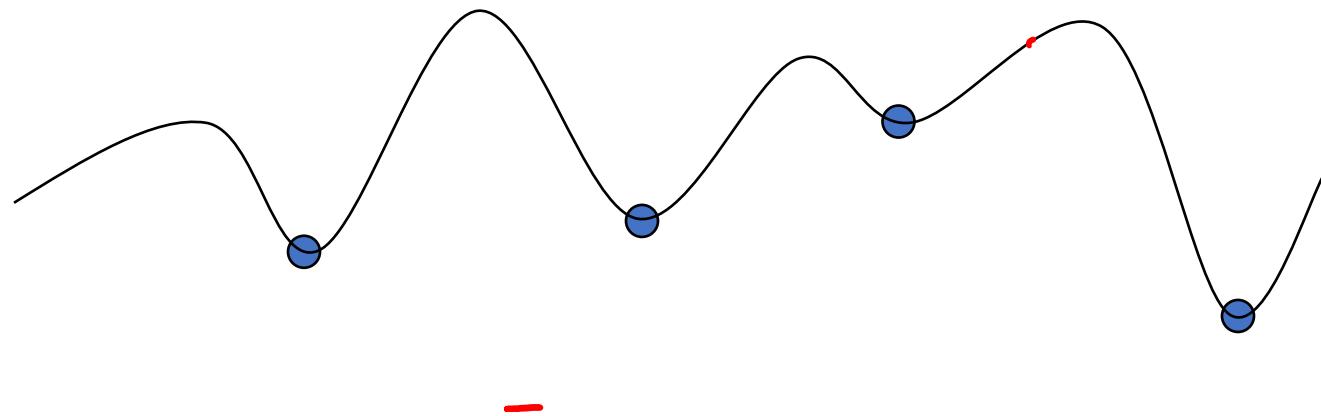


$$\nabla_{\vec{\theta}} J(\vec{\theta}) = \begin{bmatrix} \frac{\partial J(\vec{\theta})}{\partial \theta_1} \\ \frac{\partial J(\vec{\theta})}{\partial \theta_0} \end{bmatrix}$$

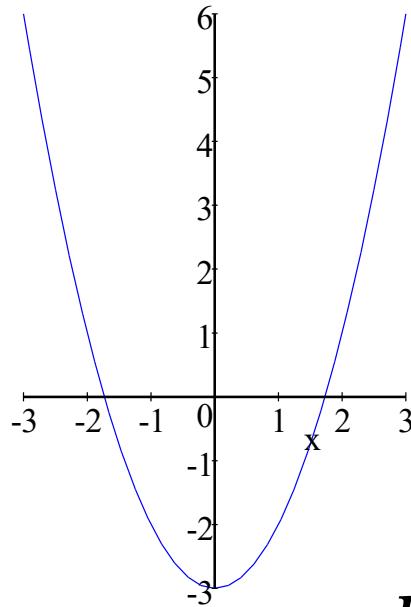
$$\vec{\theta} = \begin{bmatrix} w \\ b \end{bmatrix} = \begin{bmatrix} \theta_1 \\ \theta_2 \end{bmatrix}$$

WHY ? Optimize through Gradient Descent (iterative) Algorithms

- Works on any objective function
 - as long as we can evaluate the gradient
 - this can be very useful for minimizing complex functions



Review: Definitions of derivative (single variable case)



Review: Derivative of a Quadratic Function

$$l(w) = w^2 - 3$$

$$l'(w) = \lim_{h \rightarrow 0} \frac{(w+h)^2 - 3 - (w^2 - 3)}{h} = 2w$$

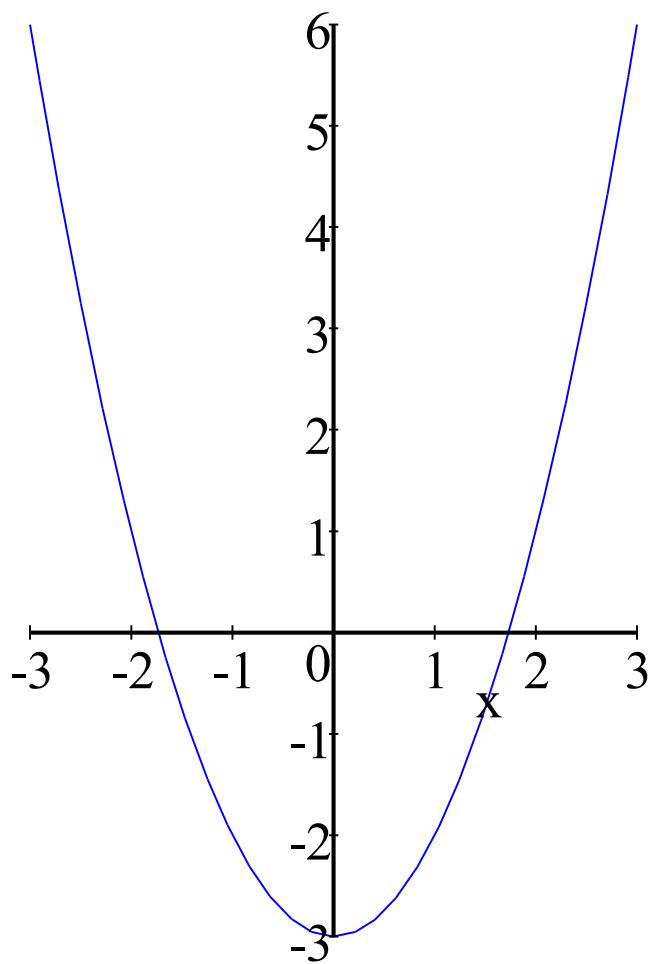
The derivative is often described as the "instantaneous rate of change",
→ the ratio of the instantaneous change in $F(x)$ to change in x

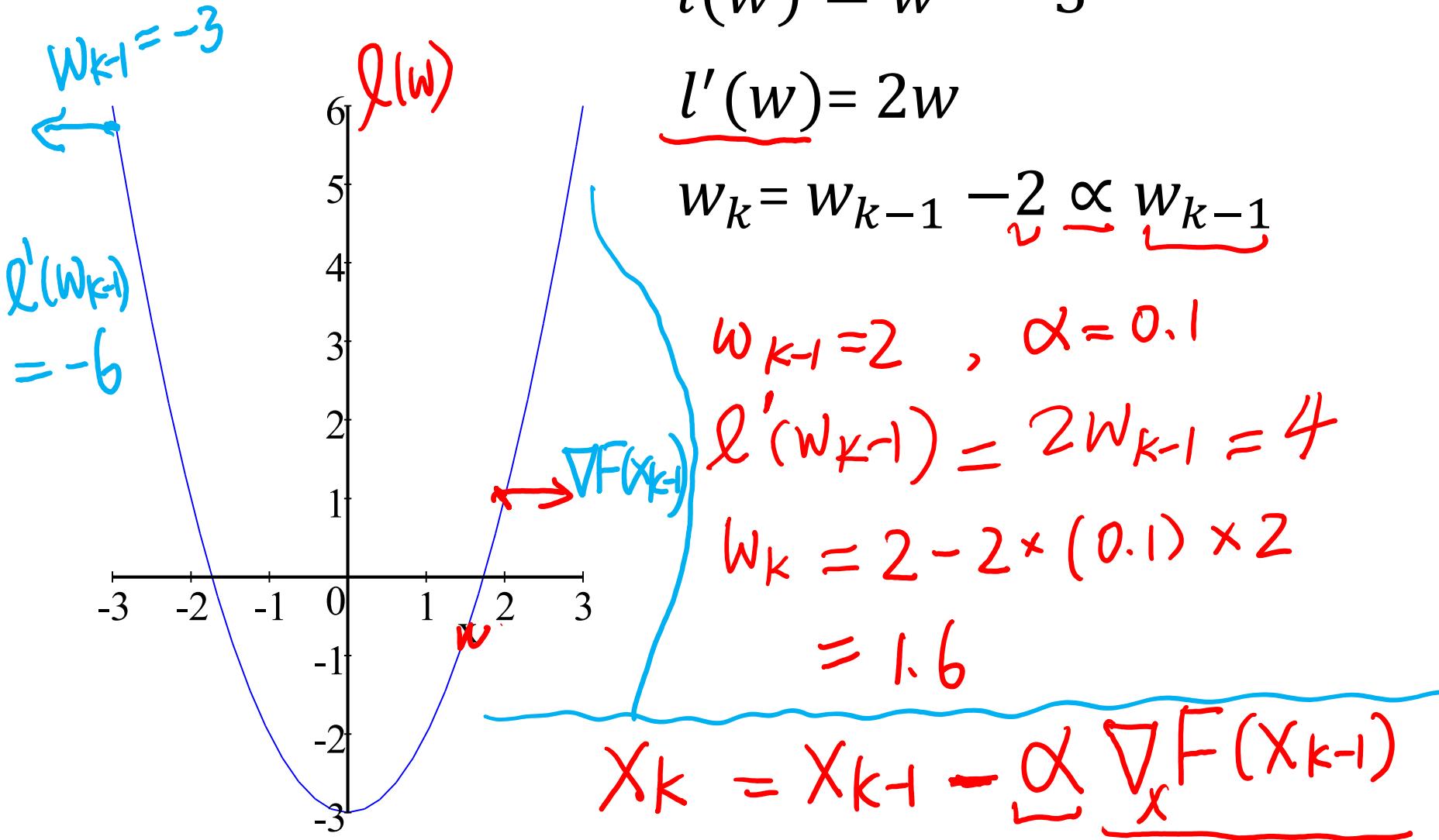
Our concrete example with Gradient Descent

$$l(w) = w^2 - 3$$

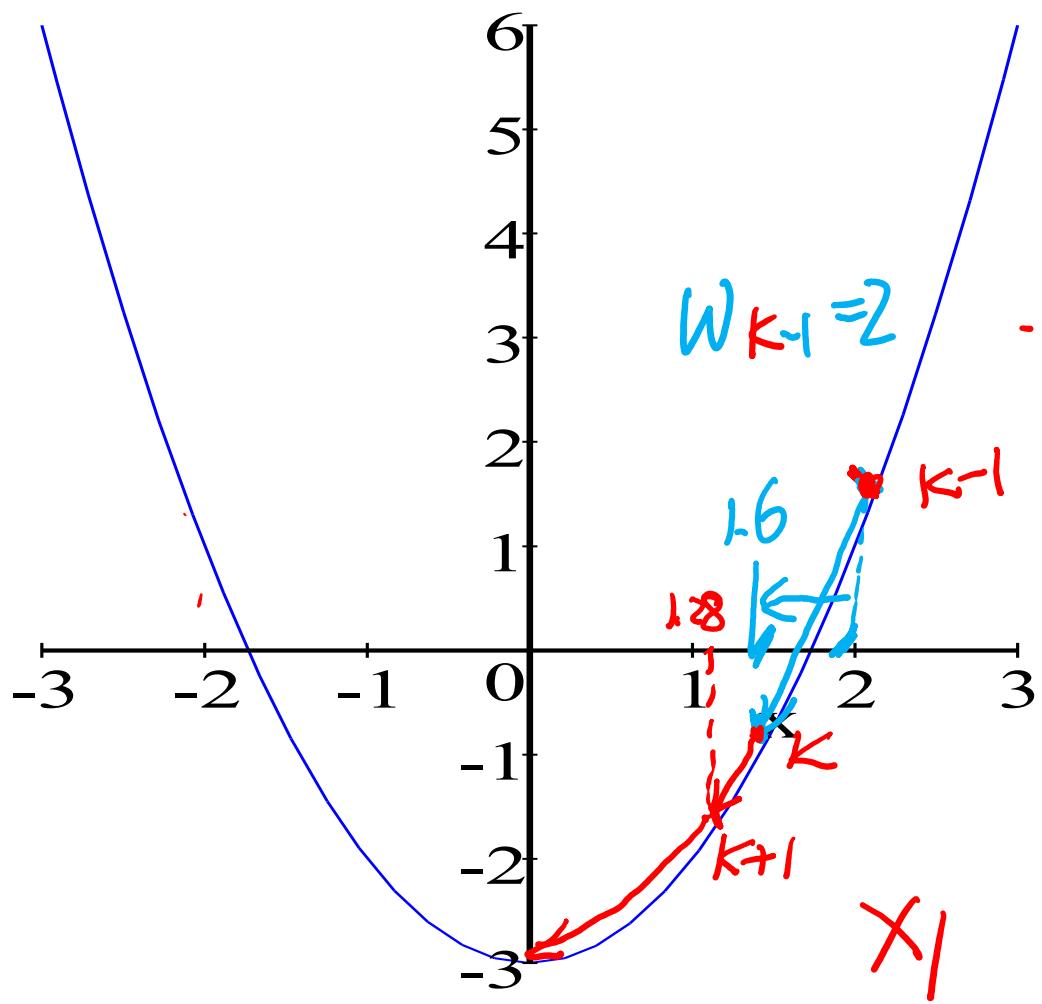
$$l'(w) = 2w$$

$$w_k = w_{k-1} - 2 \propto w_{k-1}$$





$$w_k = w_{k-1} - 2 \propto w_{k-1}$$



$$w_{k-1} = 2$$

$$\alpha = 0.1$$

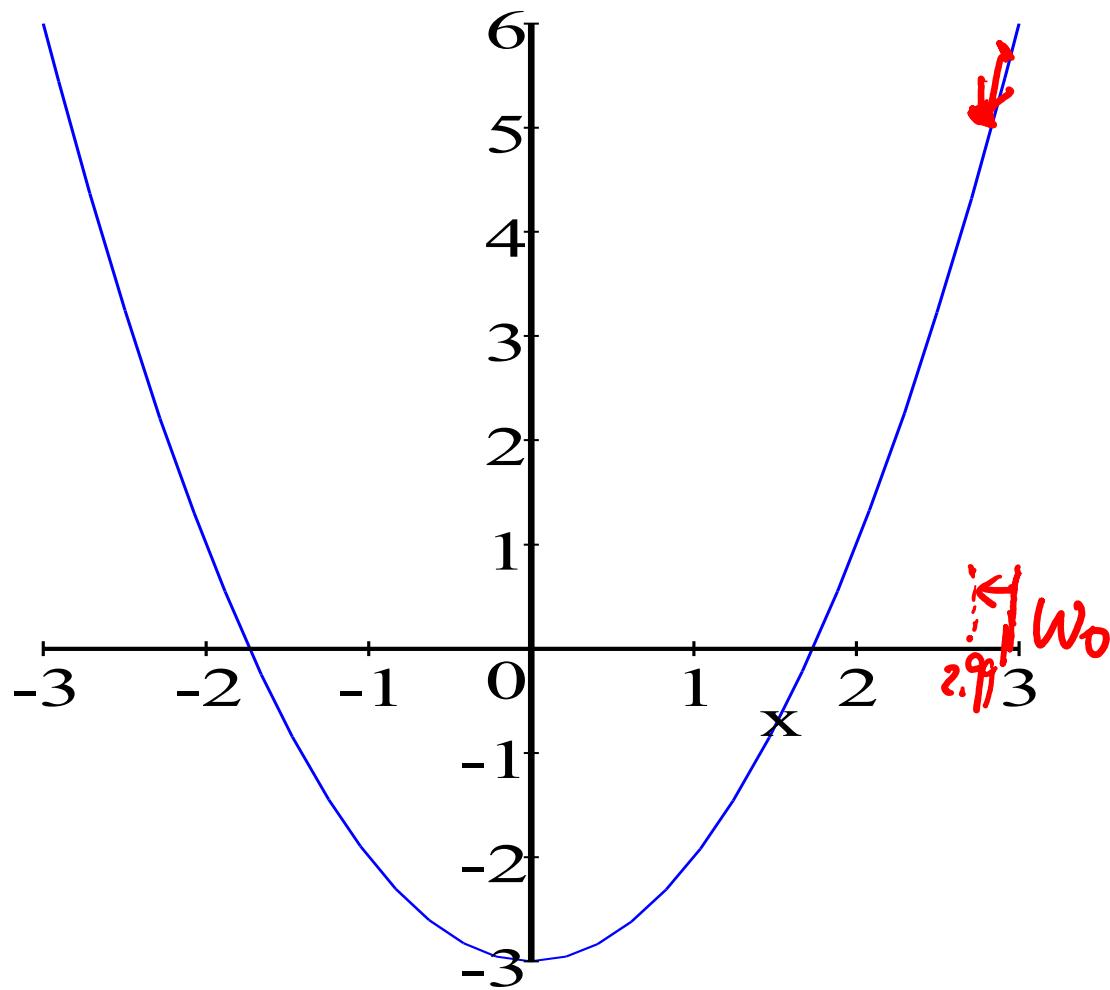
$$\begin{aligned} w_k &= 2 - 2 \times 0.1 / \sqrt{2} \\ &= 1.6 \end{aligned}$$

$$\begin{aligned} w_{k+1} &= 1.6 - \\ &\quad 2 \times 0.1 \times 1.6 \\ &= 1.28 \end{aligned}$$

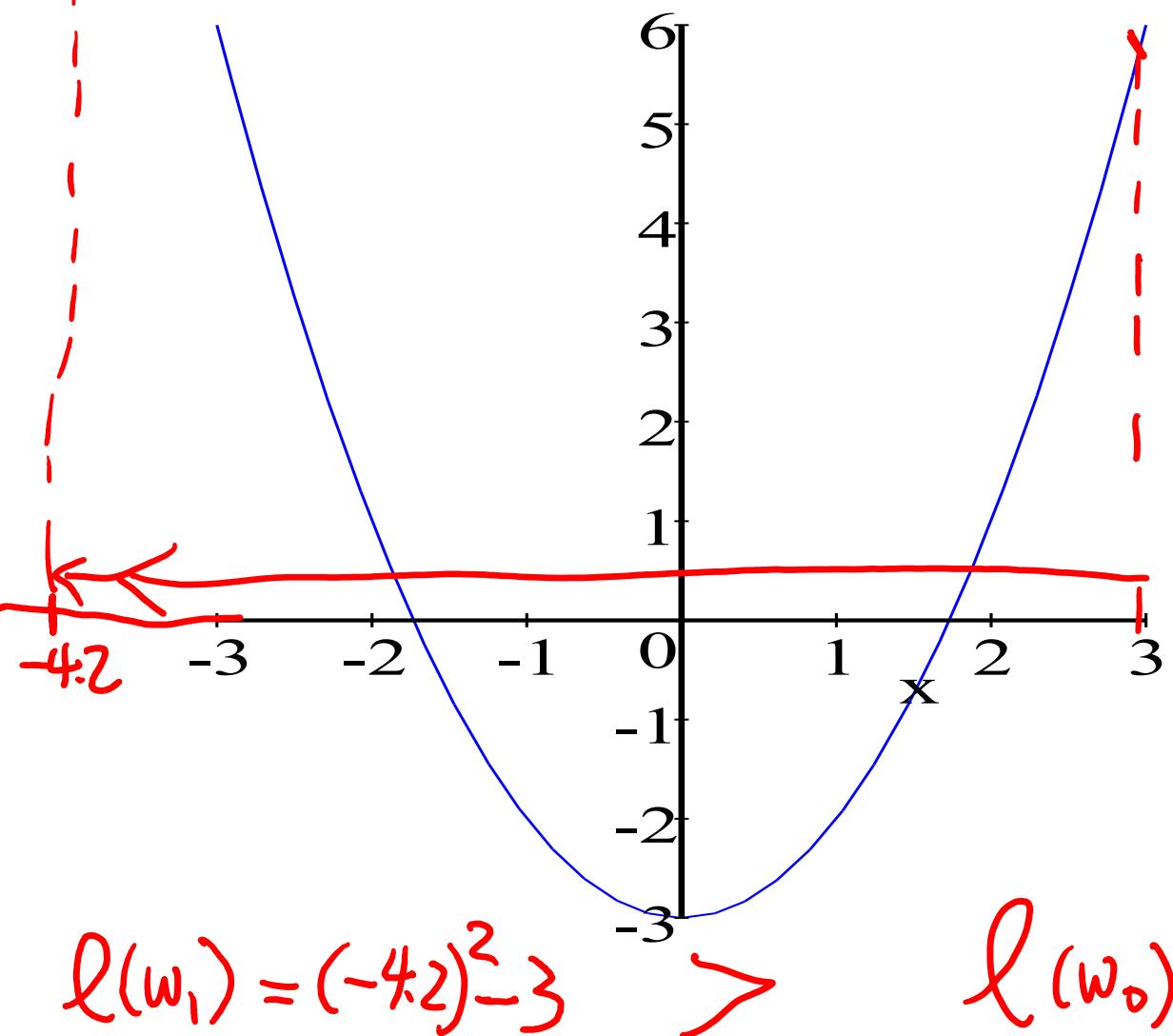
Gradient Descent (Iteratively Optimize)

- Learning Rate Matters
- Objective function matters
- Starting point matters

$$w_k = w_{k-1} - 2 \propto w_{k-1}$$



$$w_k = w_{k-1} - 2 \propto w_{k-1}$$



$$w_0 = 3$$

$$\alpha = 0.7$$

$$\alpha = 1.2$$

$$w_1 = 3 - 2 \times 1.2 \times 3$$

$$= -4.2$$

Gradient Descent (Iteratively Optimize)

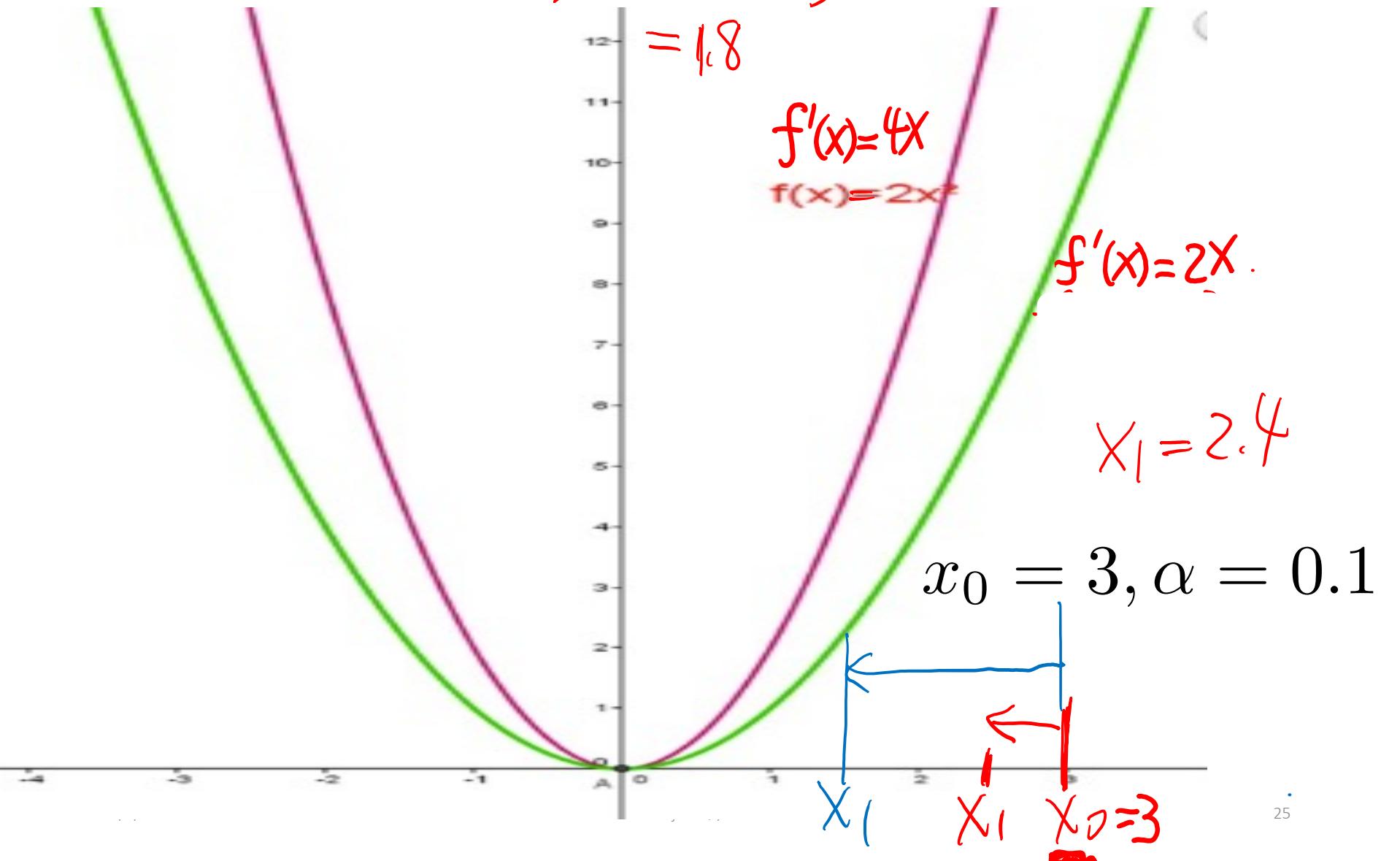
- Learning Rate Matters
- Objective function matters
- Starting point matters

$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$



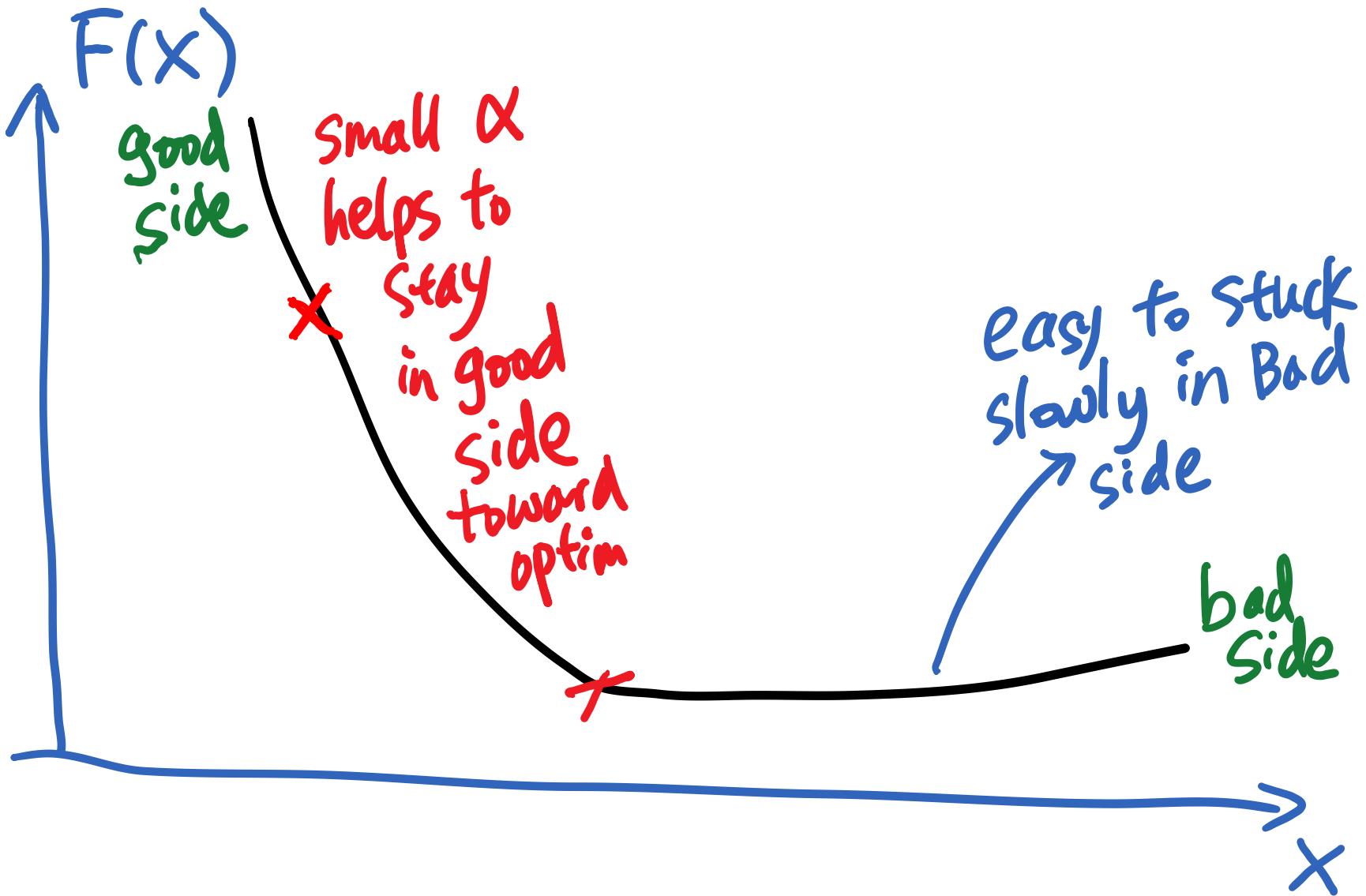
$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$

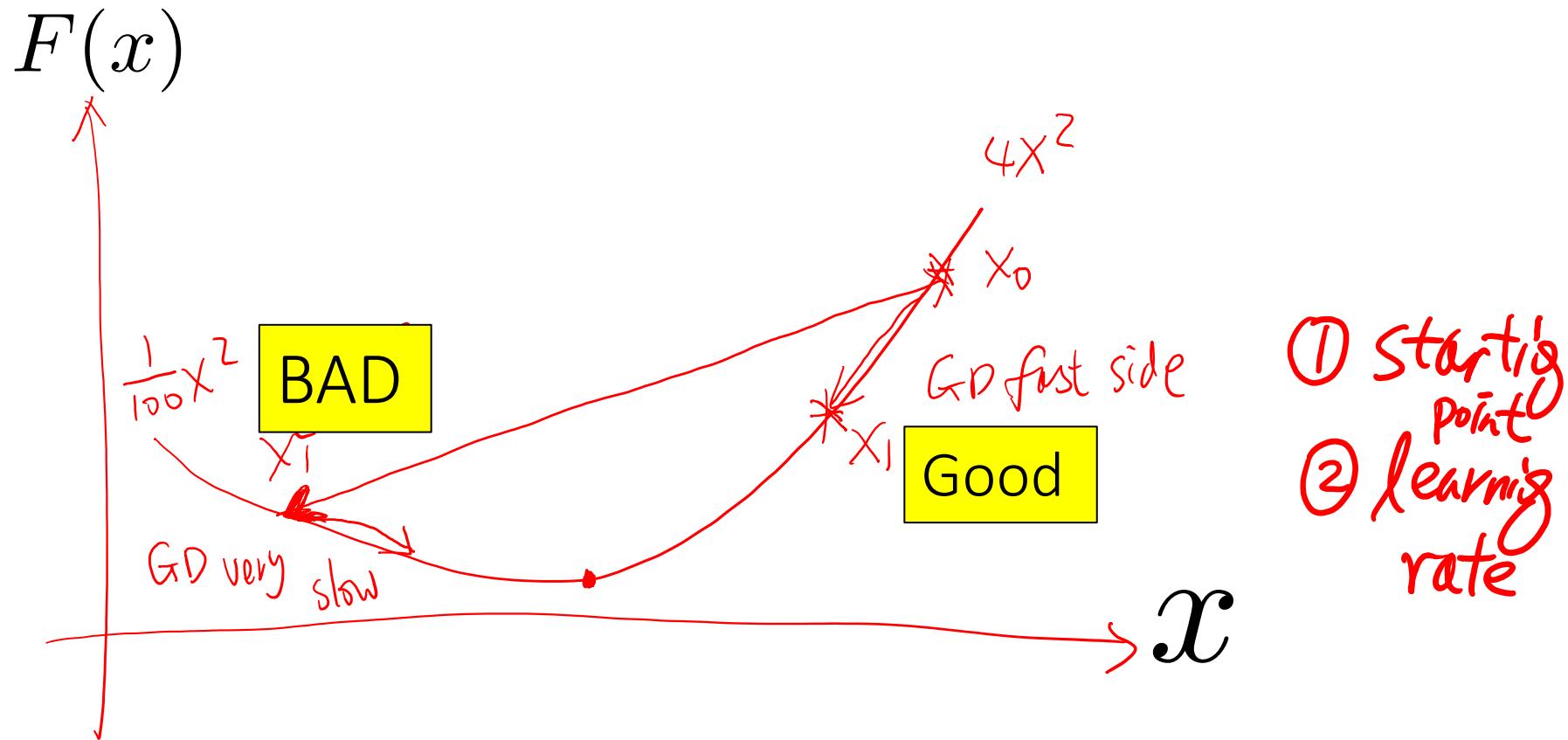
$$\begin{aligned}x_1 &= 3 - 4 \times 0.1 \times 3 \\&= 1.8\end{aligned}$$



Gradient Descent (Iteratively Optimize)

- Learning Rate Matters
- Objective function matters
- Starting point matters



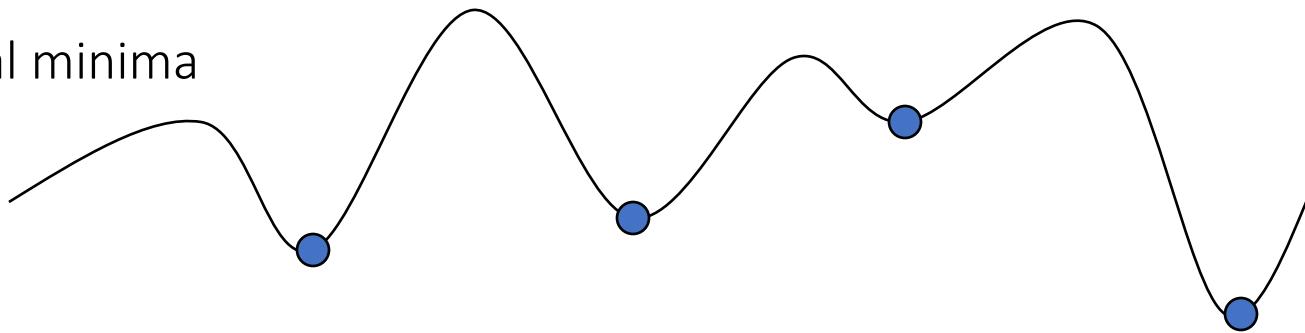


During optimization, We don't want to jump from the good side to the bad side

Comments on Gradient Descent Algorithm

- Works on any objective function $F(\underline{x})$
 - as long as we can evaluate the gradient
 - this can be very useful for minimizing complex functions

- Local minima



- Can have multiple local minima
- (note: for LR, its cost function only has a single global minimum, so this is not a problem)
- If gradient descent goes to the closest local minimum:
 - solution: random restarts from multiple places in weight space

Thank You



UVA CS 4774 : Machine Learning

Lecture 4: More optimization for Linear Regression

Module 2

Dr. Yanjun Qi

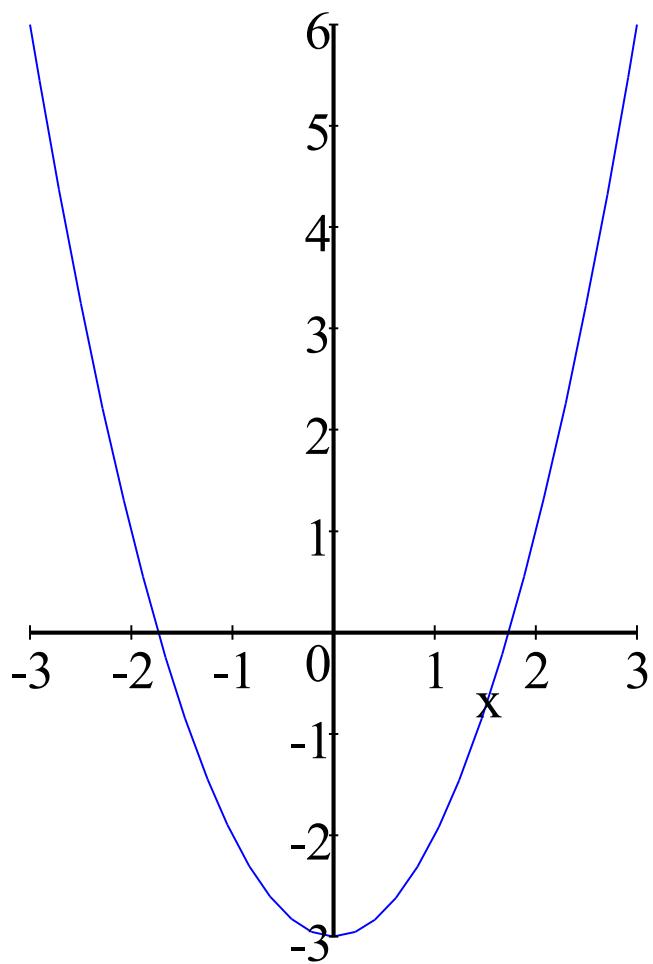
University of Virginia
Department of Computer Science

Our concrete example with Gradient Descent

$$l(w) = w^2 - 3$$

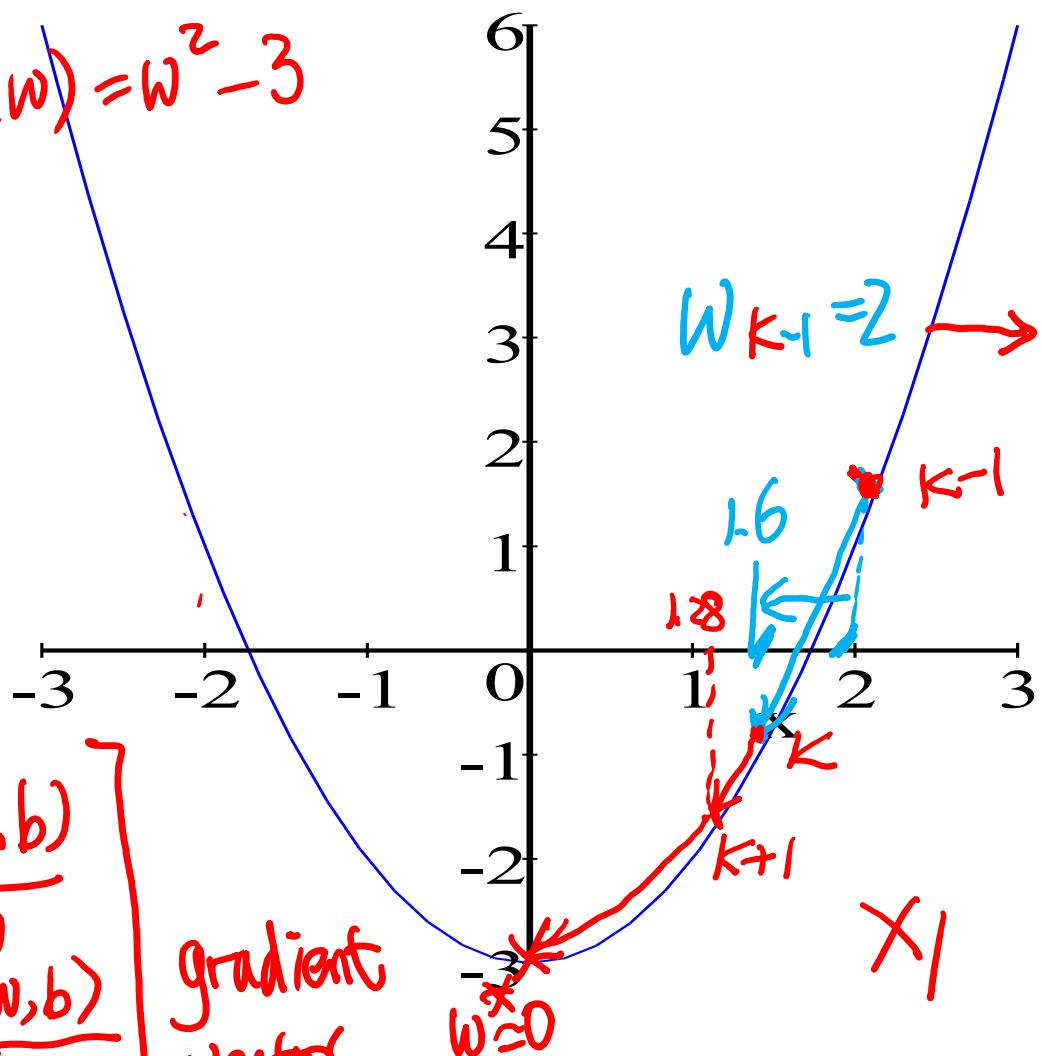
$$l'(w) = 2w$$

$$w_k = w_{k-1} - 2 \propto w_{k-1}$$



$$w_k = w_{k-1} - \alpha w_{k-1}$$

$$l(w) = w^2 - 3$$



$$w_{k-1} = 2$$

$$\alpha = 0.1$$

$$w_k = 2 - 2 \times 0.1 \times 2$$

$$= 1.6$$

$$w_{k+1} = 1.6 -$$

$$2 \times 0.1 \times 1.6$$

$$= 1.28$$

Iteratively Optimize: Gradient Descent (GD) and Stochastic GD for Linear Regression Training

Review: Loss function of Least Square LR

$$\begin{aligned} J(\theta) &= \frac{1}{2} \sum_{i=1}^n (f(\mathbf{x}_i) - y_i)^2 \\ &= \frac{1}{2} \left(\theta^T X^T X \theta - \theta^T X^T \bar{y} - \bar{y}^T X \theta + \bar{y}^T \bar{y} \right) \end{aligned}$$

Extra: more in note PDF \Rightarrow matrix calculus, partial deri \Rightarrow Gradient

$$\nabla_{\theta} (\theta^T X^T X \theta) = 2X^T X \theta \quad (\text{P24})$$

$$\nabla_{\theta} (-2\theta^T X^T Y) = -2X^T Y \quad (\text{P24})$$

$$\nabla_{\theta} (Y^T Y) = 0$$

$$\Rightarrow \nabla_{\theta} J(\theta) = \boxed{\bar{X}^T \bar{X} \theta - \bar{X}^T Y}$$

$$\nabla_{\theta} J(\theta)$$

$$= X^T X \theta - X^T \bar{y}$$

$$= X^T (X \theta - \bar{y})$$

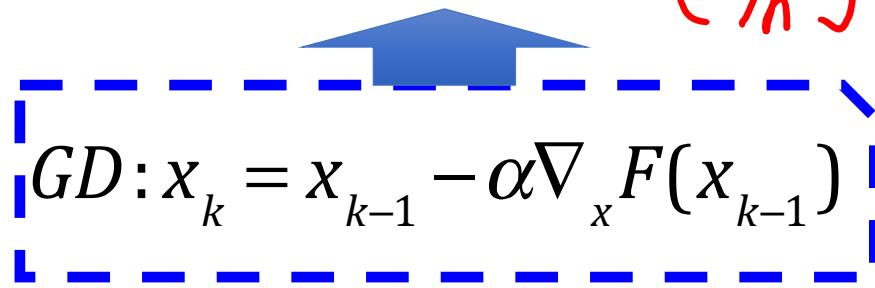
Linear Regression Trained with batch GD

- A Batch **gradient descent** algorithm:

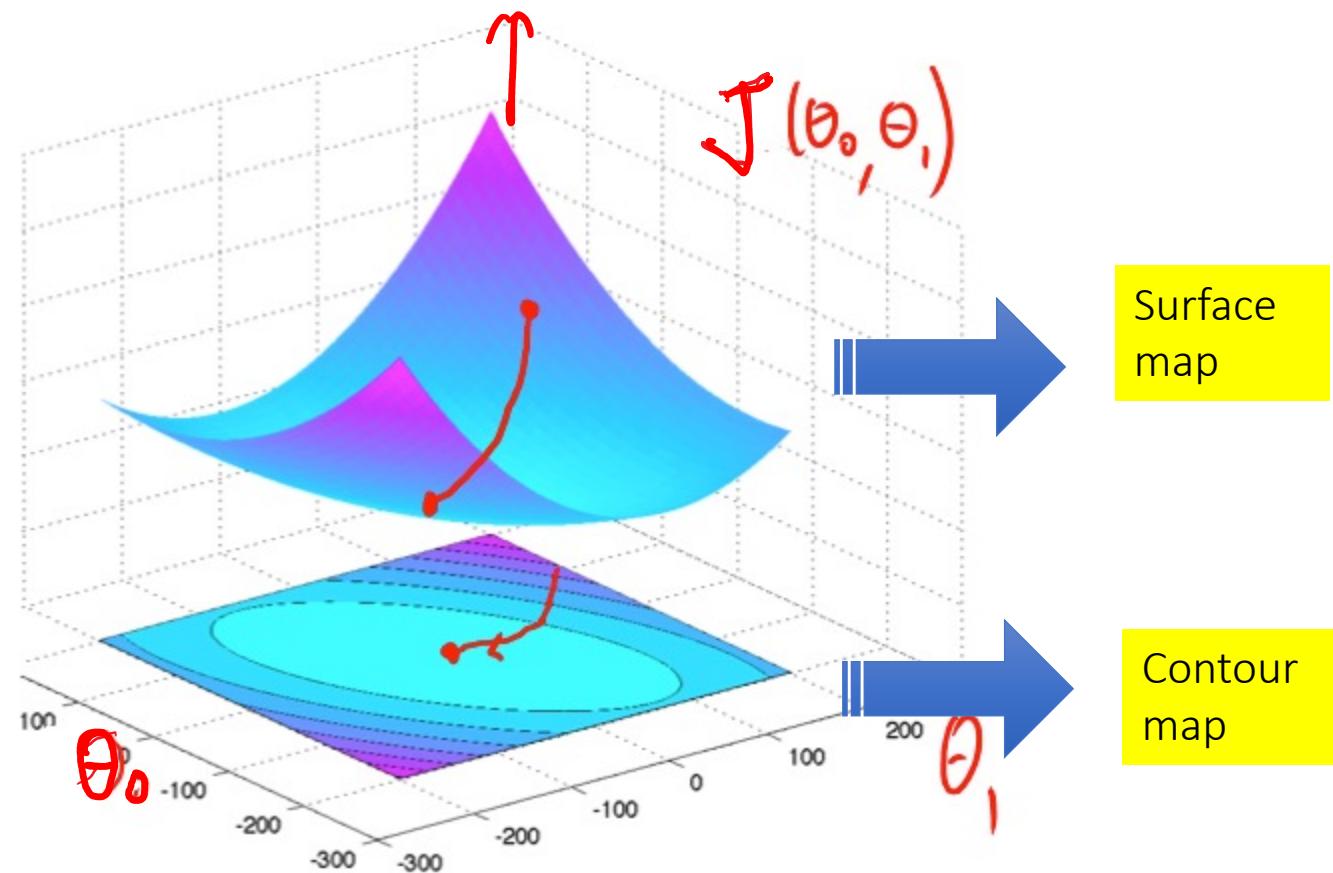
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$

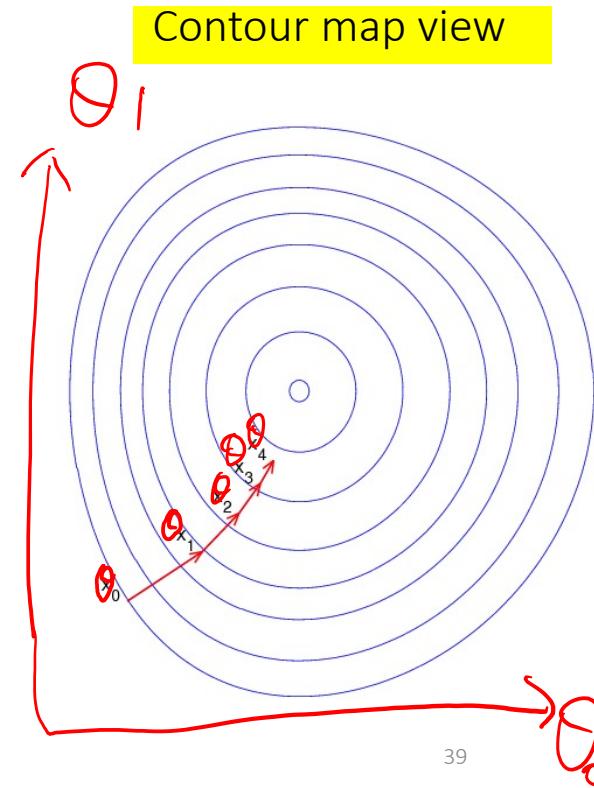
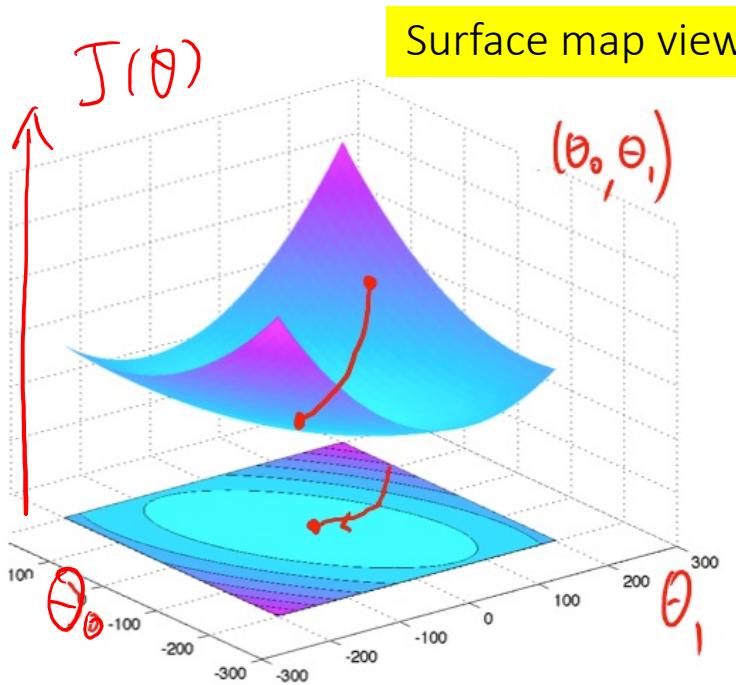
$$\begin{aligned}\theta^{t+1} &= \theta^t - \alpha \nabla_{\theta} J(\theta^t) \\ &= \theta^t + \alpha X^T (\bar{y} - \underline{\underline{X \theta^t}})\end{aligned}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} - \begin{bmatrix} -x_1^T \theta - \\ -x_2^T \theta - \\ \vdots \\ -x_n^T \theta - \end{bmatrix}$$

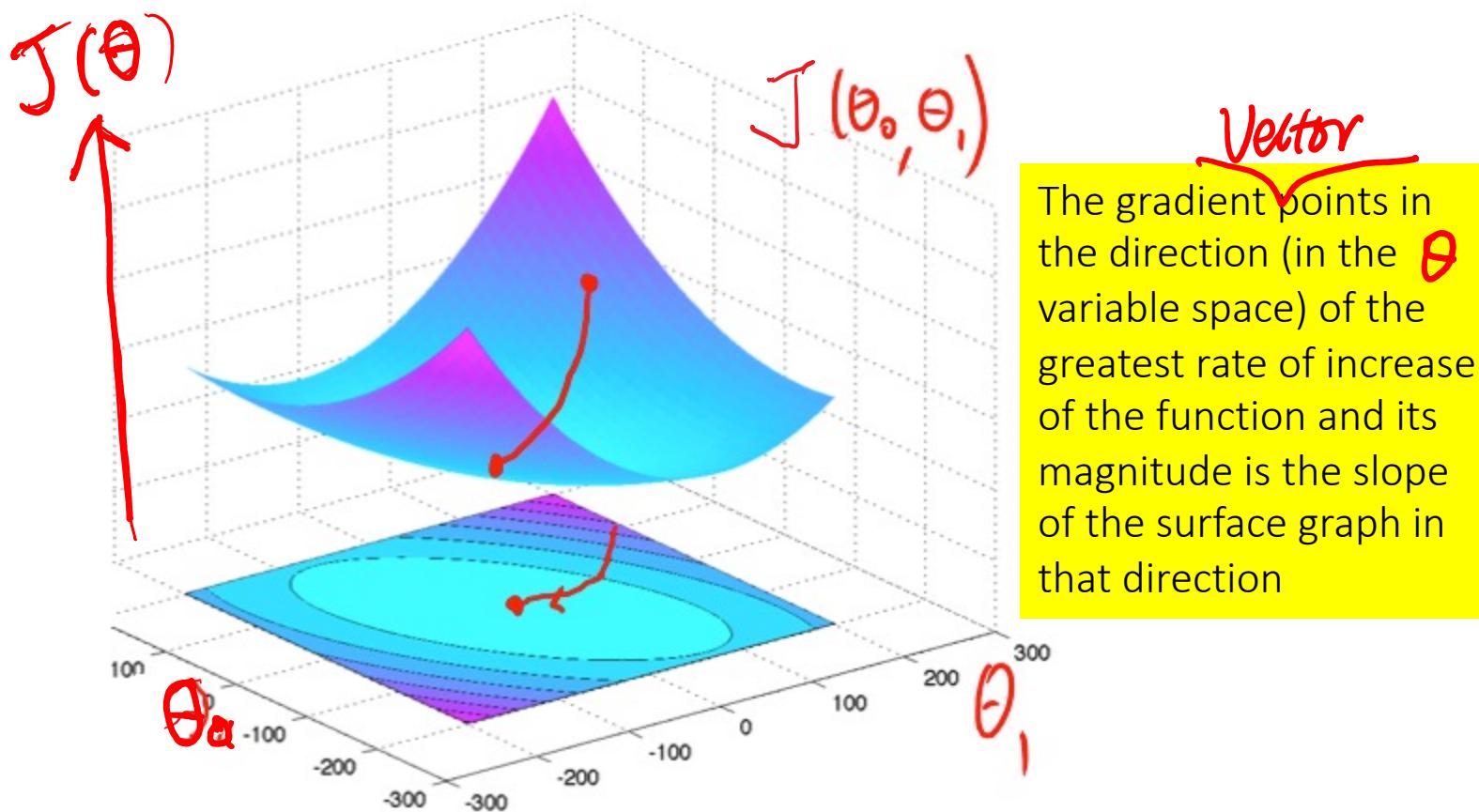

$$GD: x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$

Review: two ways of Illustrating an Objective Function (for two variables case)



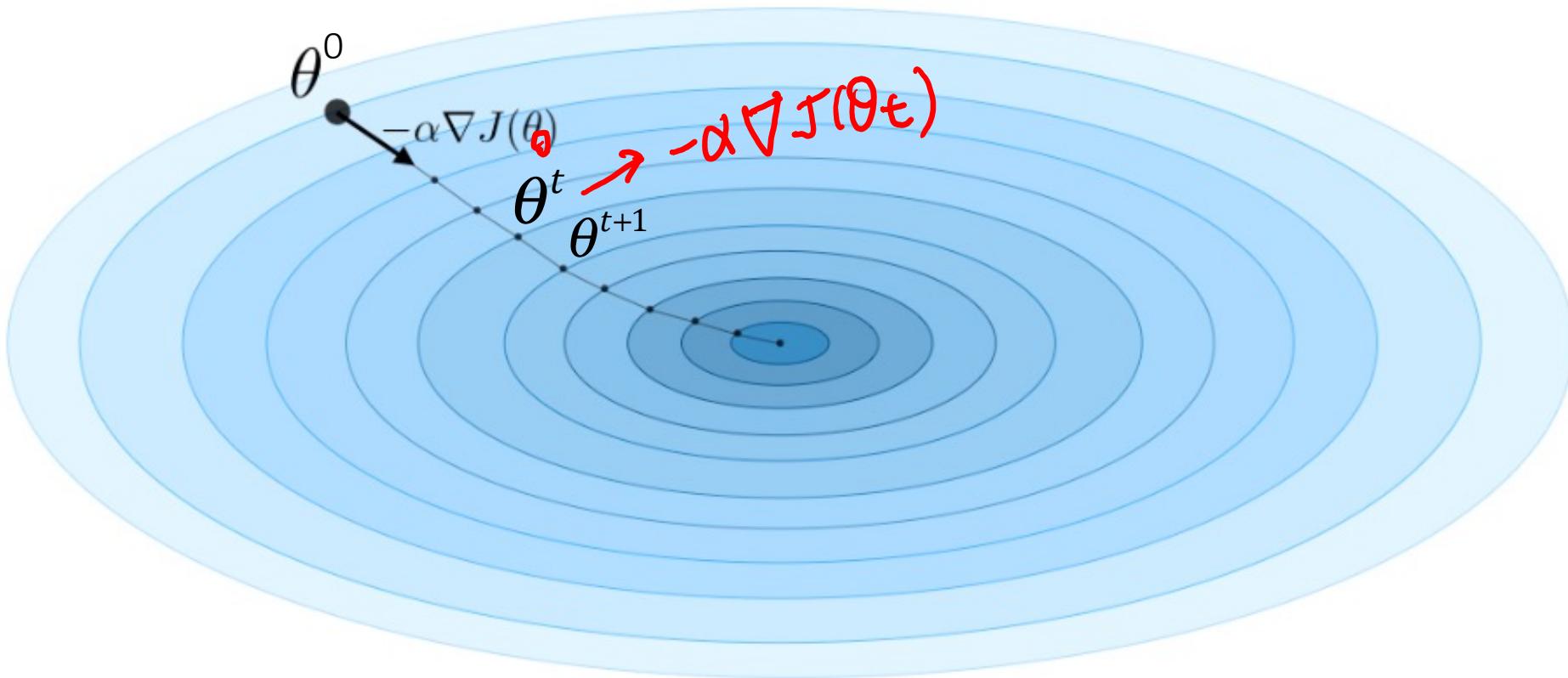


Review: two ways of Illustrating an Objective Function (for two variables case)



Gradient Descent (Steepest Descent)

$$\begin{aligned}\theta^{t+1} &= \theta^t - \alpha \nabla_{\theta} J(\theta^t) \\ &= \theta^t + \alpha X^T (\bar{y} - X \theta^t)\end{aligned}$$



To find a local minimum of a function using gradient descent, one takes **steps proportional to the negative of the gradient of the function at the current point.**

$$\theta^{t+1} = \theta^t - \alpha \nabla_{\theta} J(\theta^t)$$

https://s3.amazonaws.com/assets.datacamp.com/blog_assets/Numpy_Python_Cheat_Sheet.pdf

$$= \theta^t + \alpha X^T (\bar{y} - X\theta^t)$$

$$= \theta^t + \alpha \sum_{i=1}^n (y_i - \bar{x}_i^T \theta^t) \bar{x}_i$$

a vector of

errors on
each point @ θ_t

$$\overset{\leftarrow}{\theta}^{t+1} = \overset{\leftarrow}{\theta}^t + \alpha \underset{p \times 1}{\cancel{X}}^T (\bar{y} - \underset{n \times 1}{\cancel{X}} \overset{\leftarrow}{\theta}^t)$$

$p \times 1$

$p \times 1$

$|X|$

$p \times n$

$n \times 1$

$n \times p$

$p \times 1$

$n \times 1$

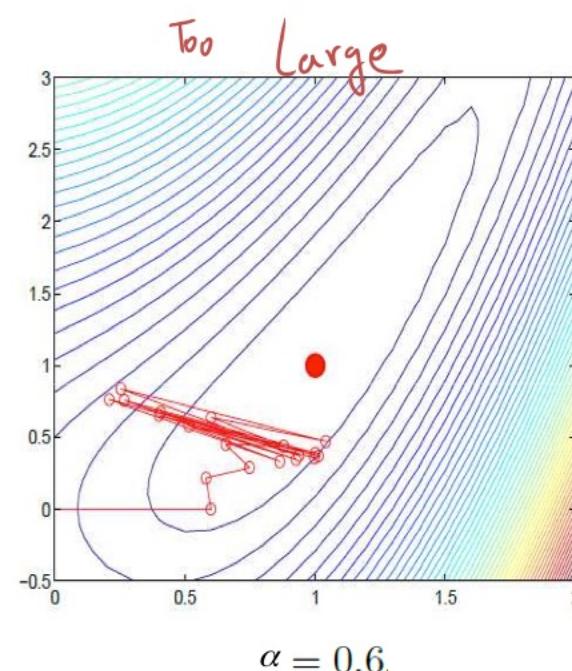
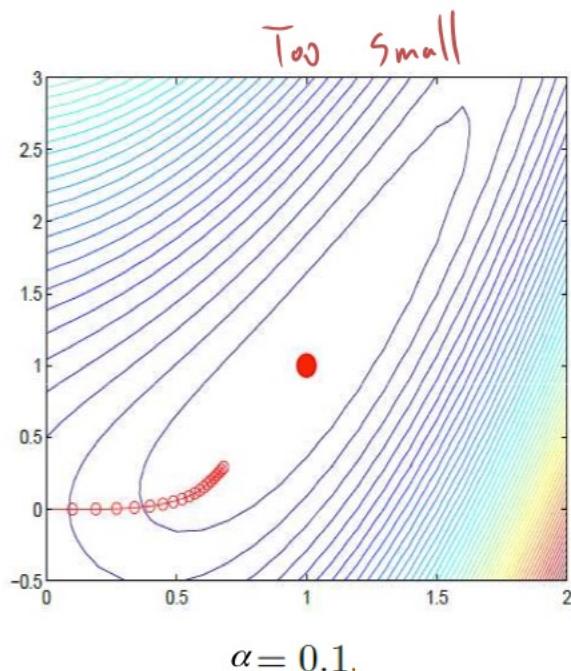
$p \times 1$

$p \times 1$

<https://numpy.org/doc/stable/reference/routines.linalg.html>

You can use panda: pandas is a thin abstraction layer on top of numpy

Choosing the Right [Learning-Rate] is critical



Training with batch GD

- Gradient of Cost Function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

- Consider a **gradient descent** algorithm and reformulate:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$
$$\begin{aligned} \theta^{t+1} &= \theta^t - \alpha \nabla_{\theta} J(\theta^t) \\ &= \theta^t + \alpha X^T (\bar{y} - X \theta^t) \\ &= \theta^t + \alpha \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta^t) \mathbf{x}_i \end{aligned}$$

Review: Gradient Vector of Linear Regression Loss

- The Cost Function:

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2$$

- Consider a **gradient descent** algorithm and reformulate:

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_p \end{bmatrix}$$
$$\begin{aligned}\theta^{t+1} &= \theta^t - \alpha \nabla_{\theta} J(\theta^t) \\ &= \theta^t + \alpha X^T (\bar{y} - X \theta^t) \\ &= \theta^t + \alpha \sum_{i=1}^n (y_i - \bar{\mathbf{x}}_i^T \theta^t) \bar{\mathbf{x}}_i\end{aligned}$$


LR with Stochastic GD \rightarrow

$$\theta^t + \alpha \bar{x}^T (y - \bar{x}\theta^t)$$

$n \sim 2M$ points
 $p \sim 3k$ features

$\sum_{i=1}^n$ $P \times n$ $n \times 1$

- Batch GD rule:

$$\theta^{t+1} = \theta^t + \alpha \sum_{i=1}^n (y_i - \bar{x}_i^T \theta^t) \bar{x}_i$$

- For a single training point (i-th), we have:

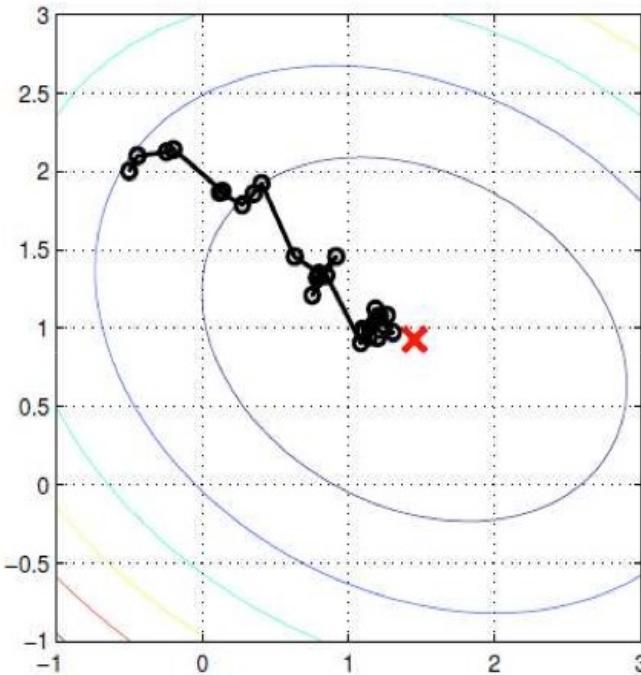
$$\theta^{t+1} = \theta^t + \alpha (y_i - \bar{x}_i^T \theta^t) \bar{x}_i$$

$\theta^0 \longrightarrow \theta^1 \longrightarrow \theta^2$

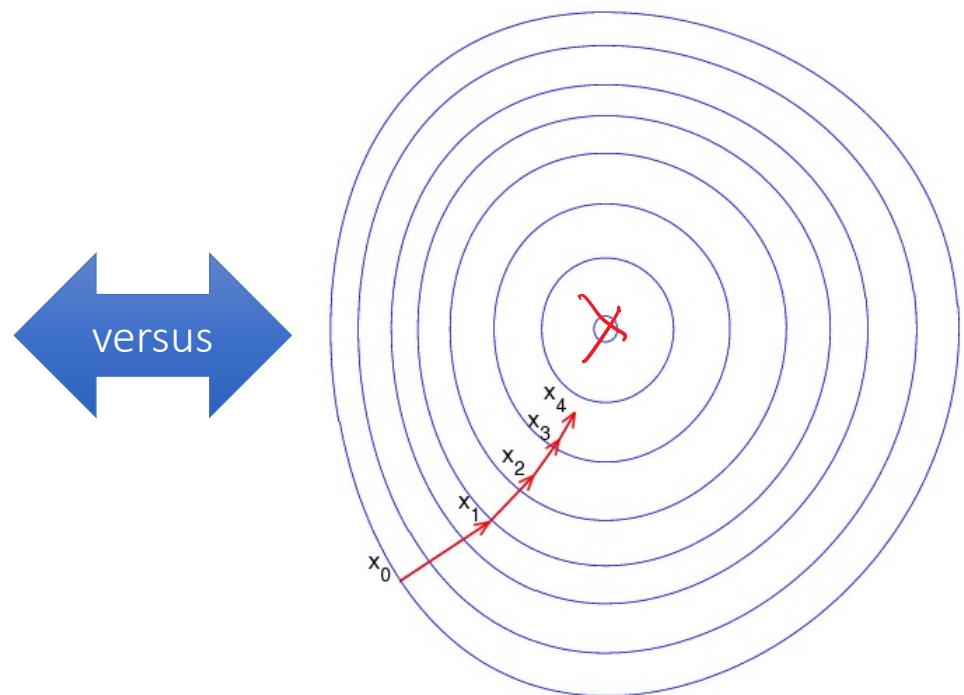
- A "stochastic" descent algorithm
- Can be used as an **on-line** algorithm

Stochastic gradient descent VS. Gradient Descent

SGD



GD



Stochastic gradient descent : More variations

- Single-sample:

$$\theta^{t+1} = \theta^t + \alpha(y_i - \bar{\mathbf{x}}_i^T \theta^t) \bar{\mathbf{x}}_i$$

- Mini-batch:

$$\theta^{t+1} = \theta^t + \alpha \sum_{j=1}^B (y_{Ij} - \bar{\mathbf{x}}_{Ij}^T \theta^t) \bar{\mathbf{x}}_{Ij}$$

e.g. $B=15$

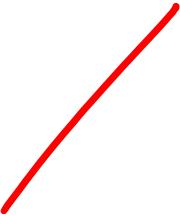
$$= \theta^t + \alpha X_B^T (\bar{y} - X \theta^t)$$

Stochastic gradient descent (more)

- Very useful when training with massive datasets , e.g. not fit in main memory
- Very useful when training data arrives online (e.g. streaming)..
- SGD can be used for offline training, by repeated cycling through the whole data
 - Each such pass over the whole data → an epoch !
- In offline case, often better to use mini-batch SGD
 - E.g. $B=64$
 - $B=1$ standard SGD
 - $B=N$ standard batch GD

Mini-batch: (stochastic gradient descent)

- Compute the gradient on a small mini-batch of samples (e.g. $B=32/64/\dots$)
- Much faster computationally than single point SGD (better use of computer architecture like GPU)

Epoch  Cover all examples

GD : one update

SGD : n_{tr} updates

miniSGD : n_{tr}/B updates

Low per-step cost, fast convergence and
perhaps less prone to local optimum

(Stochastic) Gradient Descent (Iteratively Optimize)

- Learning Rate Matters
- Starting point matters
- Objective function matters
- Stop criterion matters!

- Train MSE Error to observe:

$$J_{train_MSE}^t = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \boldsymbol{\theta}^t - y_i)^2$$

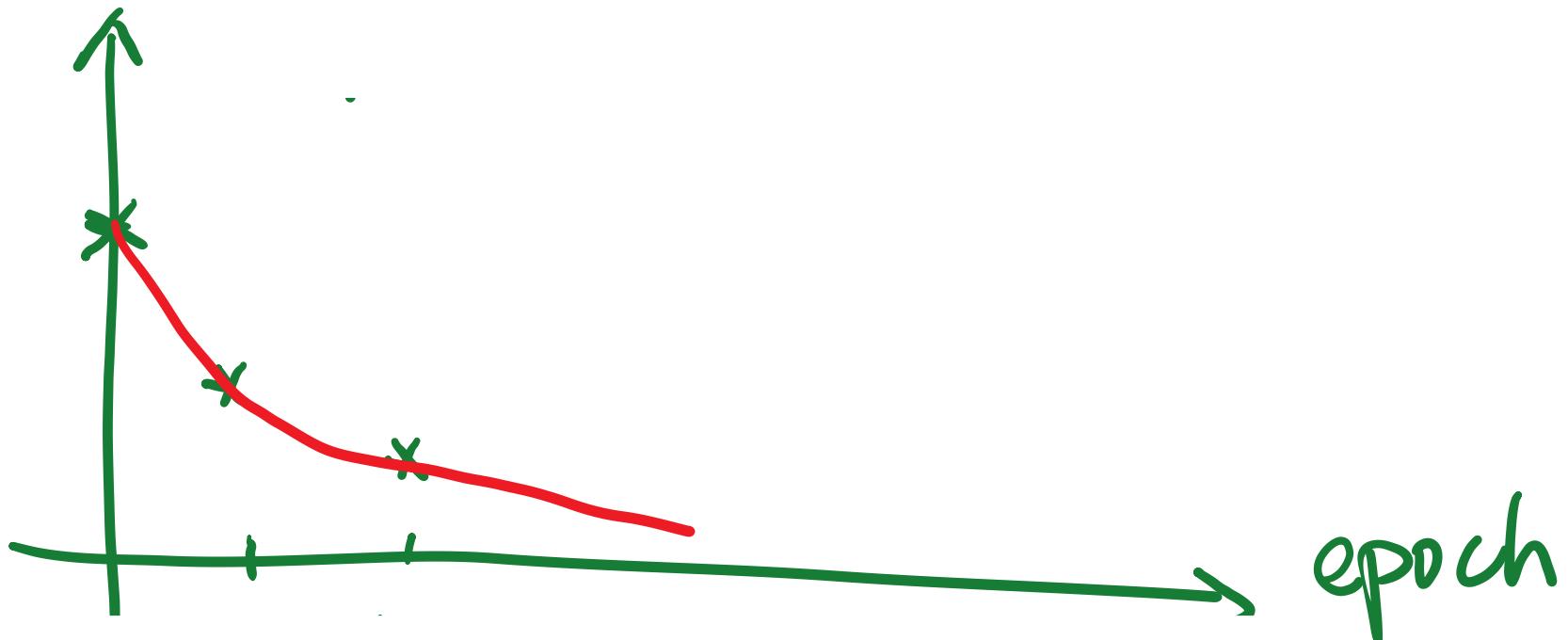
In many situations, visualizing Train-MSE can be helpful to understand the behavior of your method, e.g., how it decreases with epochs, ...

In Homework, when we ask for plots of training error, we ask for the MSE per-sample train errors; Because it is comparable to test MSE error (later to cover).

One good plot for GD / SGD

mean-training-loss

$$J_{train-MSE} = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \theta^* - y_i)^2$$



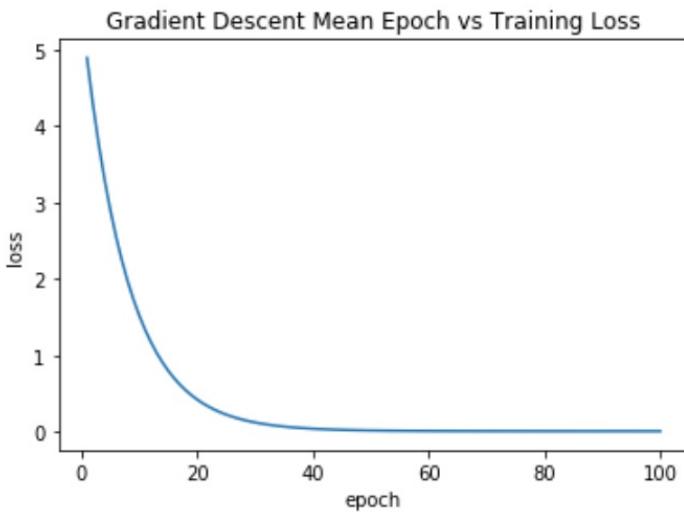
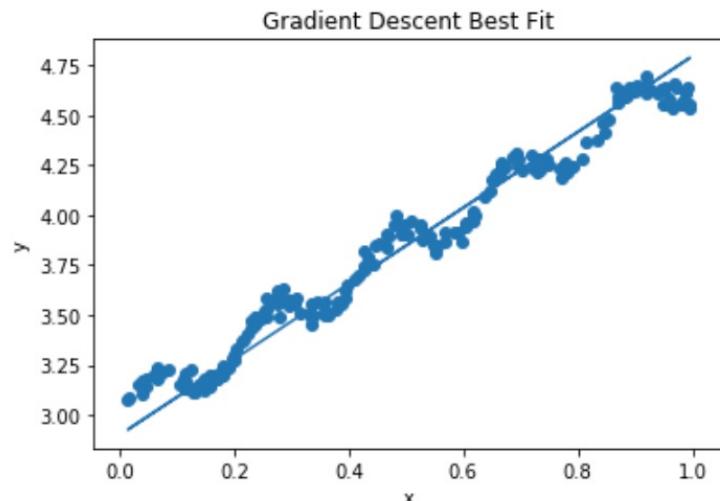
Each pass of SGD repeated cycling through all samples in the whole train → an epoch !

When to stop (S)GD ?

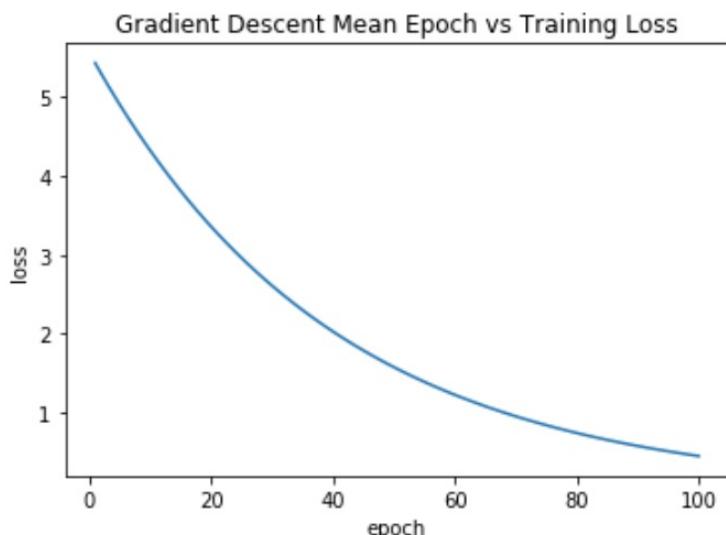
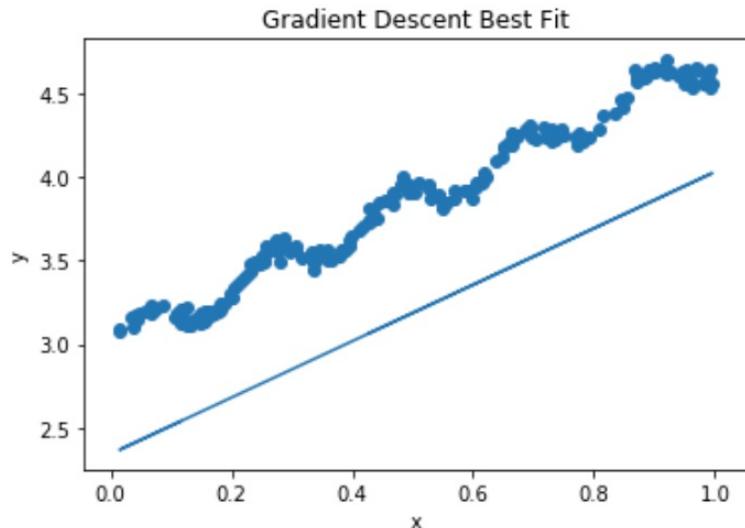
- Lots of stopping rules in the literature,
- there are advantages and disadvantages to each, depending on context
- E.g., a **predetermined maximum number of iterations**
- E.g., stop when the improvement drops below a threshold
-

e.g. HW1 discussions: Stopping and Learning Rates

```
theta = gradient_descent(X, Y, 0.05, 100)
plotPredict(X, Y, theta[-1], "Gradient Descent Best Fit")
plot_training_errors(X, Y, theta, "Gradient Descent Mean Epoch Loss")
```



```
theta = gradient_descent(X, Y, 0.01, 100)
plotPredict(X, Y, theta[-1], "Gradient Descent Best Fit")
plot_training_errors(X, Y, theta, "Gradient Descent Mean Epoch Loss")
```



Summary so far: Four ways to learn LR

- Normal equations

$$\theta^* = (X^T X)^{-1} X^T \bar{y}$$

- Pros: a single-shot algorithm! Easiest to implement.
- Cons: need to compute pseudo-inverse $(X^T X)^{-1}$, expensive, numerical issues (e.g., matrix is singular ..), although there are ways to get around this ...

- GD

$$\theta^{t+1} = \theta^t + \alpha X^T (\bar{y} - X\theta) = \theta^t + \alpha \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta^t) \mathbf{x}_i$$

- Pros: easy to implement, conceptually clean, guaranteed convergence
- Cons: batch, often slow converging

- Stochastic GD and miniBatch

$$\theta^{t+1} = \theta^t + \alpha (y_i - \mathbf{x}_i^T \theta^t) \mathbf{x}_i$$

frequent $\xrightarrow{\text{N updates}}$ *noise*

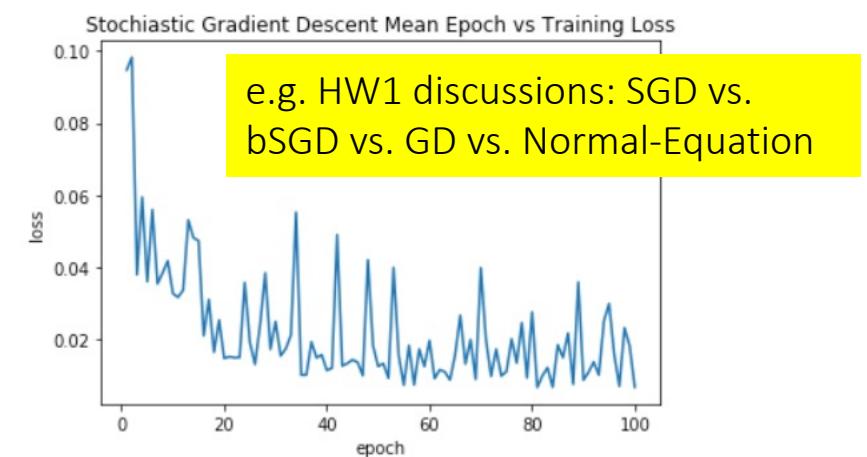
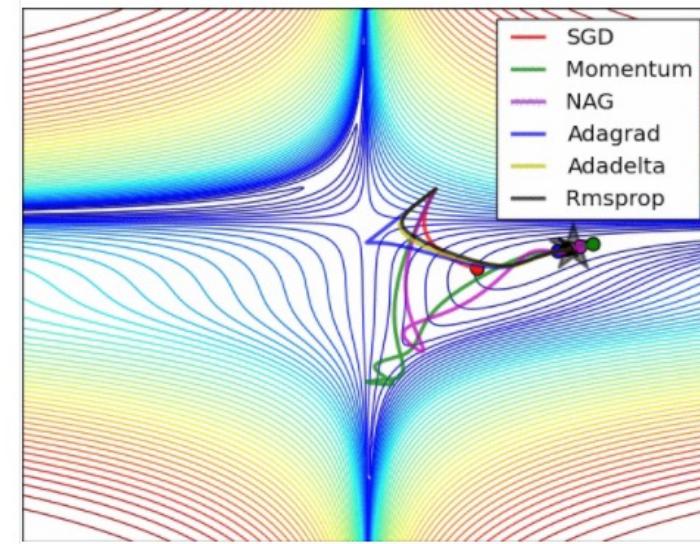
$$\theta^{t+1} = \theta^t + \alpha \sum_{j=1}^B (y_{Ij} - \mathbf{x}_{Ij}^T \theta^t) \mathbf{x}_{Ij}$$

- Pros: on-line, low per-step cost, fast convergence and perhaps less prone to local optimum
- Cons: convergence to optimum not always guaranteed

- Challenges
- Gradient descent optimization algorithms
 - Momentum
 - Nesterov accelerated gradient
 - Adagrad
 - Adadelta
 - RMSprop
 - Adam
 - AdaMax
 - Nadam
 - AMSGrad
 - Other recent optimizers
 - Visualization of algorithms
 - Which optimizer to use?
- Parallelizing and distributing SGD
 - Hogwild!
 - Downpour SGD
 - Delay-tolerant Algorithms for SGD
 - TensorFlow
 - Elastic Averaging SGD
- Additional strategies for optimizing SGD
 - Shuffling and Curriculum Learning
 - Batch normalization
 - Early Stopping
 - Gradient noise

- More about SGD:

- Popular optimization for NOW
- Many advanced variations
- <https://ruder.io/optimizing-gradient-descent/>
- <https://arxiv.org/abs/1609.04747>



GD and SGD code cell run @

<https://colab.research.google.com/drive/1Gchka0n69mTwRvZUEpPBKK1BgUq94qPk?usp=sharing>

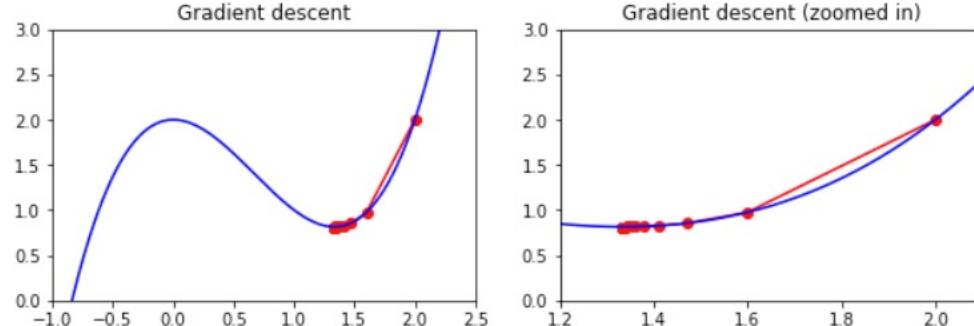
Revised from:

<http://cs229.stanford.edu/>

https://github.com/dtnewman/stochastic_gradient_descent

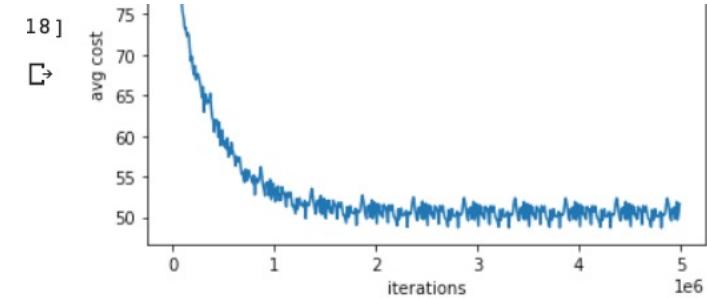
```
# use scipy fmin function to find ideal step size.  
n_k = fmin(f2, 0.1, (x_old, s_k), full_output = False, disp = False)
```

```
plt.title("Gradient descent (zoomed in)")  
plt.show()
```

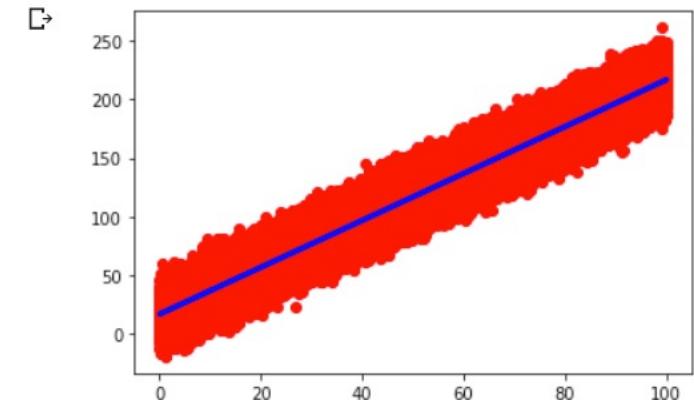


2/3/22

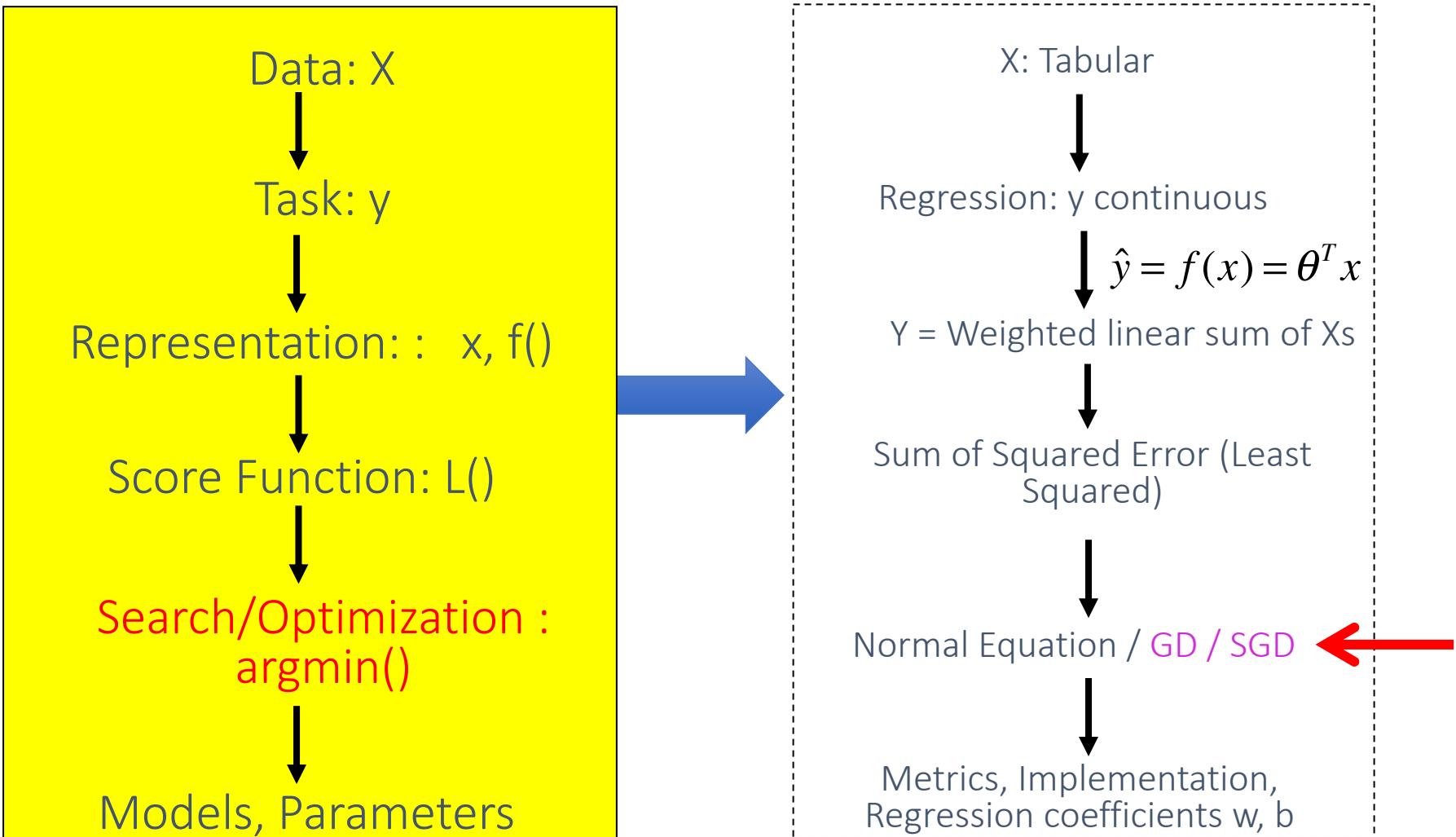
Dr. Yanjun Qi / UVA CS



```
f = lambda x: x*x+17+np.random.randn(len(x))*10  
  
x = np.random.random(500000)*100  
y = f(x)  
hy = theta_new[0] + theta_new[1]*x  
plt.plot(x,hy,c="b", linewidth=3)  
plt.scatter(x,y,c="r")  
plt.show()
```



Recap : GD and SGD for Multivariate Linear Regression



Thank You



EXTRA I

In Case you are interested in more advanced details!

B

Varying the value B In

$$\theta^{t+1} = \theta^t + \alpha \sum_{j=1}^B (y_{Ij} - \vec{x}_{Ij}^T \theta) \vec{x}_{Ij}$$

I	$I < B < n$	n
<u>SGD</u> <u>very noisy GD update</u>	<u>miniB-SGD</u> <u>a bit noisy GD update</u>	<u>GD</u> <u>precise GD update</u>
<u>(low memory cost)</u>	<u>middle memory cost</u>	<u>high memory cost</u>
<u>$O(1)$ one gradient calc cost</u>	<u>if multi-parallel B threads</u> <u>$O(1)$ one gradient calc cost</u>	<u>very costly gradient calculation</u>

Extra: Direct (normal equation) vs. Iterative (GD, SGD,) methods

- Direct methods: we can achieve the solution in a single step by solving the normal equation
 - Using Gaussian elimination or QR decomposition, we converge in a finite number of steps
 - It can be infeasible when data are streaming in real time, or of very large amount
- Iterative methods: stochastic GD or GD
 - Converging in a limiting sense
 - But more attractive in large practical problems
 - Caution is needed for deciding the learning rate

Stochastic gradient descent (Pros)

- Efficiency: Good approximation of Gradient:
 - Intuitively fairly good estimation of the gradient by looking at just a few examples
 - Carefully evaluating precise gradient using large set of examples is often a waste of time (because need to calculate the gradient of the next t anyway)
 - Better to get a noisy estimate and move rapidly in the parameter space
- SGD is often less prone to stuck in shallow local minima
 - Because of the certain “noise”,
 - popular for nonconvex optimization cases

Extra: Convergence rate

- **Theorem:** the steepest descent / GD equation algorithm converge to the minimum of the cost characterized by normal equation:

$$\theta^{(\infty)} = (X^T X)^{-1} X^T y$$

If the learning rate parameter satisfy \rightarrow

$$0 < \alpha < 2/\lambda_{\max}[X^T X]$$

- A formal analysis of GD-LR need more math; in practice, one can use a small a , or gradually decrease a .

$$\alpha_0 = 0.05$$

Extra: Computational Cost (Naïve..)

$$\hat{\theta}^* = \left(\begin{matrix} X^T \\ \Sigma_{n \times p} \end{matrix} \right)^{-1} X^T \vec{y}$$

$$X^T X : O(p^2 n)$$

$$(X^T X)^{-1} : O(p^3)$$

mostly about Memory Cost →

Interesting discussion in:
<https://stackoverflow.com/questions/10326853/why-does-lm-run-out-of-memory-while-matrix-multiplication-works-fine-for-coeffic>

$$O(n p^2 + p^3)$$

when $n \gg p$, matrix multi
slower than inversion

$$\nabla_{\theta} J(\theta) = \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{Y}$$

$$= \mathbf{X}^T (\mathbf{X} \theta - \mathbf{Y})$$

$$\mathbf{X} = \begin{bmatrix} \cdots & \mathbf{x}_1^T & \cdots \\ \cdots & \mathbf{x}_2^T & \cdots \\ \vdots & \vdots & \vdots \\ \cdots & \mathbf{x}_n^T & \cdots \end{bmatrix} = \mathbf{X}^T \left(\begin{bmatrix} -\bar{x}_1 - \\ -\bar{x}_2 - \\ \vdots \\ -\bar{x}_n - \end{bmatrix} \theta - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \right).$$

$n \times p \quad p \times 1$

$$= \mathbf{X}^T \begin{bmatrix} \bar{x}_1^T \theta - y_1 \\ \bar{x}_2^T \theta - y_2 \\ \dots \\ \bar{x}_n^T \theta - y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & 1 & 1 \\ \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_n \end{bmatrix} \begin{bmatrix} \bar{x}_1^T \theta - y_1 \\ \bar{x}_2^T \theta - y_2 \\ \dots \\ \bar{x}_n^T \theta - y_n \end{bmatrix}_{n \times 1}.$$

$$= \sum_{i=1}^n \bar{x}_i \cdot (\bar{x}_i^T \theta - y_i)$$

1*1

Illustration of Gradient Descent (2D case)

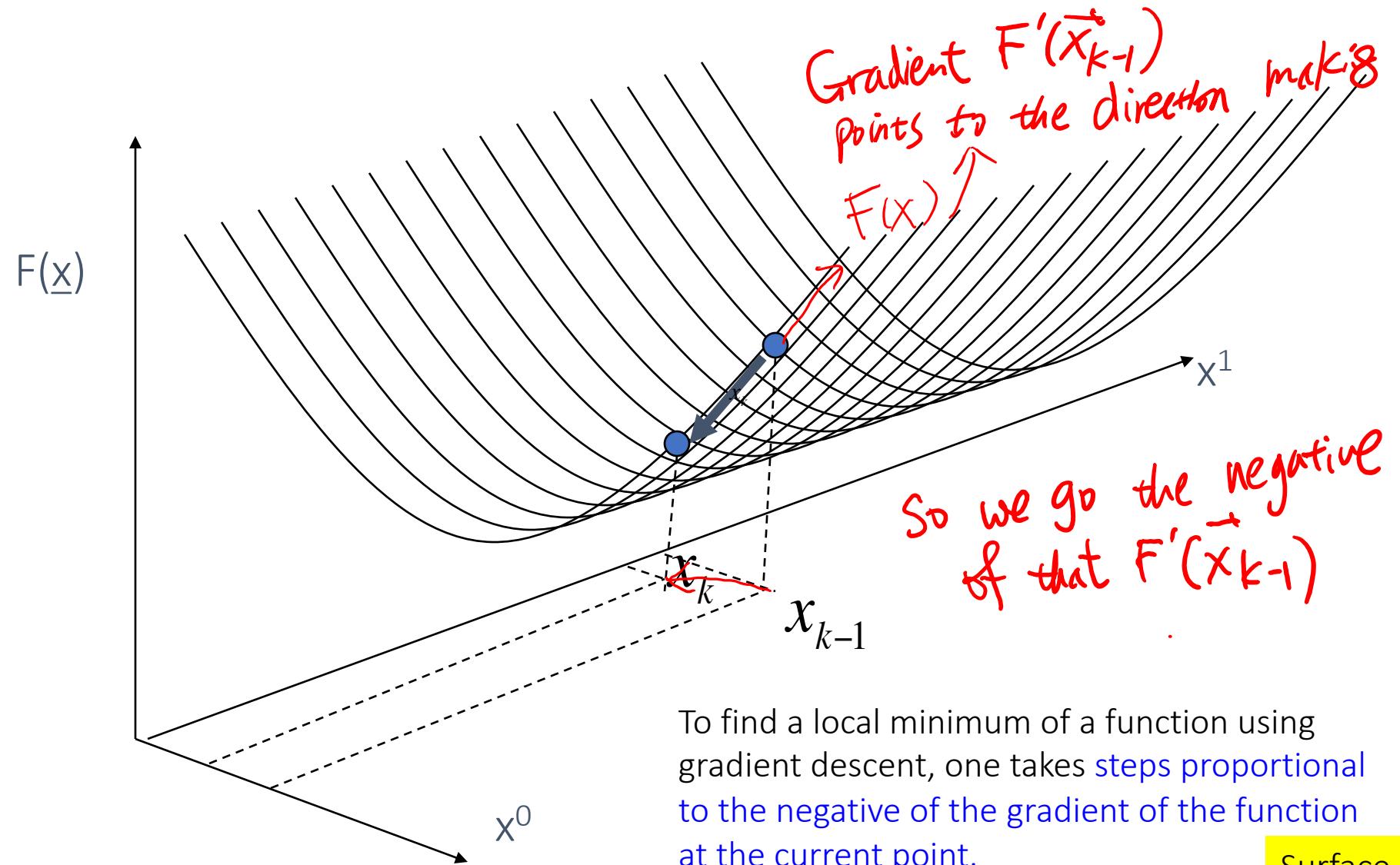
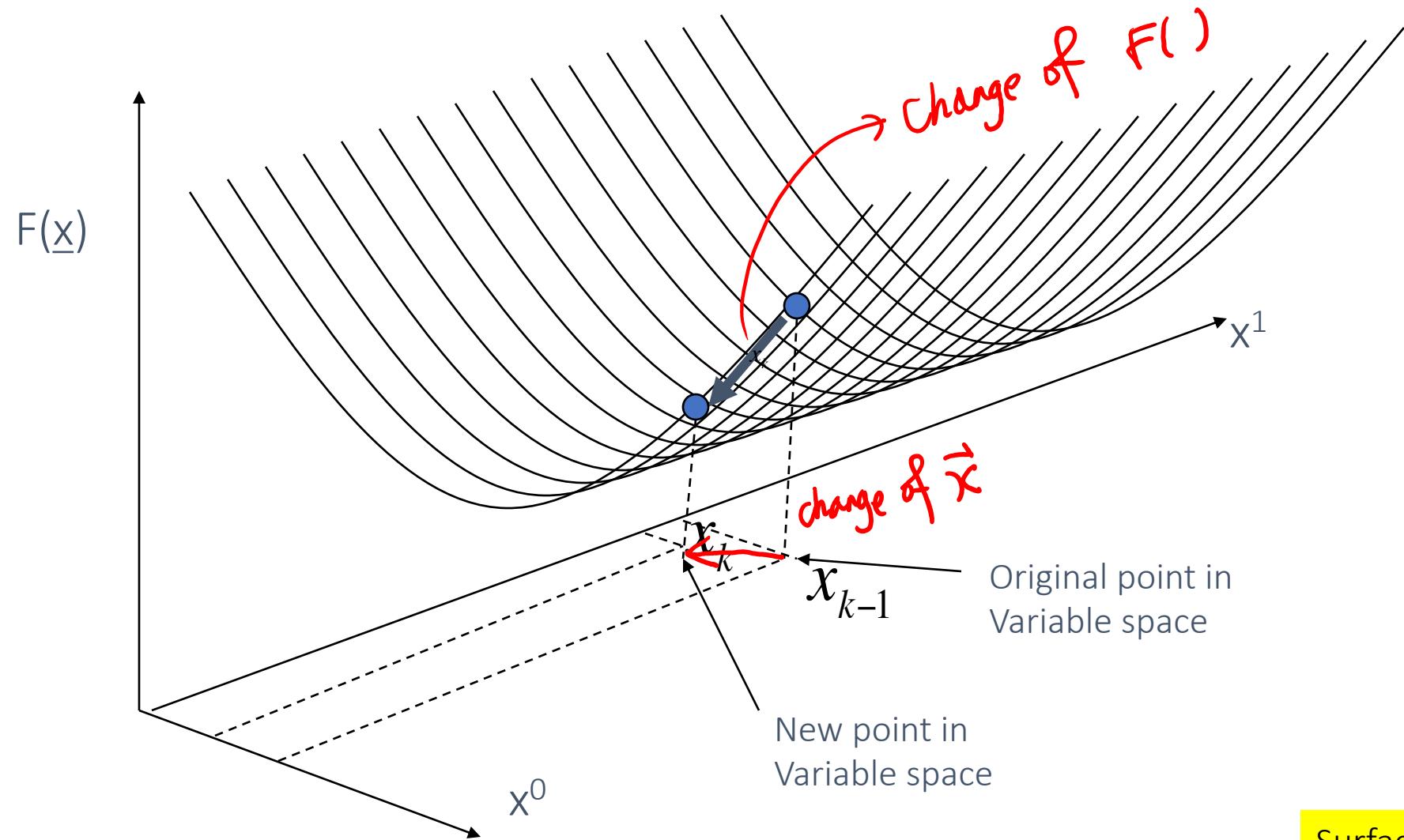


Illustration of Gradient Descent (2D case)



LR with batch GD / Per Feature View

- Note that:

$$\nabla_{\theta} J = \left[\frac{\partial}{\partial \theta_1} J, \dots, \frac{\partial}{\partial \theta_k} J \right]^T$$

- For its j-th variable:

$$\theta^{t+1} = \theta^t + \alpha \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta^t) \mathbf{x}_i$$

$$\theta_j^{t+1} = \theta_j^t + \alpha \sum_{i=1}^n (y_i - \mathbf{x}_i^T \theta^t) x_{i,j}$$



LR with Stochastic GD / Per Feature View

- For a single training point (i-th), we have:
 - For its j-th variable:

$$\theta^{t+1} = \theta^t + \alpha(y_i - \bar{\mathbf{x}}_i^T \theta^t) \bar{\mathbf{x}}_i$$

$$\theta_j^{t+1} = \theta_j^t + \alpha(y_i - \bar{\mathbf{x}}_i^T \theta^t) x_{i,j}$$

Update Rule Per Feature
(Variable-Wise)

EXTRA II

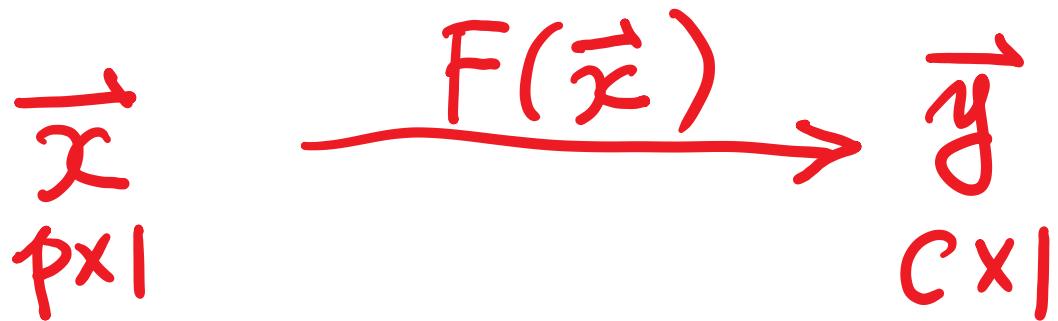
In Case you are interested in more advanced details!

Extra: Newton's Method and

Connecting to Normal Equation

Review: Single Var-Func to Multivariate

Single Var- Function	Multivariate Calculus
Derivative	Partial Derivative
Second-order derivative	Gradient
	Directional Partial Derivative
	Vector Field
	Contour map of a function
	Surface map of a function
	Hessian matrix
	Jacobian matrix (vector in / vector out)



$p=1$	$C=1$	derivative	2nd-order derivative
$p > 1$	$C=1$	gradient vector	Hessian matrix
$p > 1$	$C > 1$	Jacobian matrix	

Newton's method for optimization

- The most basic second-order optimization algorithm
- Updating parameter with

$$\text{GD: } \theta_{k+1} = \theta_k - \alpha g_k$$

$$\text{Newton: } \theta_{k+1} = \theta_k - H_K^{-1} g_k$$

$$\underbrace{\begin{matrix} P \times P & P \times 1 \\ \hline P \times 1 \end{matrix}}$$

Review: Hessian Matrix / n==2 case

Singlevariate

→ multivariate

$$f(x, y)$$

- 1st derivative to gradient,
- 2nd derivative to Hessian

$$g = \nabla f = \begin{pmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{pmatrix}$$

$$H = \begin{pmatrix} \frac{\partial^2 f}{\partial x^2} & \frac{\partial^2 f}{\partial x \partial y} \\ \frac{\partial^2 f}{\partial x \partial y} & \frac{\partial^2 f}{\partial y^2} \end{pmatrix}$$

Review: Hessian Matrix

Suppose that $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is a function that takes a vector in \mathbb{R}^n and returns a real number. Then the **Hessian** matrix with respect to x , written $\nabla_x^2 f(x)$ or simply as H is the $n \times n$ matrix of partial derivatives,

$$\nabla_x^2 f(x) \in \mathbb{R}^{n \times n} = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \dots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix}.$$

Newton's method for optimization

- Making a quadratic/second-order Taylor series approximation

$$\hat{f}_{quad}(\boldsymbol{\theta}) = f(\boldsymbol{\theta}_k) + \mathbf{g}_k^T (\boldsymbol{\theta} - \boldsymbol{\theta}_k) + \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_k)^T \mathbf{H}_k (\boldsymbol{\theta} - \boldsymbol{\theta}_k)$$

Finding the minimum solution of the above right quadratic approximation (quadratic function minimization is easy !)

$$\hat{f}(\theta) = f(\theta_K) + g_K^\top (\theta - \theta_K) +$$

$$\frac{1}{2} (\theta - \theta_K)^\top H_K (\theta - \theta_K)$$
$$\frac{1}{2} (\theta^\top H_K \theta - 2\theta^\top H_K \theta_K + \theta_K^\top H_K \theta_K)$$

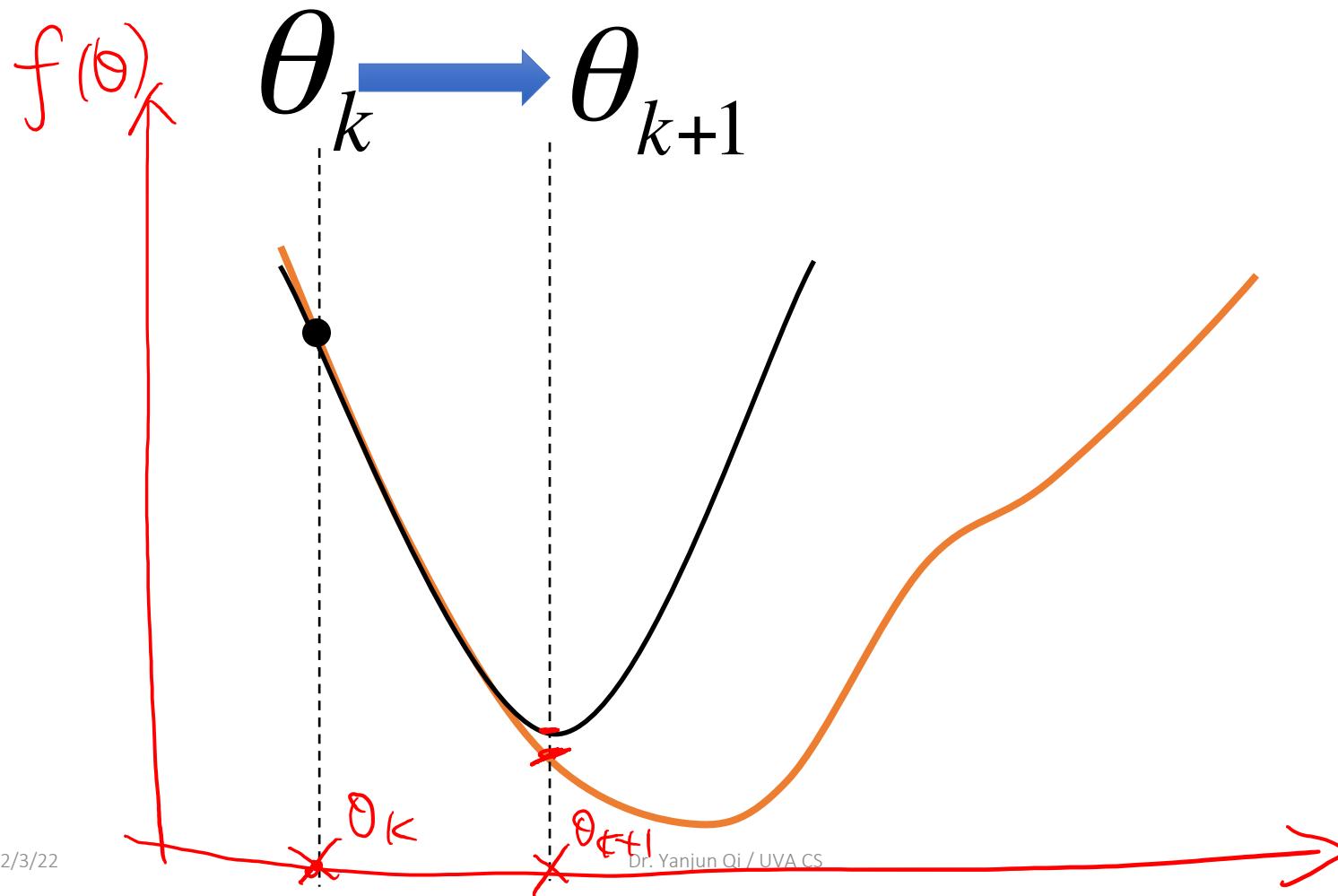
$$\frac{\partial \hat{f}(\theta)}{\partial \theta} = 0 + g_K + \underbrace{\frac{2}{2} H_K \theta - \frac{2}{2} H_K \theta_K}_{\text{See P24 handout}} := 0$$

$$g_K + H_K (\theta - \theta_K) = 0$$

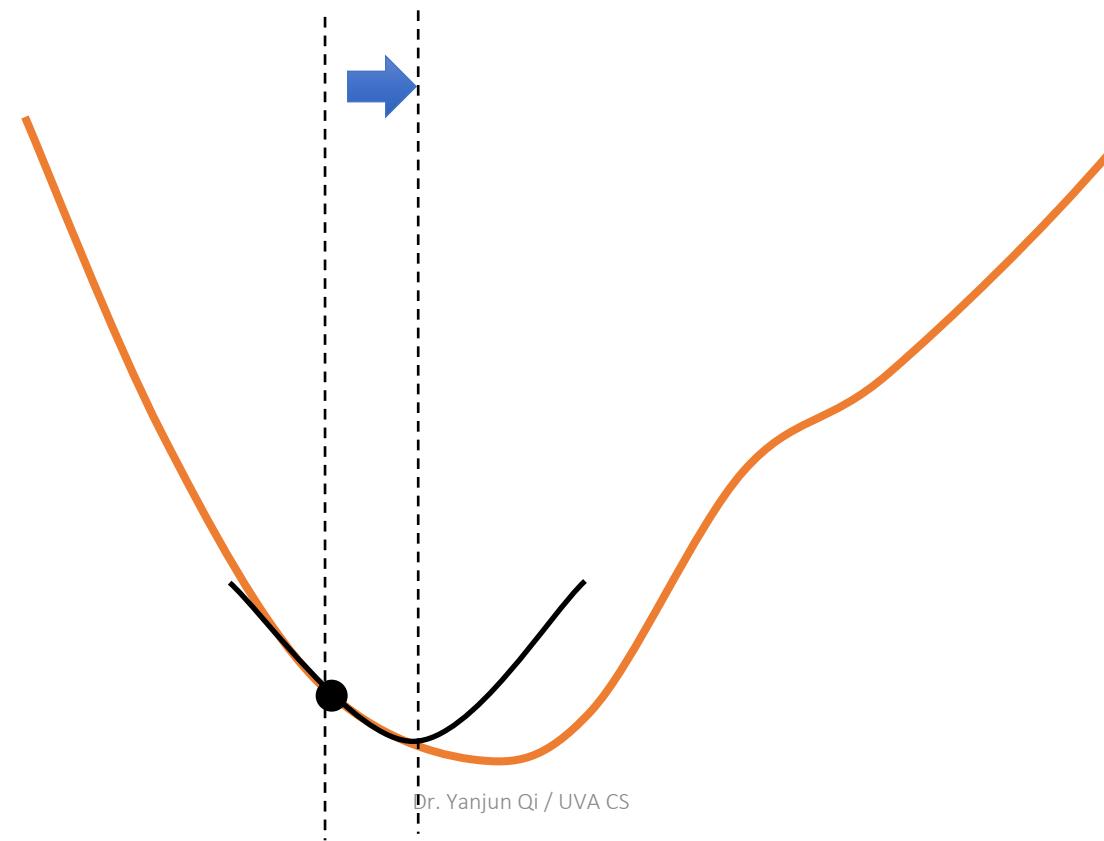
$$\Rightarrow \theta = \theta_K - H_K^{-1} g_K$$

where $\begin{cases} H_K \in \mathbb{R}^{P \times P} \\ g_K \in \mathbb{R}^P \end{cases}$
 80

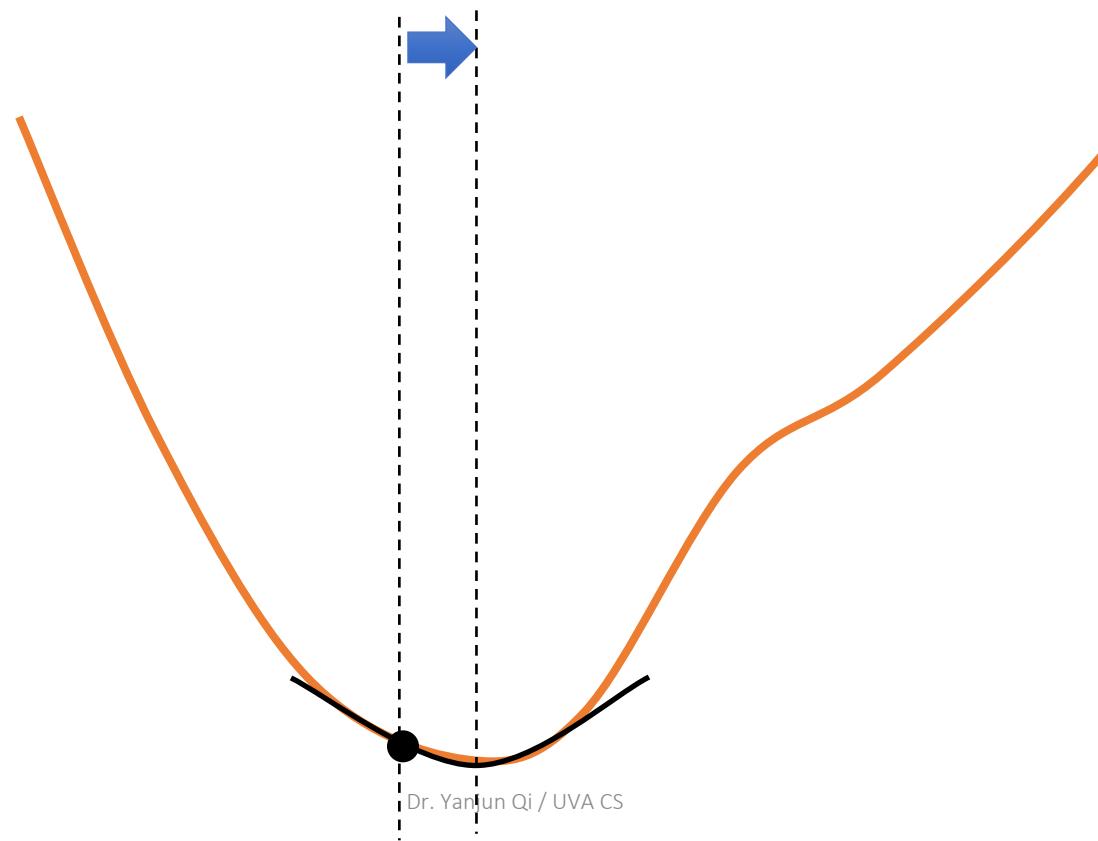
Newton's Method / second-order Taylor series approximation



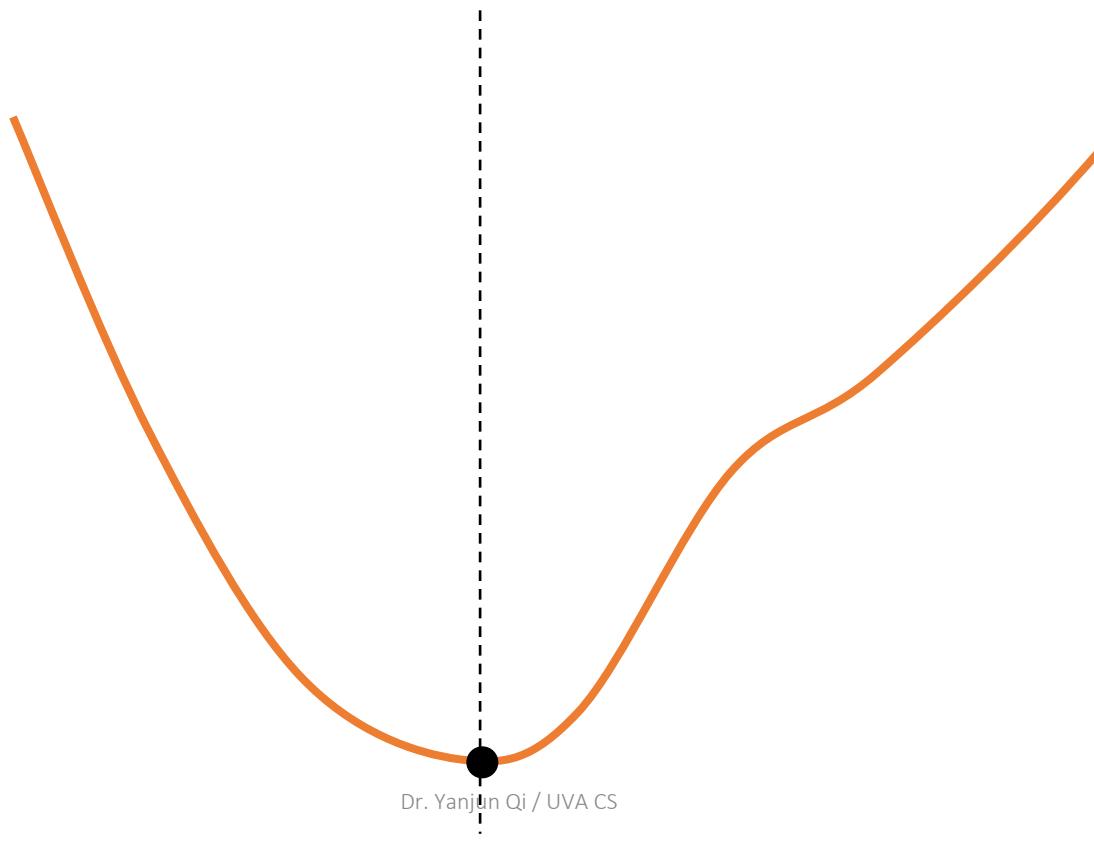
Newton's Method / second-order Taylor series approximation



Newton's Method / second-order Taylor series approximation



Newton's Method / second-order Taylor series approximation



Newton's Method

- At each step:

$$\theta_{k+1} = \theta_k - \frac{f'(\theta_k)}{f''(\theta_k)}$$

$$\theta_{k+1} = \theta_k - H^{-1}(\theta_k) \nabla f(\theta_k)$$

- Requires 1st and 2nd derivatives
- Quadratic convergence
- → However, finding the inverse of the Hessian matrix is often expensive

Newton vs. GD for optimization

- Newton: a quadratic/second-order Taylor series approximation

$$\Theta_{k+1} = \Theta_k - \frac{1}{H(\Theta_k)} g(\Theta_k)$$

$$\widehat{f}_{quad}(\theta) = f(\theta_k) + \mathbf{g}_k^T (\theta - \theta_k) + \frac{1}{2} (\theta - \theta_k)^T \mathbf{H}_k (\theta - \theta_k)$$

Finding the minimum solution of
the above right quadratic
approximation (quadratic
function minimization is easy !)

$$\widehat{f}_{quad}(\theta) = f(\theta_k) + \mathbf{g}_k^T (\theta - \theta_k) + \frac{1}{2} (\theta - \theta_k)^T \frac{1}{\alpha} (\theta - \theta_k)$$

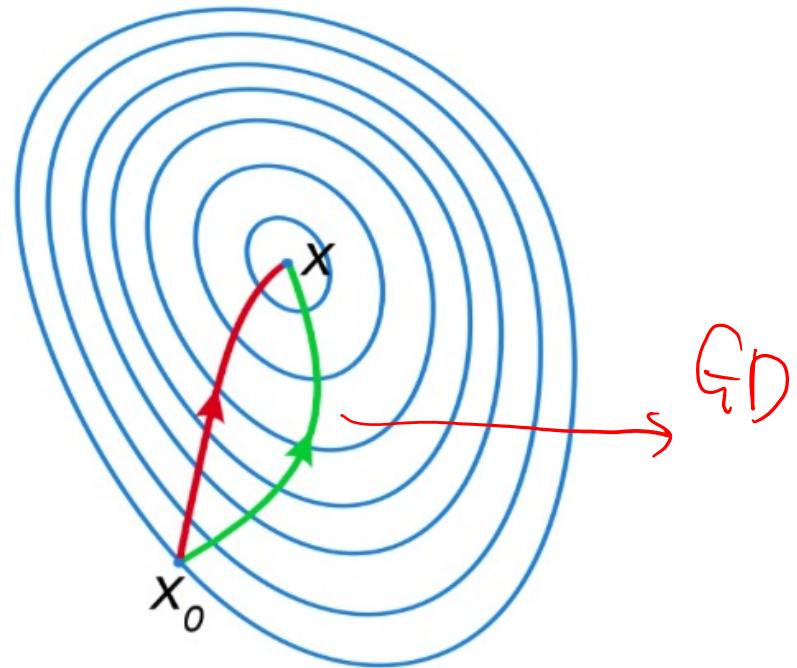
$$\Downarrow \quad \Theta_{k+1} = \Theta_k - \frac{1}{\alpha} g(\Theta_k)$$

Comparison

- Newton's method vs. Gradient descent

A comparison of gradient descent (green) and Newton's method (red) for minimizing a function (with small step sizes).

Newton's method uses curvature information to get a more direct route ...



$$J(\theta) = \frac{1}{2} (\mathbf{y} - \mathbf{X}\theta)^T (\mathbf{y} - \mathbf{X}\theta)$$

$$\nabla_{\theta} J(\theta) = \mathbf{X}^T \mathbf{X} \theta - \mathbf{X}^T \mathbf{y}$$

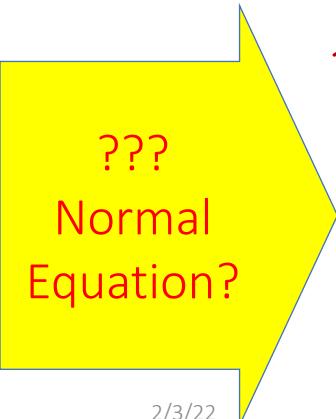
$$H = \nabla_{\theta}^2 J(\theta) = \mathbf{X}^T \mathbf{X}$$

$$\Rightarrow \theta^t = \theta^{t-1} - H^{-1} \nabla J(\theta^{t-1}) \quad \text{Newton}$$

$$= \theta^{t-1} - (\mathbf{X}^T \mathbf{X})^{-1} [\mathbf{X}^T \mathbf{X} \theta^{t-1} - \mathbf{X}^T \mathbf{y}]$$

$$= [\theta^{t-1} - \theta^{t-1}] + (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

$$= (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$



References

- Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- http://en.wikipedia.org/wiki/Matrix_calculus
- Prof. Nando de Freitas's tutorial slide
- An overview of gradient descent optimization algorithms,
<https://arxiv.org/abs/1609.04747>