

UVA CS 6316: Machine Learning

Lecture 4: More optimization for Linear Regression

Dr. Yanjun Qi

University of Virginia
Department of Computer Science

Course Content Plan →

Six major sections of this course

- Regression (supervised)
- Classification (supervised)
- Unsupervised models
- Learning theory
- Graphical models
- Reinforcement Learning

Y is a continuous

Y is a discrete

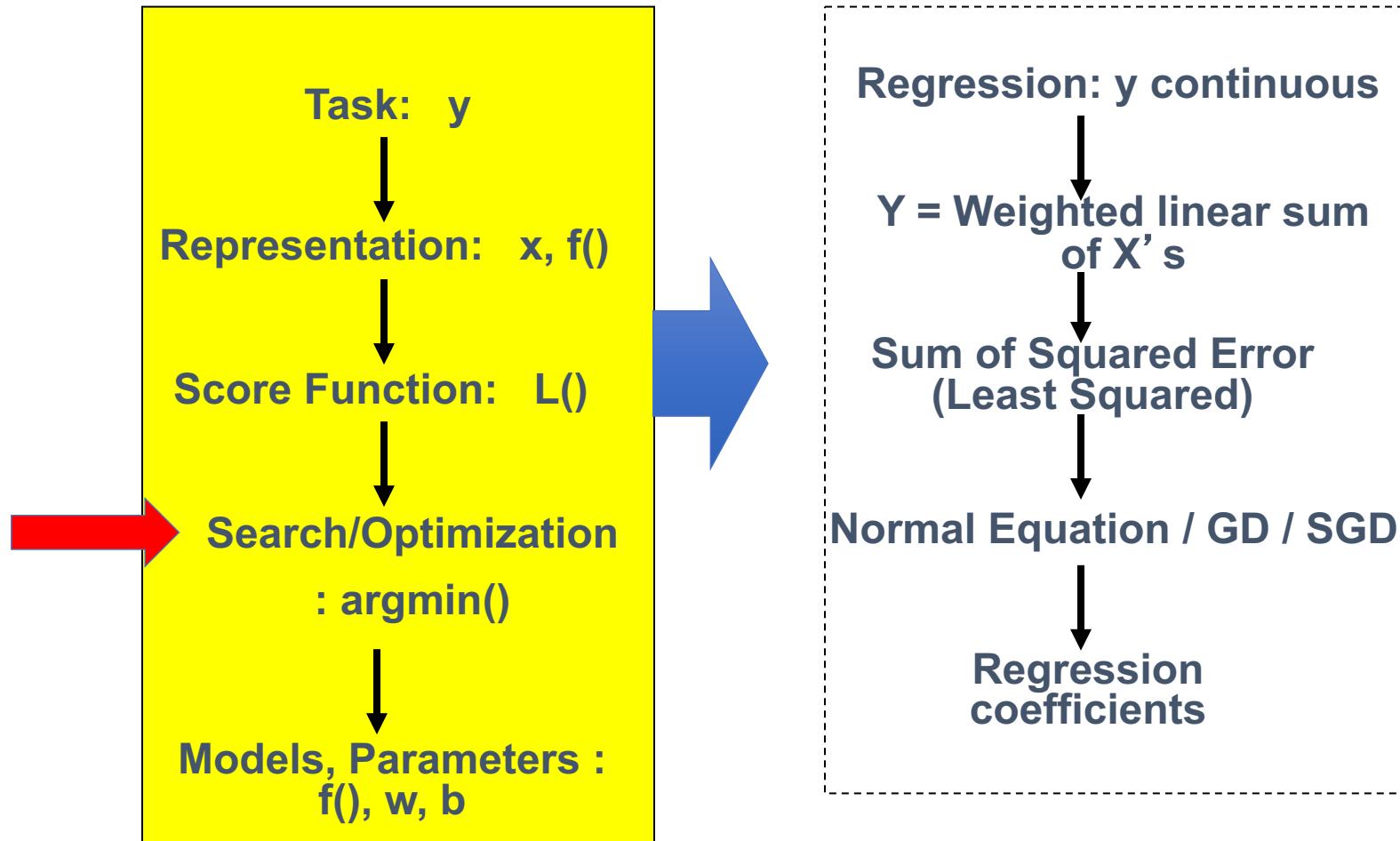
NO Y

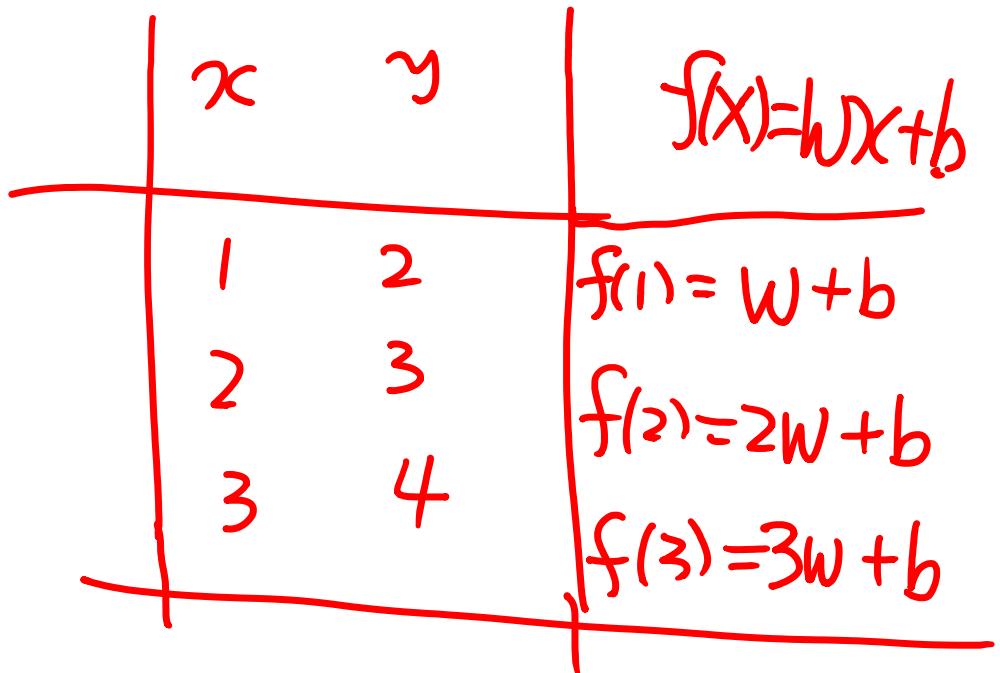
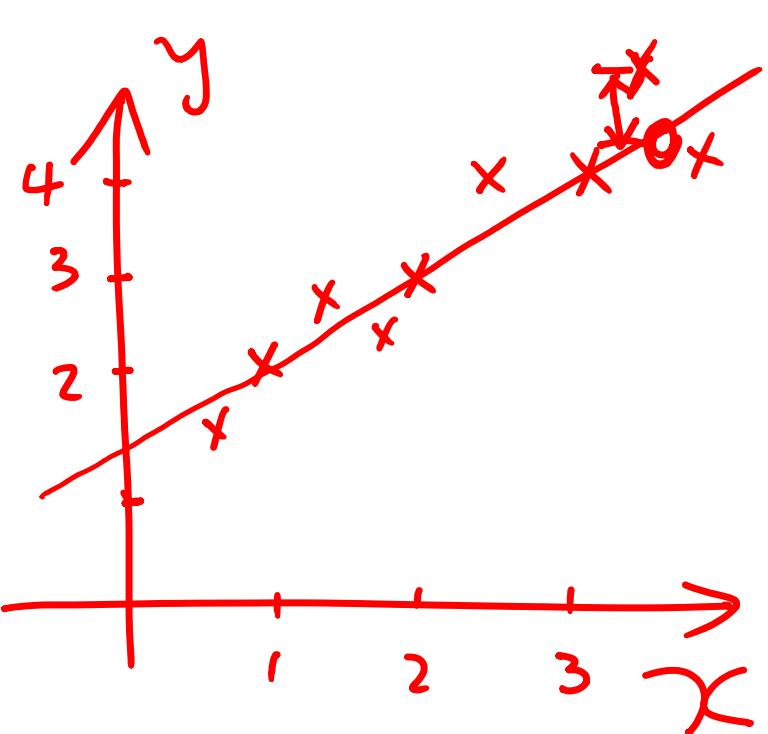
About $f()$

About interactions among X_1, \dots, X_p

Learn program to Interact with its environment

Today: Multivariate Linear Regression in a Nutshell





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

$$\text{loss}(w, b) = \frac{1}{2} J(w, b)$$

$$J(w, b) = J(\theta_1, \theta_0)$$

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

Method I: normal equations

- Write the cost function in matrix form:

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

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

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

$$\Rightarrow \boxed{\mathbf{X}^T \mathbf{X}\theta = \mathbf{X}^T \bar{\mathbf{y}}}$$

The normal equations

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

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

Learning Regression Models (supervised)

- ❑ Four ways to train / perform optimization for learning linear regression models
 - ❑ ~~Normal Equation~~
 - ❑ Gradient Descent (GD)
 - ❑ Stochastic GD / Mini-Batch
 - ❑ Connecting to Newton's method

A little bit more about [Optimization]

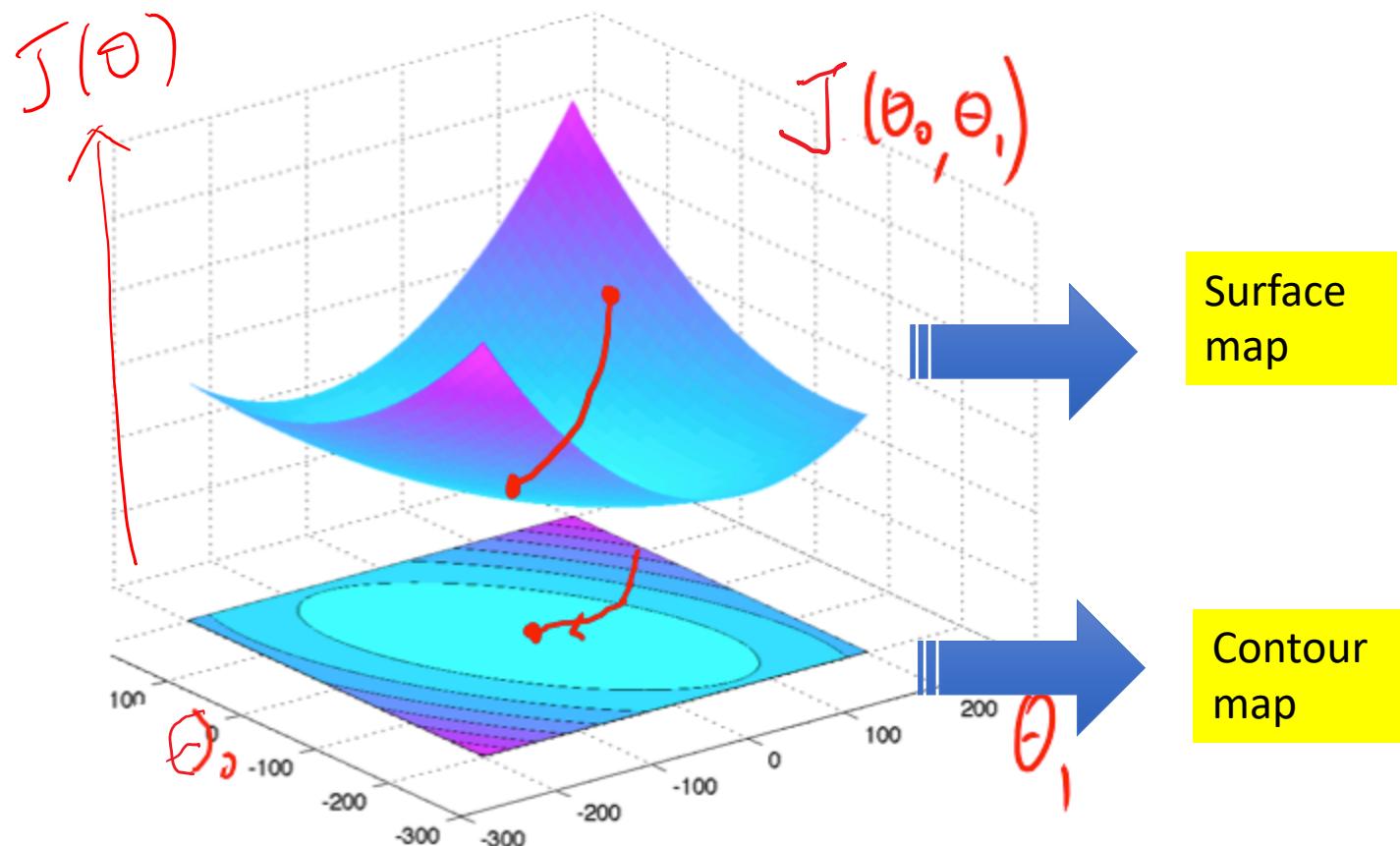
- Objective function $F(x) \rightarrow J(\theta)$
- Variables $x \rightarrow \theta$
- Constraints $\rightarrow \theta \in \mathbb{R}^P$

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

Today

- ❑ More ways to train / perform optimization for linear regression models
- ❑ Review: Gradient Descent
 - ❑ Gradient Descent (GD) for LR
 - ❑ Stochastic GD (SGD) for LR

Review: two ways of illustrating an Objective Function (e.g. 2D case)



Gradient vector points to the direction of greatest rate of increase of the objective function and its magnitude is the slope of the surface graph in that direction..

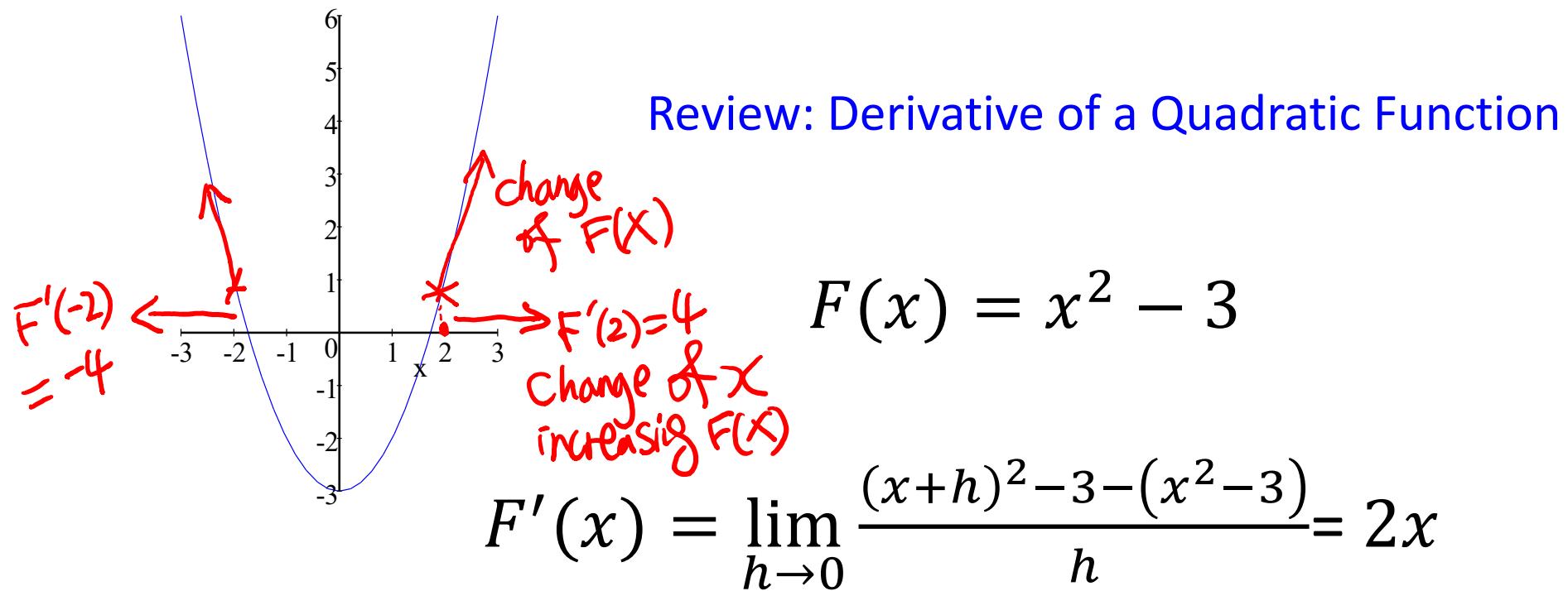
Review: Definitions of gradient (in L2-note)

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

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

A vector whose entries respectively contain the p partial derivatives

Review: Definitions of derivative (1D case)



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

Gradient Descent (GD): An iterative Algorithm

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

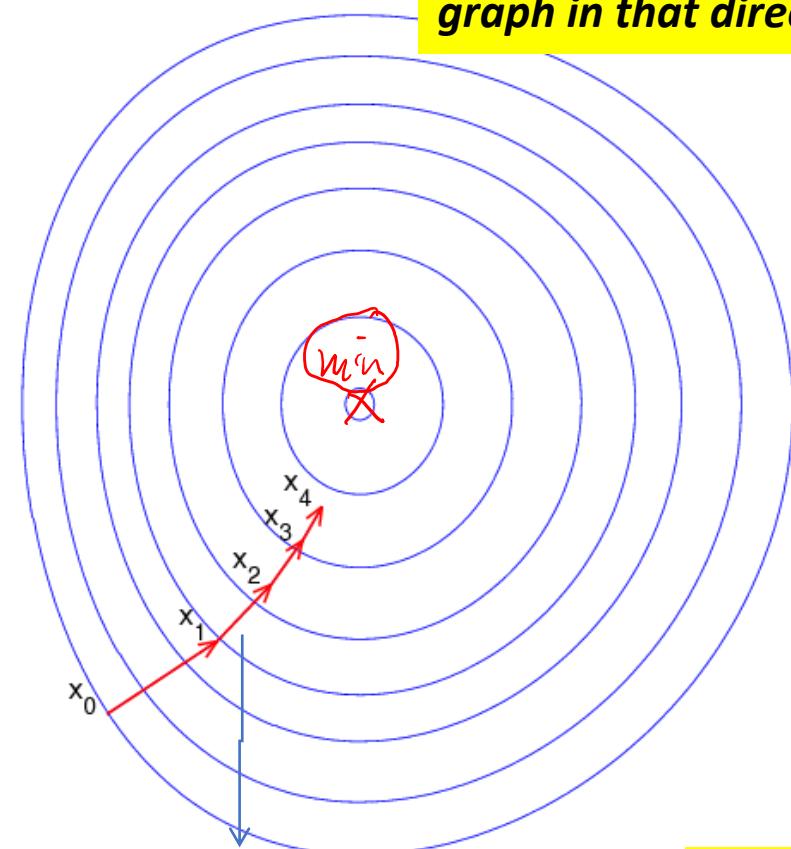
$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$

Gradient Descent (Steepest Descent) – contour map view

A first-order optimization algorithm.

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.

The gradient (in the variable space) points in the direction of the greatest rate of increase of the function and its magnitude is the slope of the surface graph in that direction



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

Contour map view

Illustration of Gradient Descent (2D case)

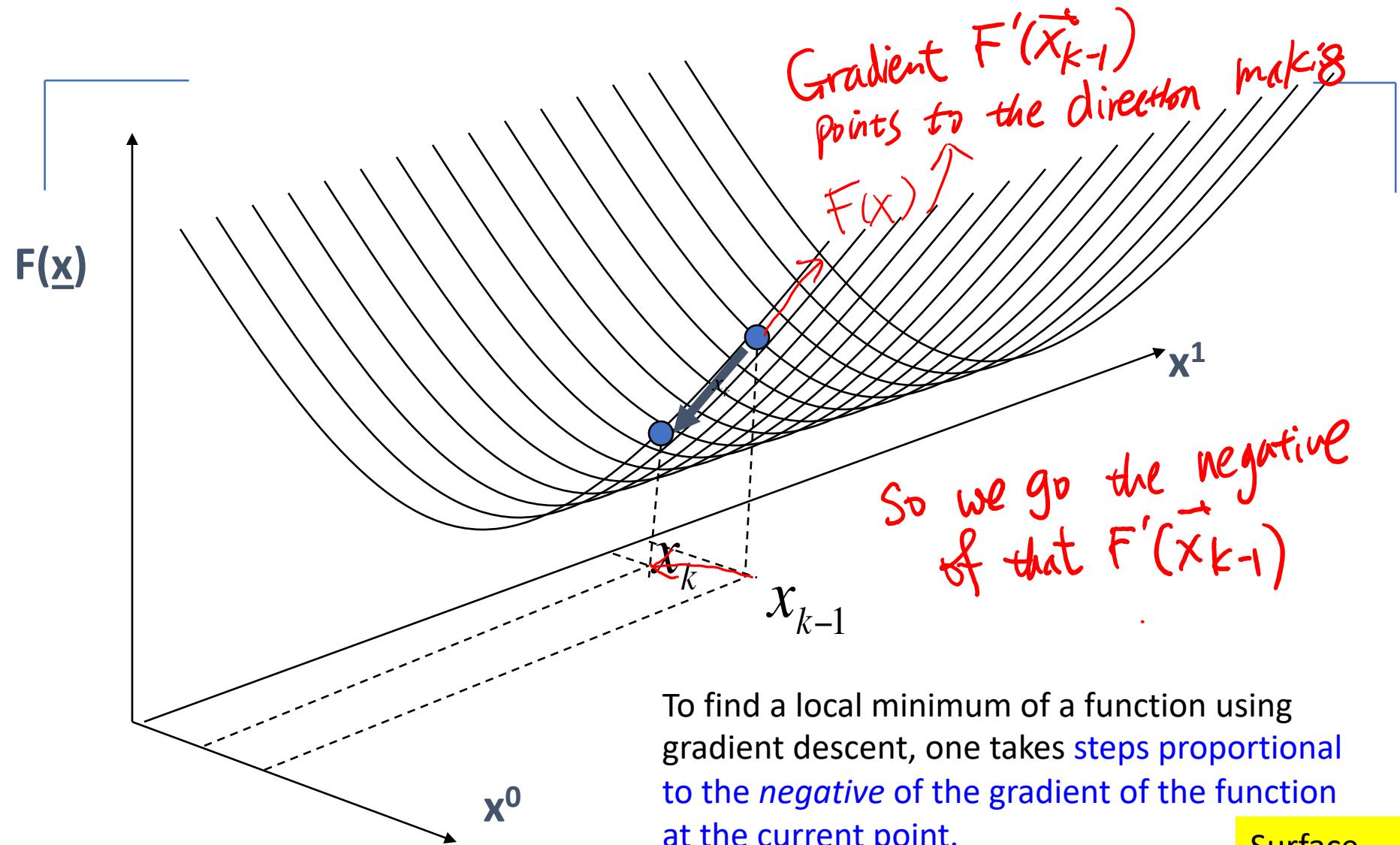
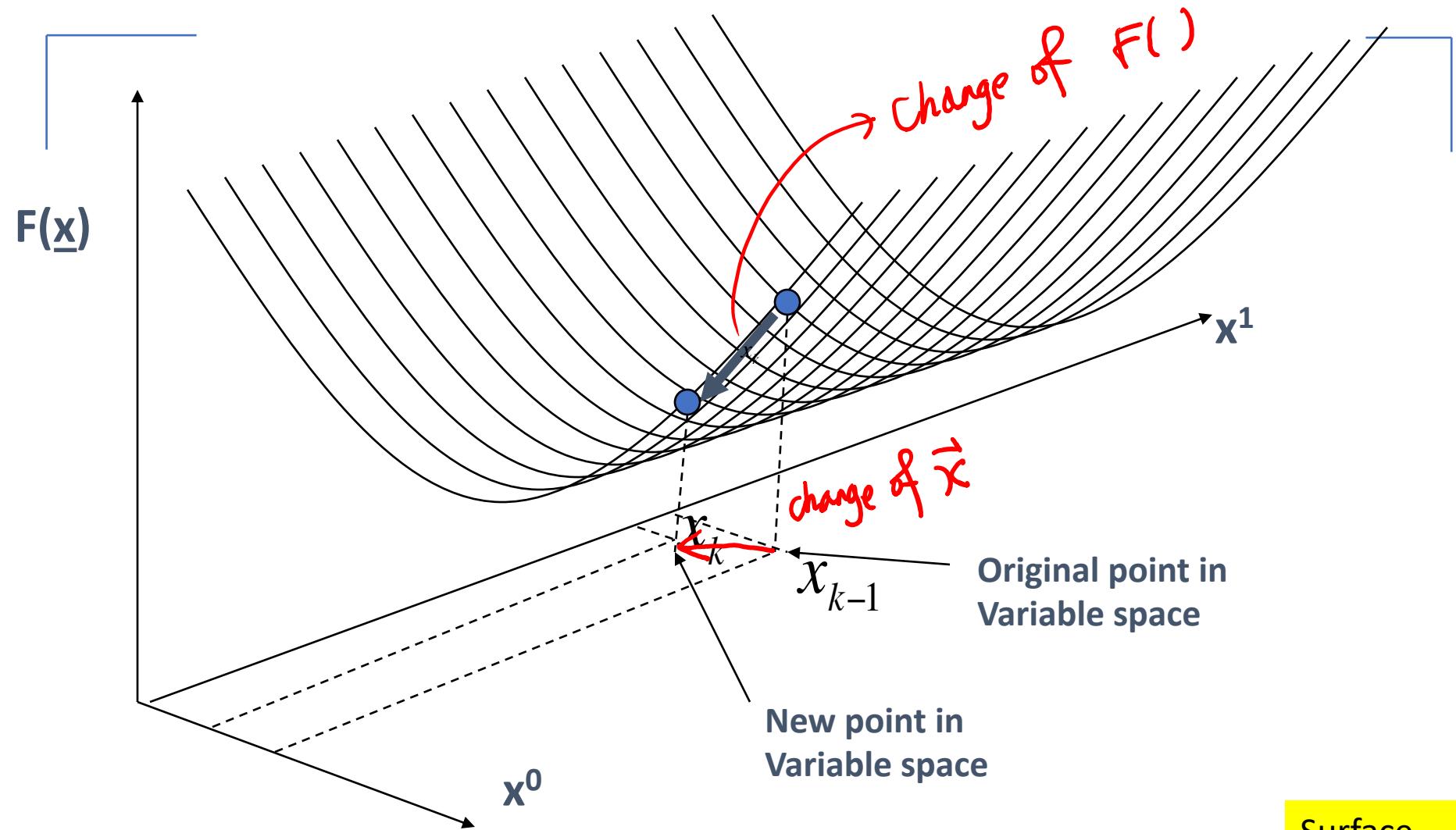
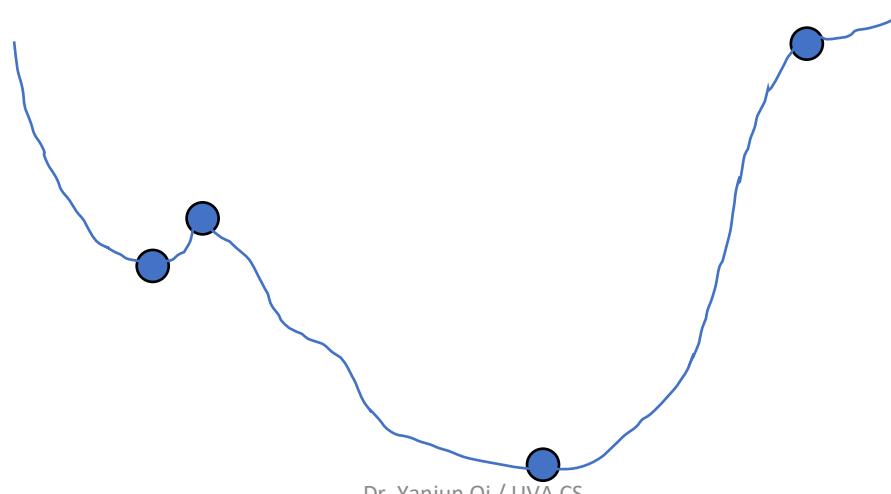


Illustration of Gradient Descent (2D case)

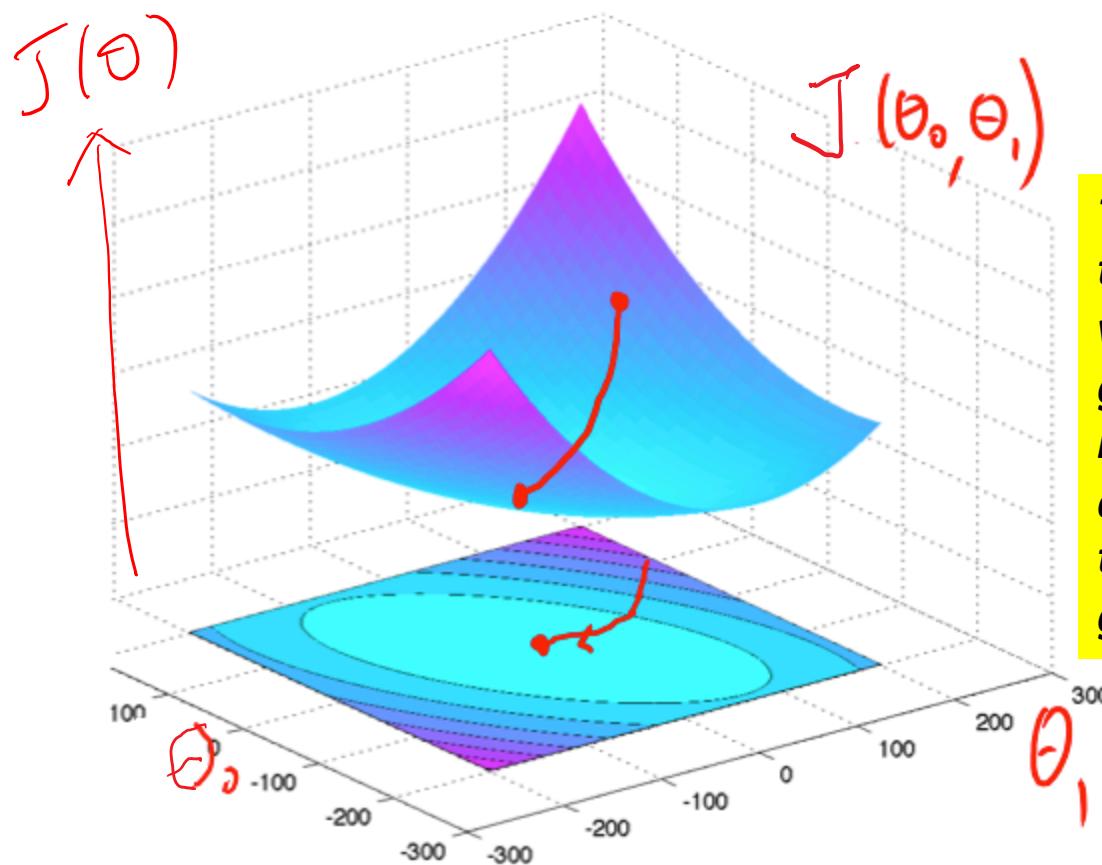


WHY ? Optimize through Gradient Descent (iterative) Algorithms

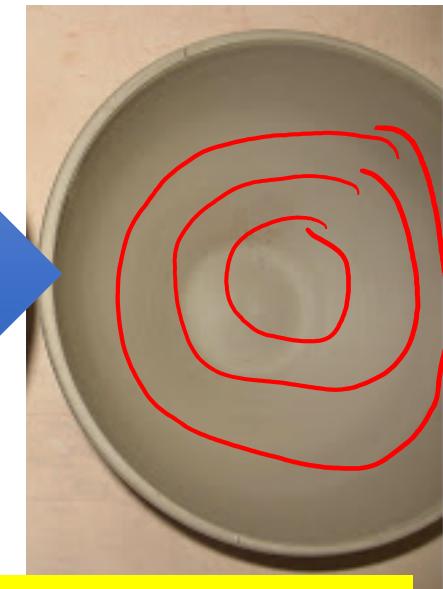
- Works on any objective function $F(x)$
 - as long as we can evaluate the gradient
 - this can be very useful for minimizing complex functions



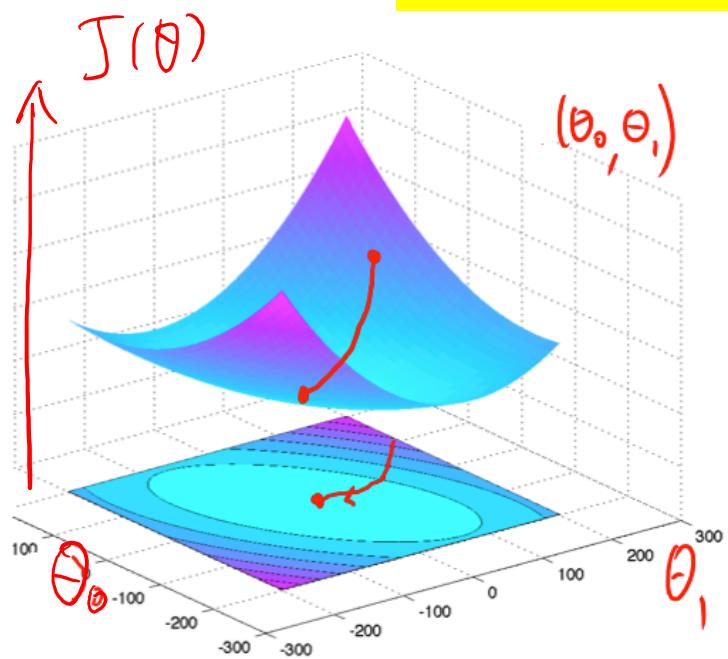
Two ways of Illustrating the Objective Function and Gradient Descent (e.g. , 2D case)



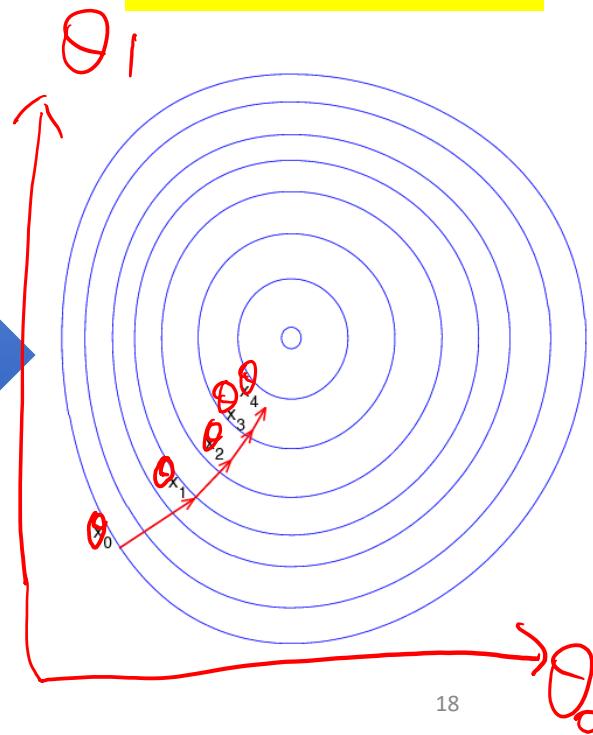
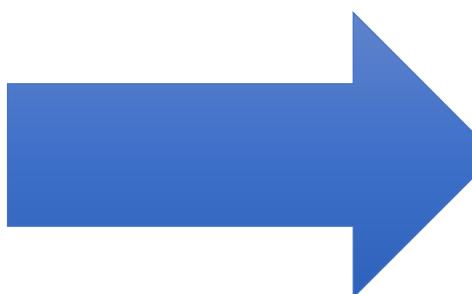
The gradient points in the direction (in the variable space) of the greatest rate of increase of the function and its magnitude is the slope of the surface graph in that direction

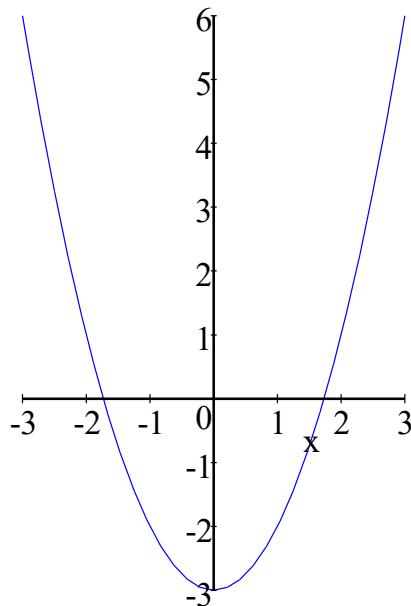


Surface map view



Contour map view





Review: Derivative of a Quadratic Function

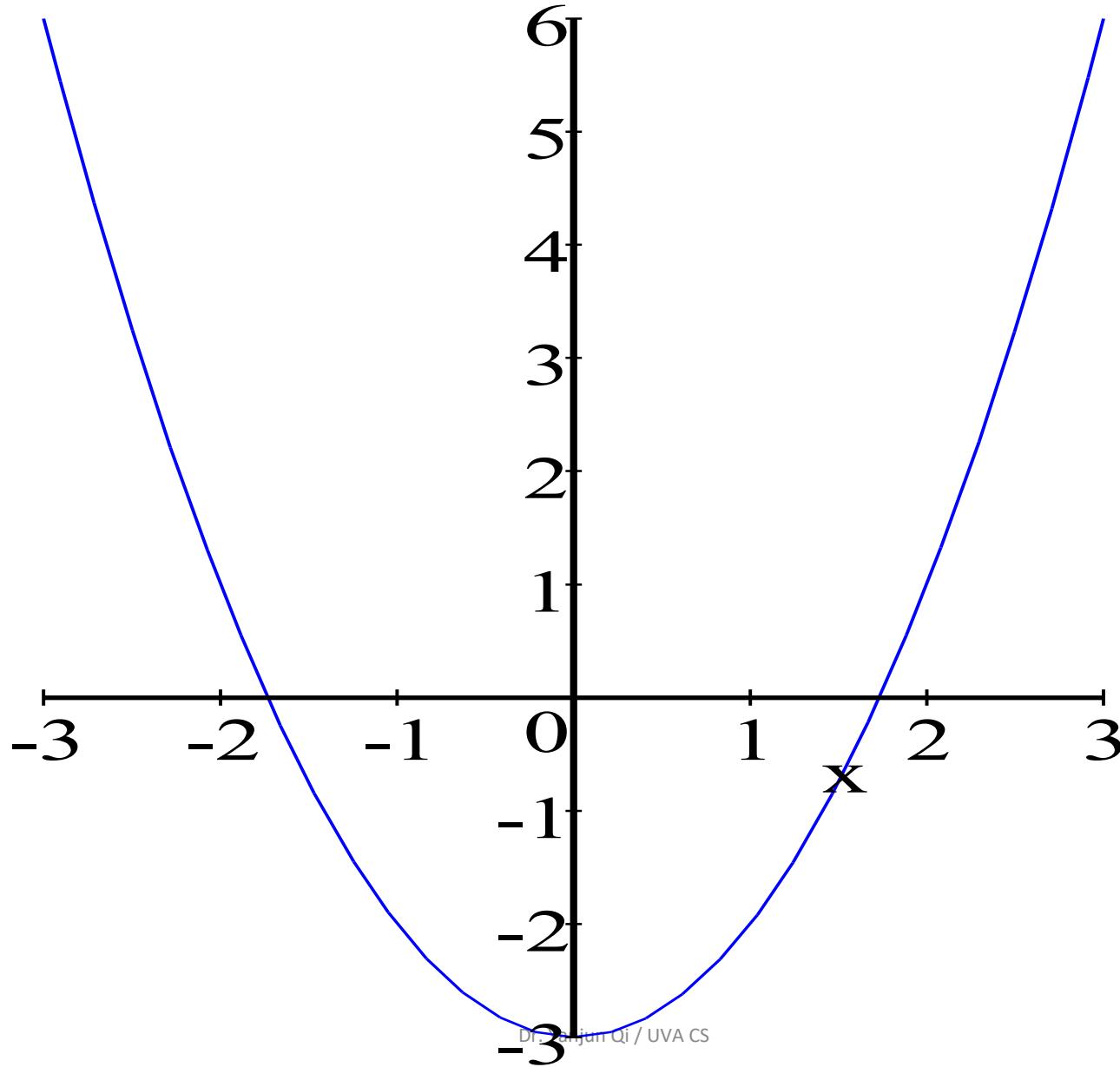
$$F(x) = x^2 - 3$$

$$\nabla_x F(x) = F'(x) = 2x$$

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

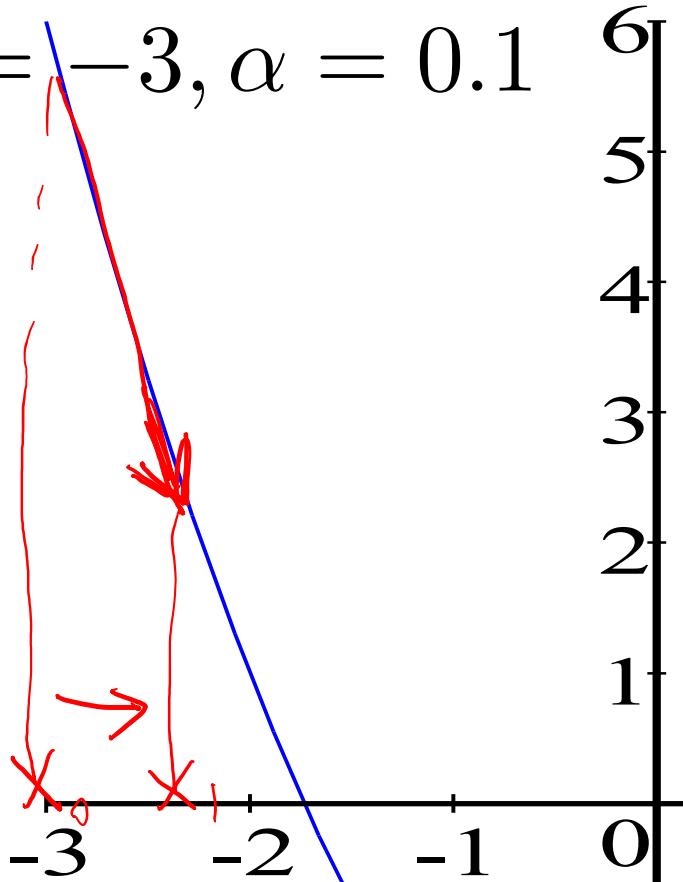
$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$

$$\nabla_x F(x) = F'(x) = 2x$$



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

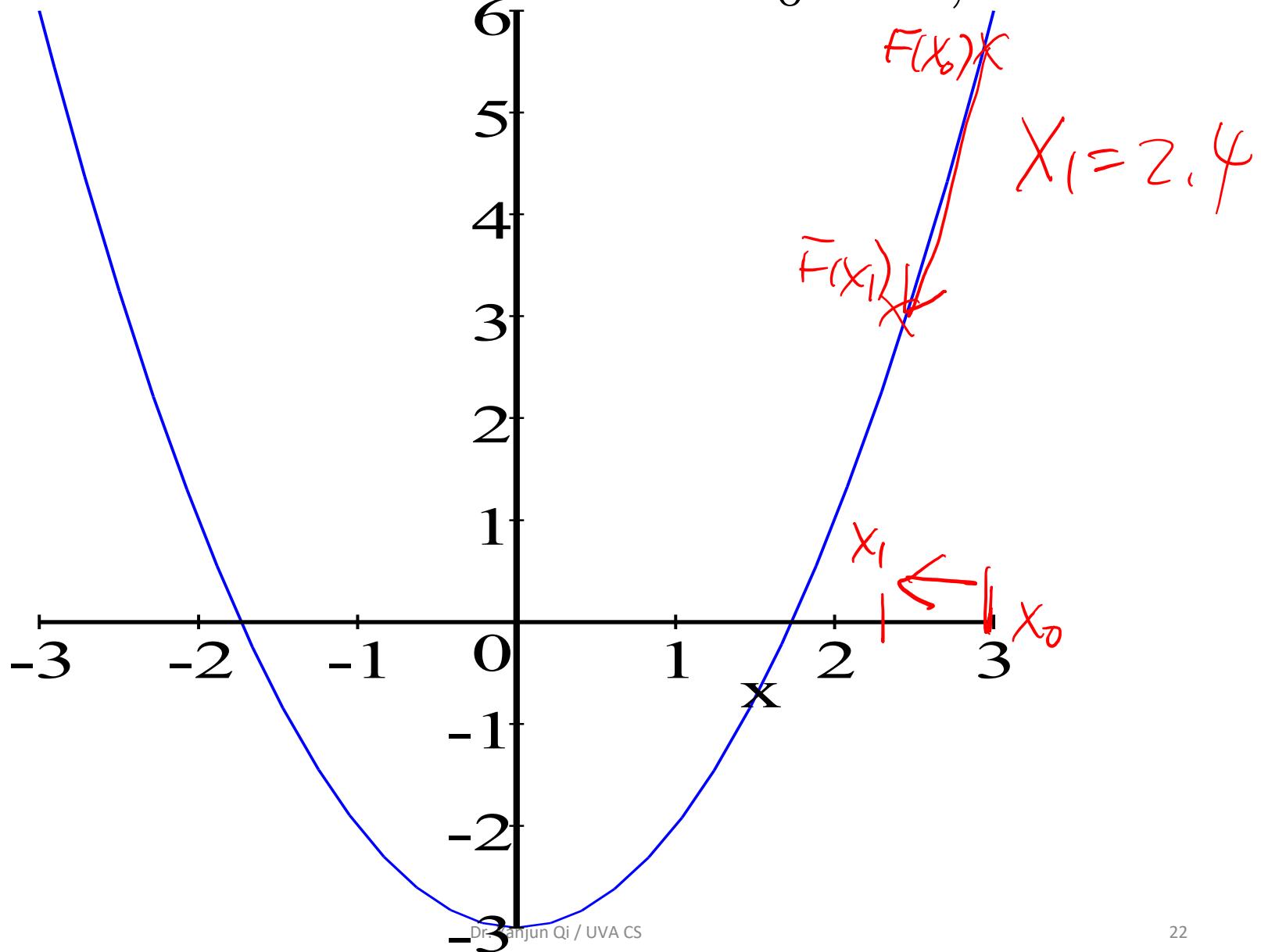
$$x_0 = -3, \alpha = 0.1$$



$$\begin{aligned} x_1 &= x_0 - \alpha \nabla F(x_0) \\ &= -3 \\ &- 2 \times 0.1 \times (-3) \\ &= -2.4 \end{aligned}$$

$$\nabla_x F(x) = F'(x) = 2x$$

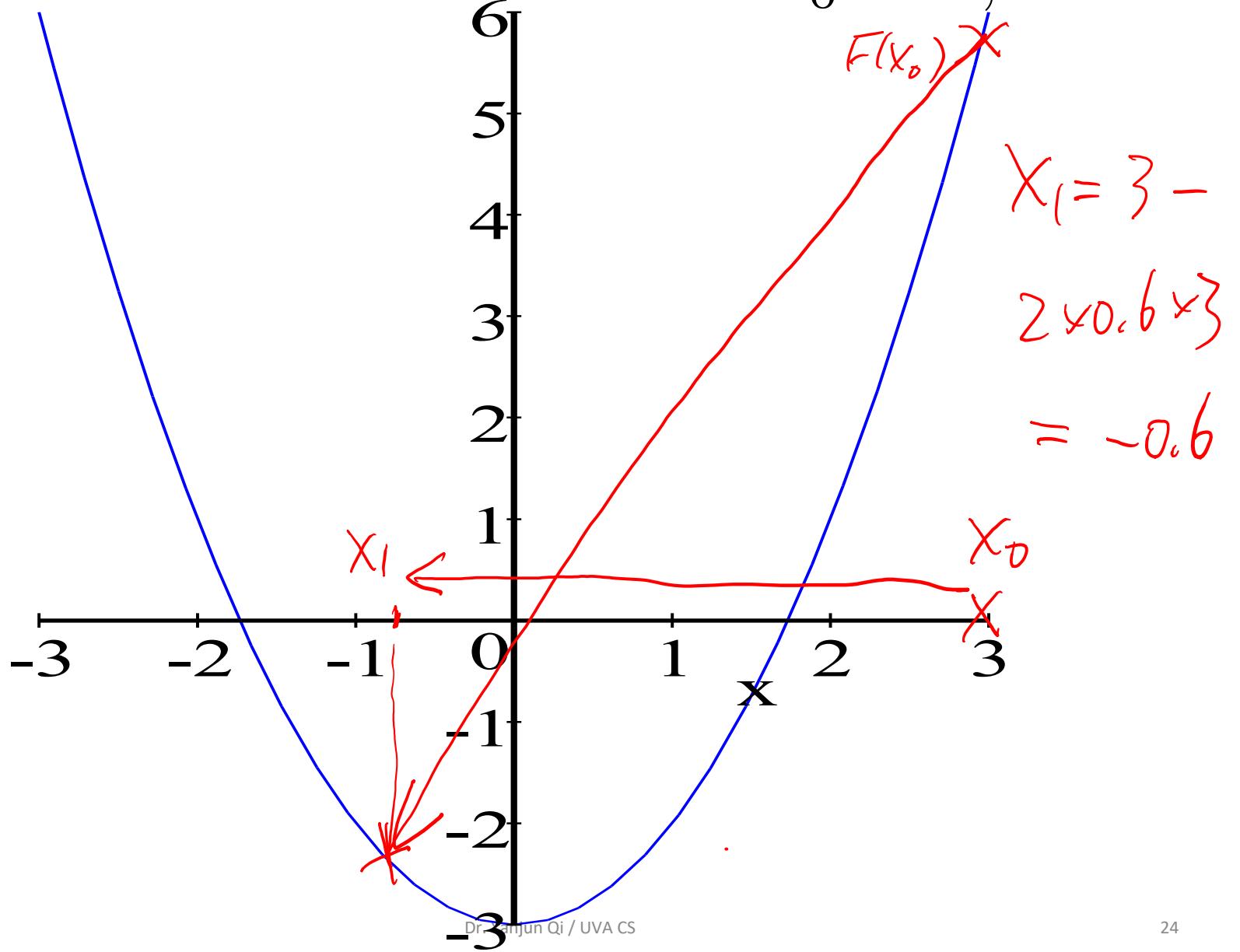
$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$



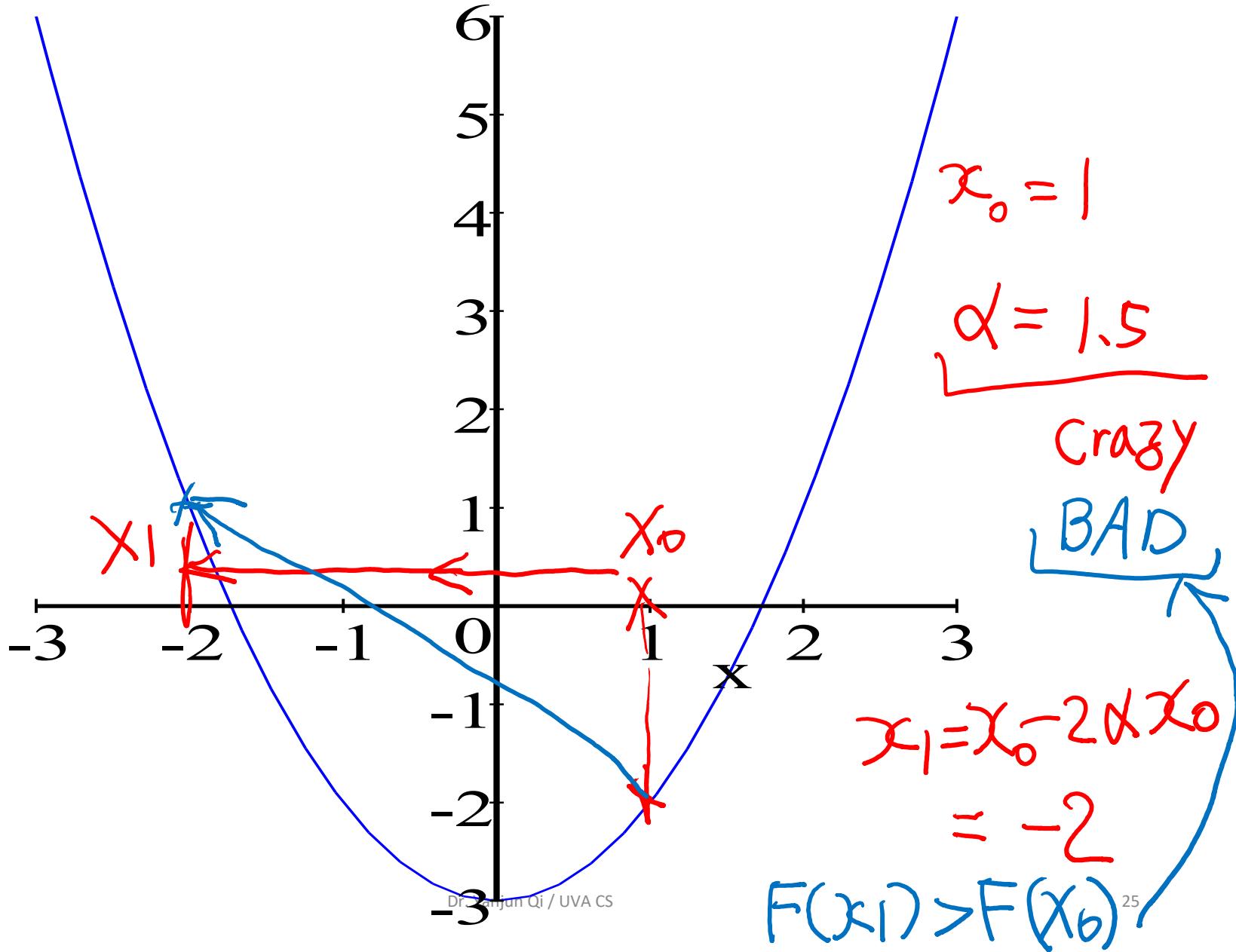
Gradient Descent (Iteratively Optimize)

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

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



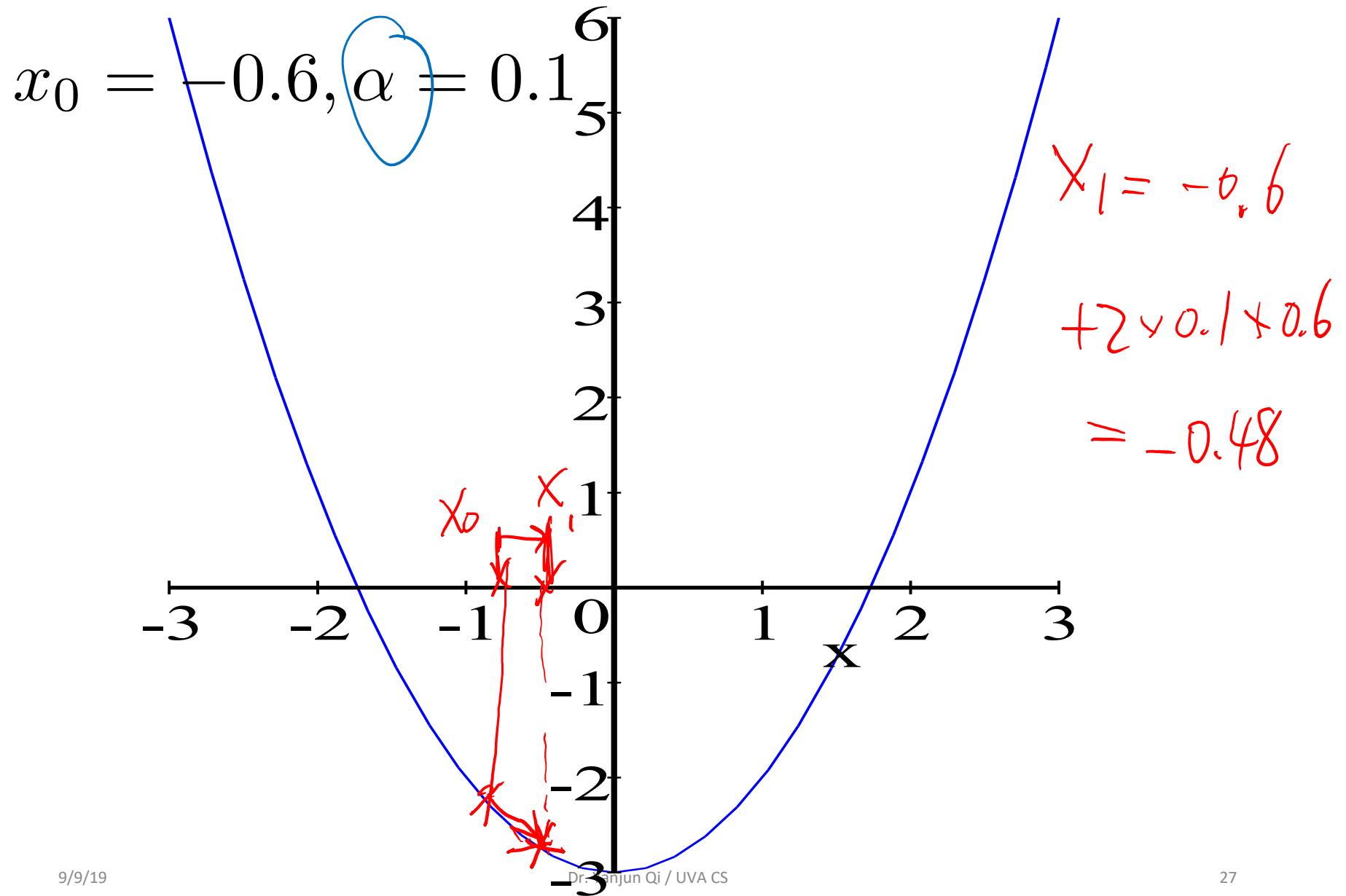
$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$



Gradient Descent (Iteratively Optimize)

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

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



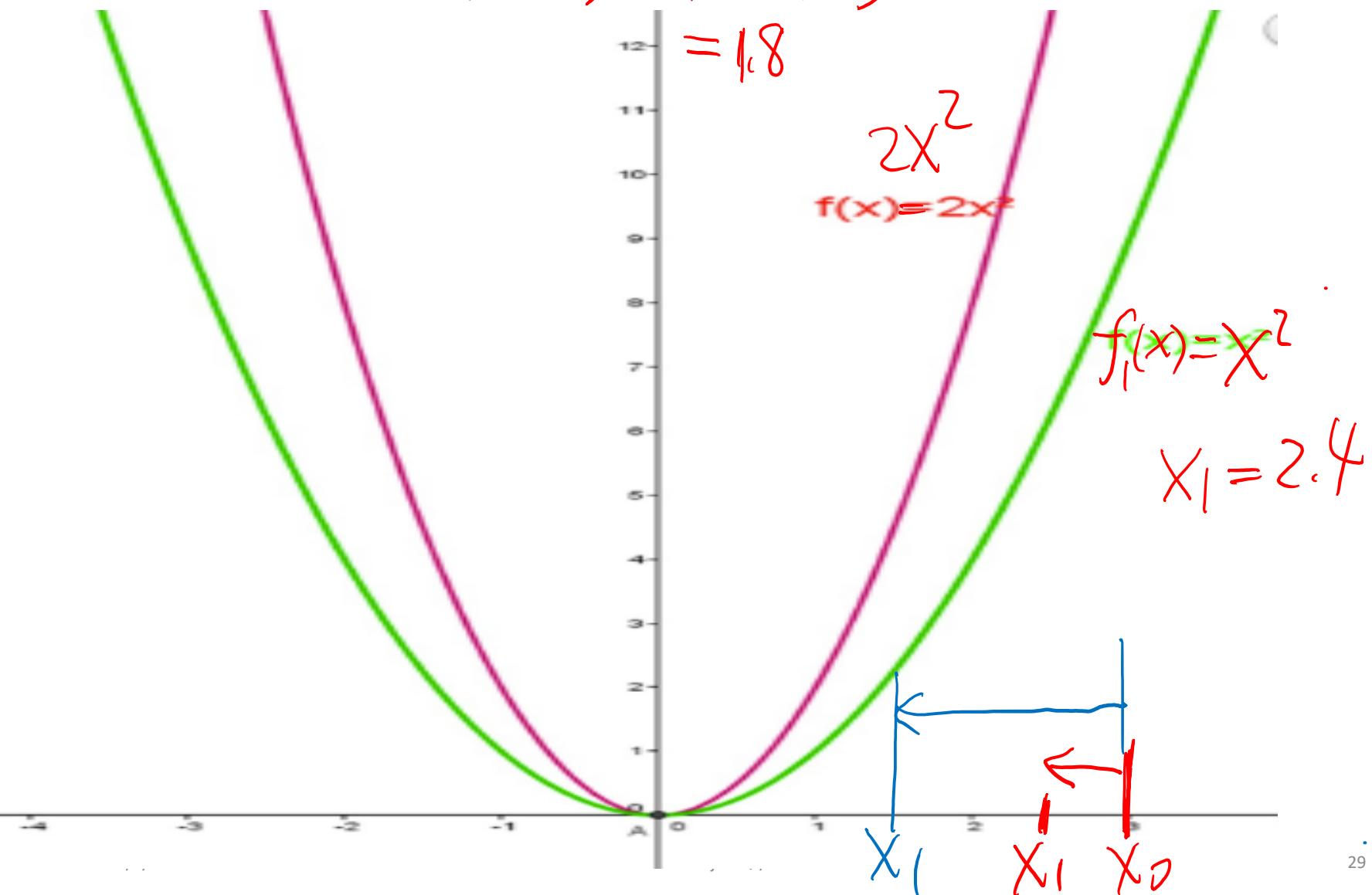
Gradient Descent (Iteratively Optimize)

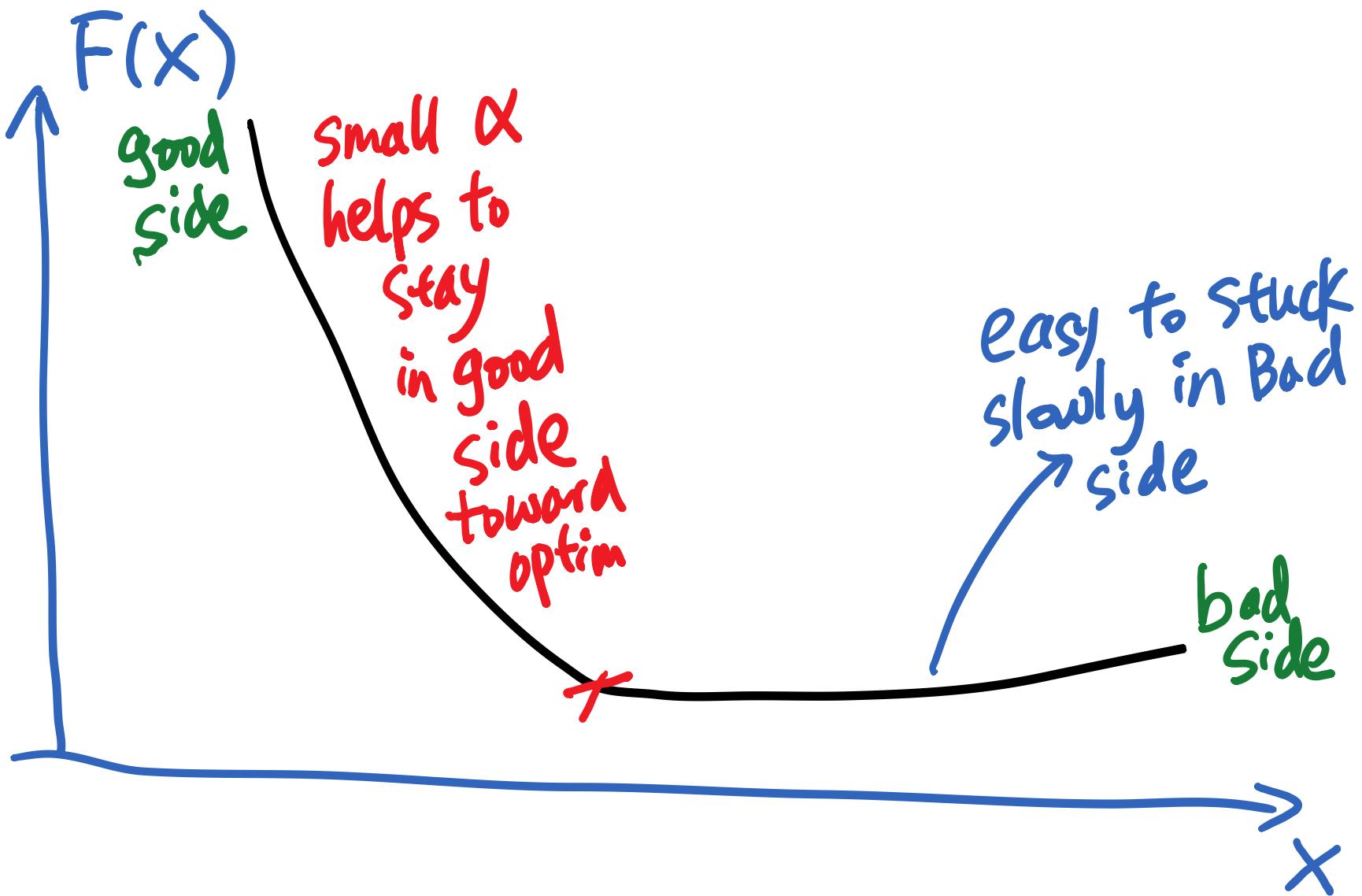
- Learning Rate Matters
- Starting point matters
- Objective function matters

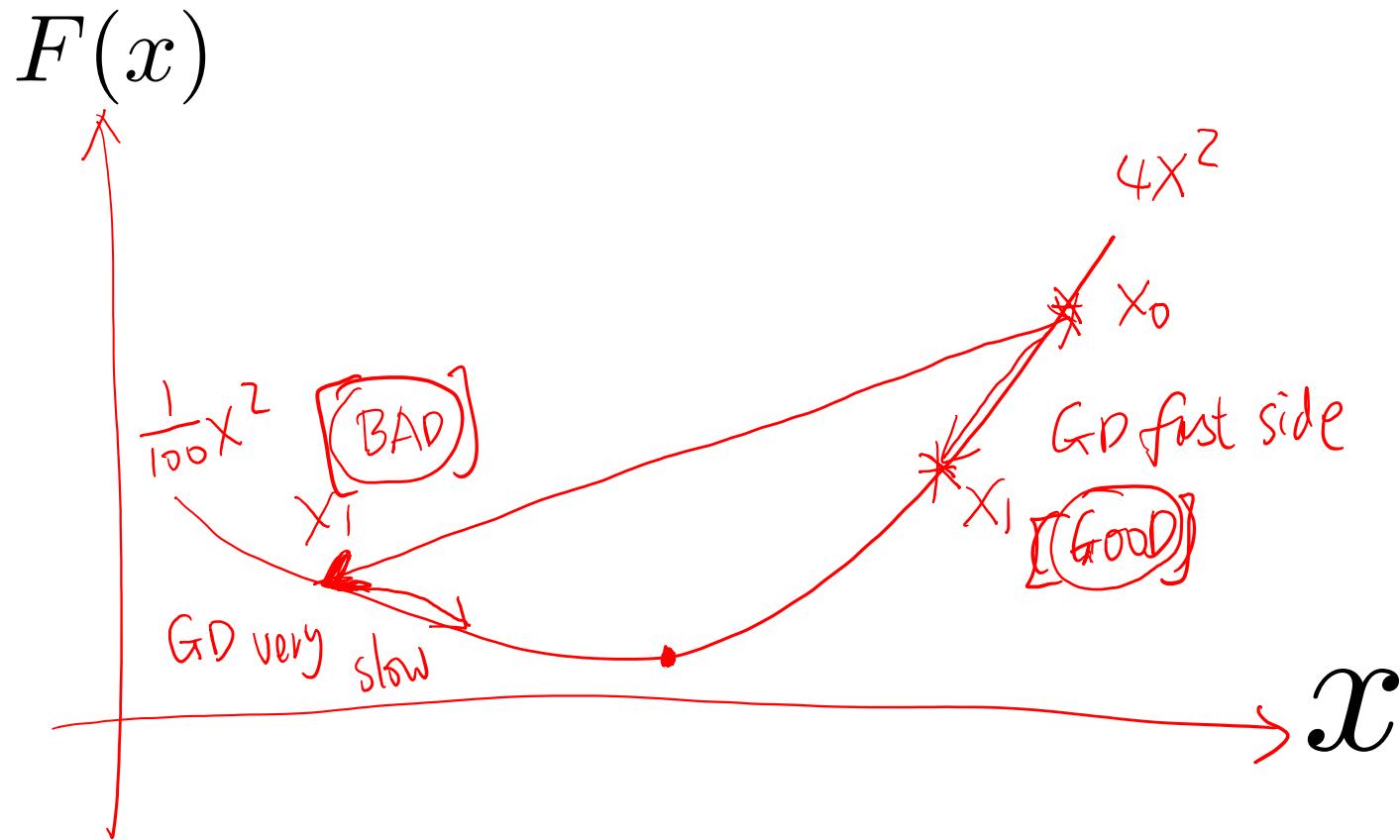
$$x_k = x_{k-1} - \alpha \nabla_x F(x_{k-1})$$

$$x_0 = 3, \alpha = 0.1$$

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







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

x_t	α	$x_{(t+1)}$	$f(x)$
-3	0.1	-2.4	x^2
3	0.1	2.4	x^2
3 -0.6	0.6 0.6	-0.6 0.12	x^2
-0.6	0.1	-0.48	x^2

3

0.1

1.8

 $2x^2$

α

① Small

{

② Smaller

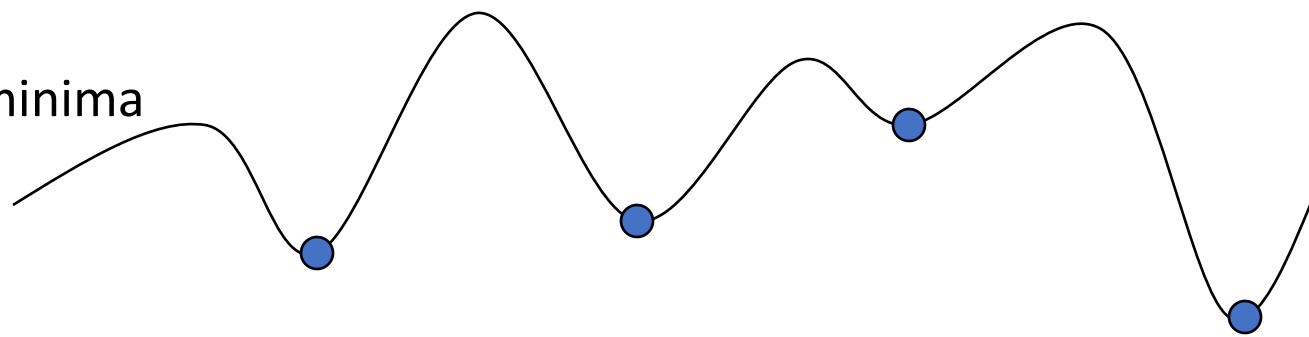
↓

$\alpha \downarrow t \uparrow$

Comments on Gradient Descent Algorithm

- Works on any objective function $F(\mathbf{x})$
 - 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

Today

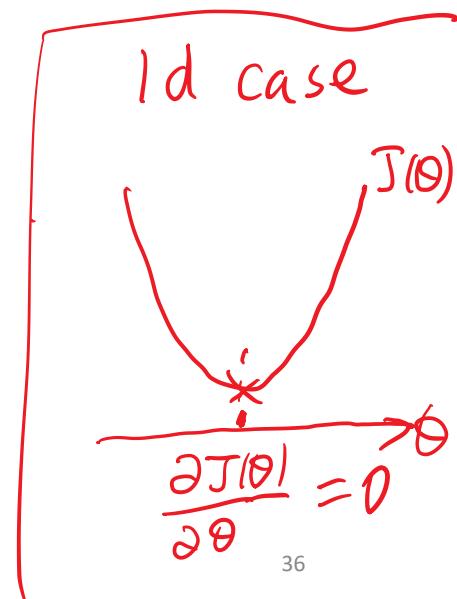
- ❑ More ways to train / perform optimization for linear regression models
 - ❑ Review: Gradient Descent
 - ❑ Gradient Descent (GD) for LR
 - ❑ Stochastic GD (SGD) for LR

Review: Loss function of Least Square LR

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

$$\begin{aligned}
 J(\theta) &= (\underline{x}\theta - y)^T (\underline{x}\theta - y) \frac{1}{2} \\
 &= ((x\theta)^T - y^T)(\underline{x}\theta - y) \frac{1}{2} \\
 &= (\theta^T \underline{x}^T - y^T)(\underline{x}\theta - y) \frac{1}{2} \\
 &= (\underbrace{\theta^T \underline{x}^T \underline{x}\theta - \theta^T \underline{x}^T y - y^T \underline{x}\theta + y^T y}_{\text{since } \theta^T \underline{x}^T y = y^T \underline{x}\theta}) \frac{1}{2} \\
 &= (\underbrace{\theta^T \underline{x}^T \underline{x}\theta}_{\langle x\theta, y \rangle} - \underbrace{2\theta^T \underline{x}^T y}_{\langle y, \underline{x}\theta \rangle} + \underbrace{y^T y}_{J(\theta)}) \frac{1}{2}
 \end{aligned}$$

$\Rightarrow J(\theta)$ quadratic func of θ ;



See handout 4.1 + 4.3 \Rightarrow matrix calculus, partial deri \Rightarrow Gradient

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

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

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

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

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

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

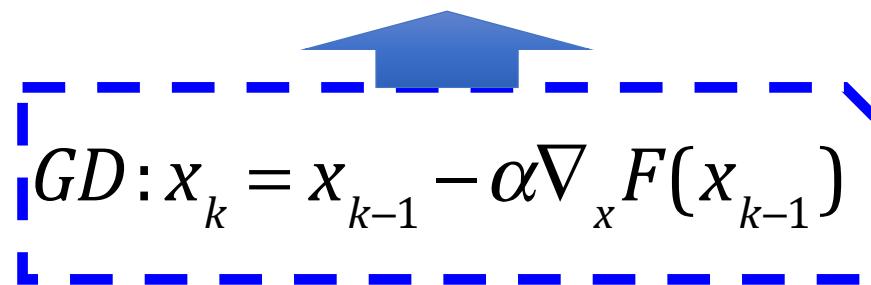
$$= X^T (X \theta - \bar{y})$$

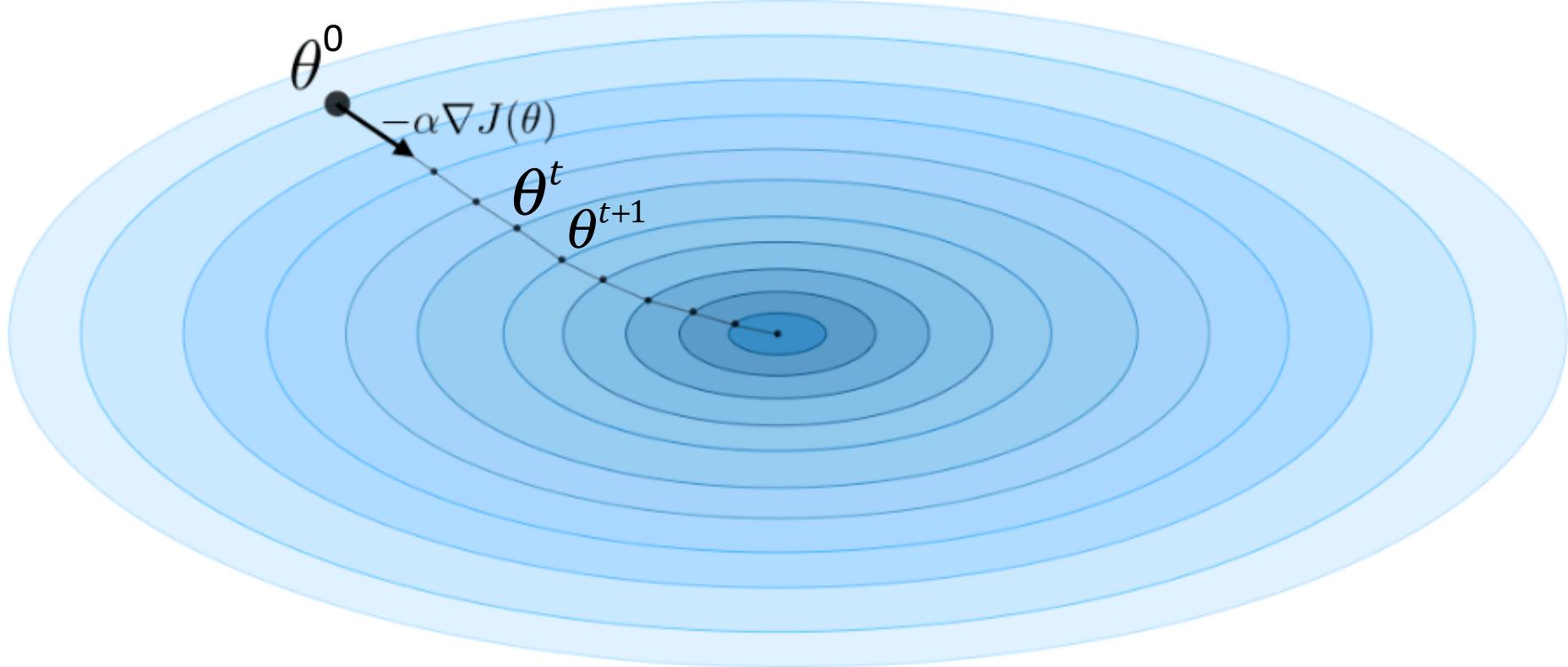
LR with batch GD

- A Batch **gradient descent** algorithm:

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

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


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



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

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

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

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

a vector of

errors on
each point @ θ_t

$$\vec{\theta}^{t+1} = \vec{\theta}^t + \alpha \vec{X}^T (\vec{\bar{y}} - \vec{X} \vec{\theta}^t)$$

$p \times 1$

$p \times 1$

$n \times 1$

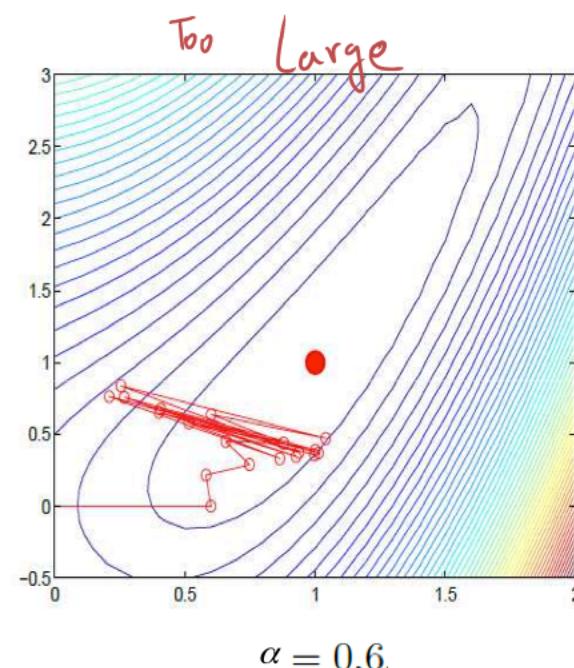
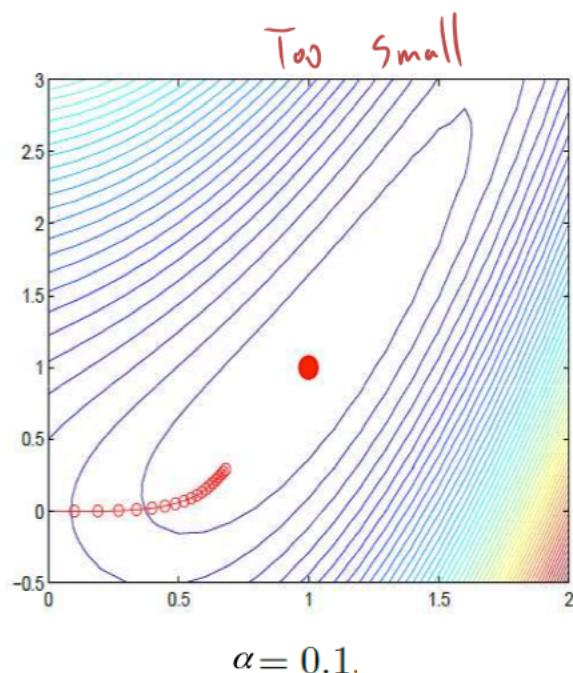
$n \times p$

$p \times 1$

$n \times 1$

$n \times 1$

Choosing the Right Step-Size /Learning-Rate is critical



Today

- ❑ More ways to train / perform optimization for linear regression models
 - ❑ Review: Gradient Descent
 - ❑ Gradient Descent (GD) for LR
 - ❑ Stochastic GD (SGD) for LR

LR with batch GD

- The Cost Function:

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

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

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

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

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

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

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

LR with Stochastic GD →

- Batch GD rule:

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

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

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

➤ A "**stochastic**" descent algorithm, can be used as an **on-line** algorithm

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

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

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

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

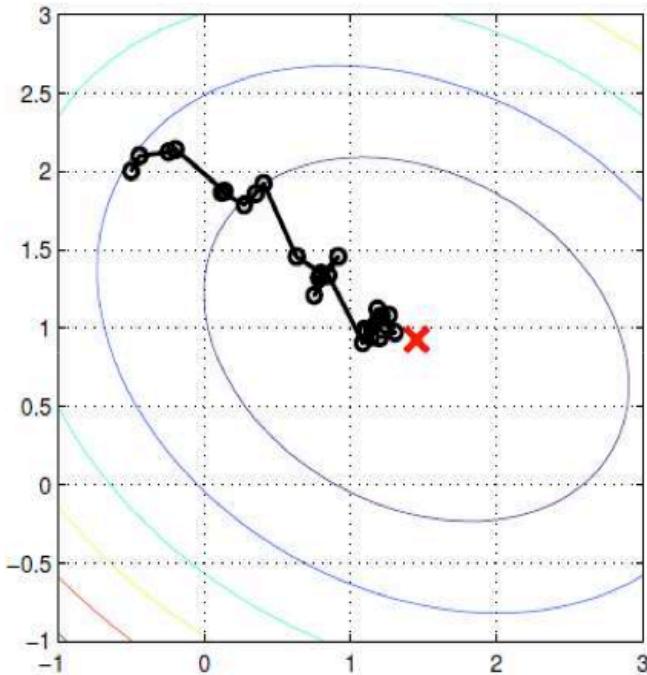
$$= \mathbf{X}^T \begin{bmatrix} \bar{x}_1^T \theta - y_1 \\ \bar{x}_2^T \theta - y_2 \\ \dots \\ \bar{x}_n^T \theta - y_n \end{bmatrix}_{n \times 1} = \begin{bmatrix} 1 & 1 & 1 \\ \bar{x}_1 & \bar{x}_2 & \dots & \bar{x}_n \end{bmatrix} \begin{bmatrix} \bar{x}_1^T \theta - y_1 \\ \bar{x}_2^T \theta -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

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



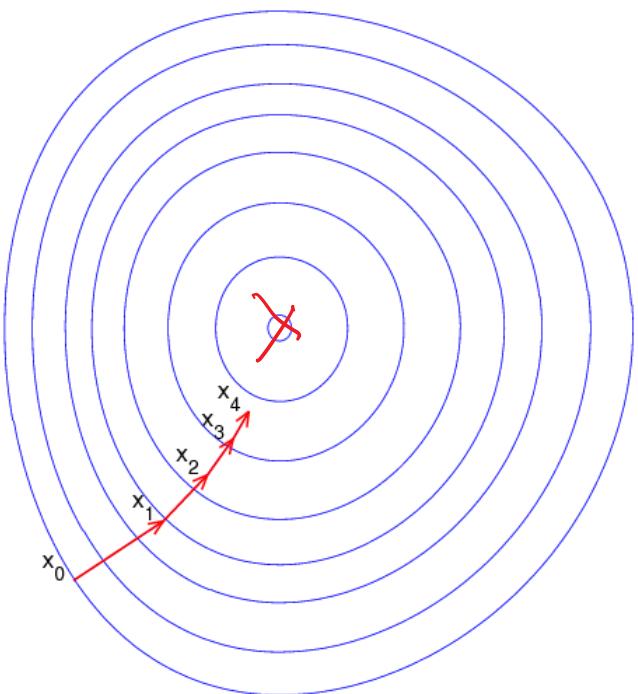
Stochastic gradient descent / Online Learning Algorithm

SGD



GD

versus



Stochastic gradient descent : More variations

- Single-sample:

$$\theta^{t+1} = \theta^t + \alpha (y_i - \vec{X}_i^\top \theta^t) \vec{X}_i$$

- Mini-batch:

$$\theta^{t+1} = \theta^t + \alpha \sum_{j=1}^B (y_j - \vec{X}_j^\top \theta^t) \vec{X}_j$$

e.g. $B=15$

Stochastic gradient descent (1)

- 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
 - $B=1$ standard SGD
 - $B=N$ standard batch GD
 - E.g. $B=50$

Stochastic gradient descent (2)

- 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

B

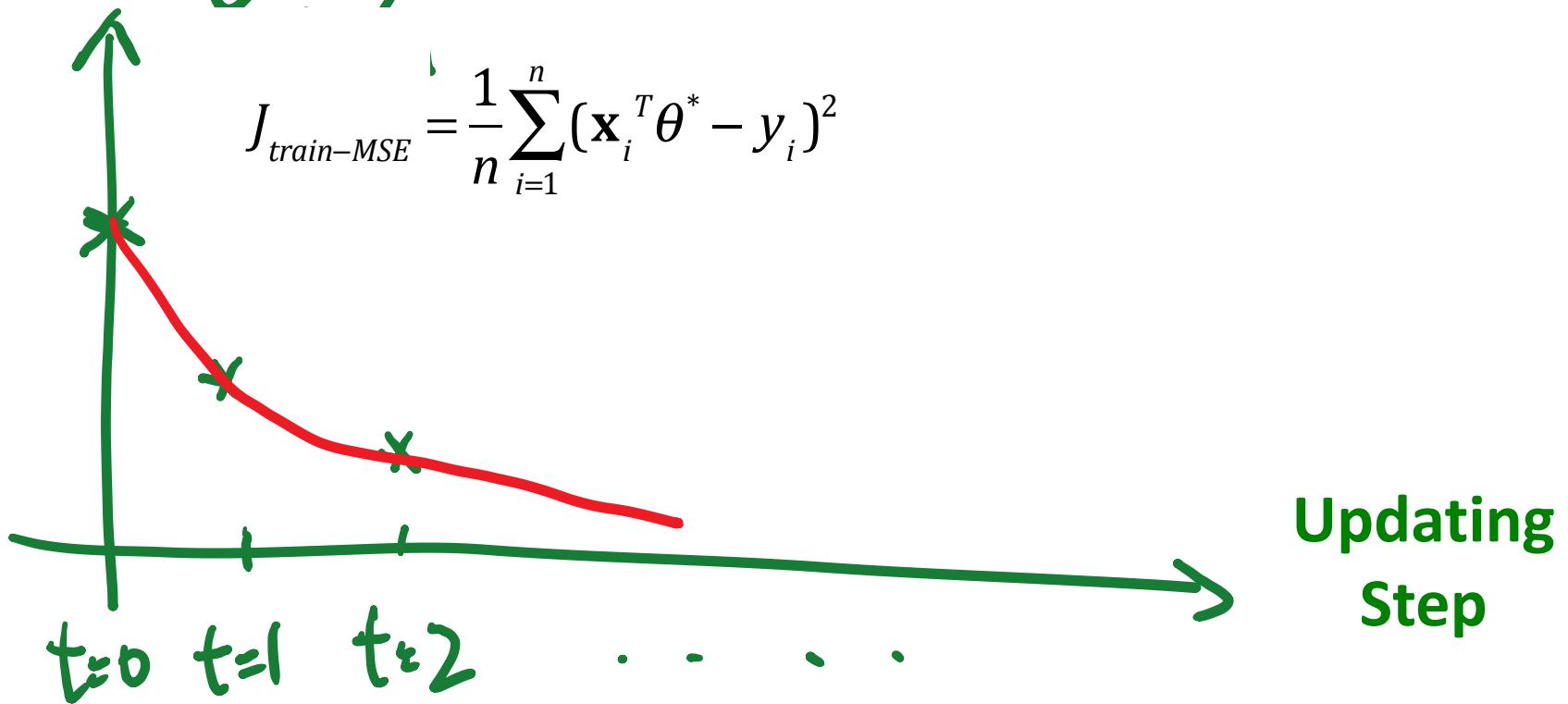
Varying the value B In

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

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

One good plot for GD / SGD

mean-training-loss



- Train MSE Error to observe:

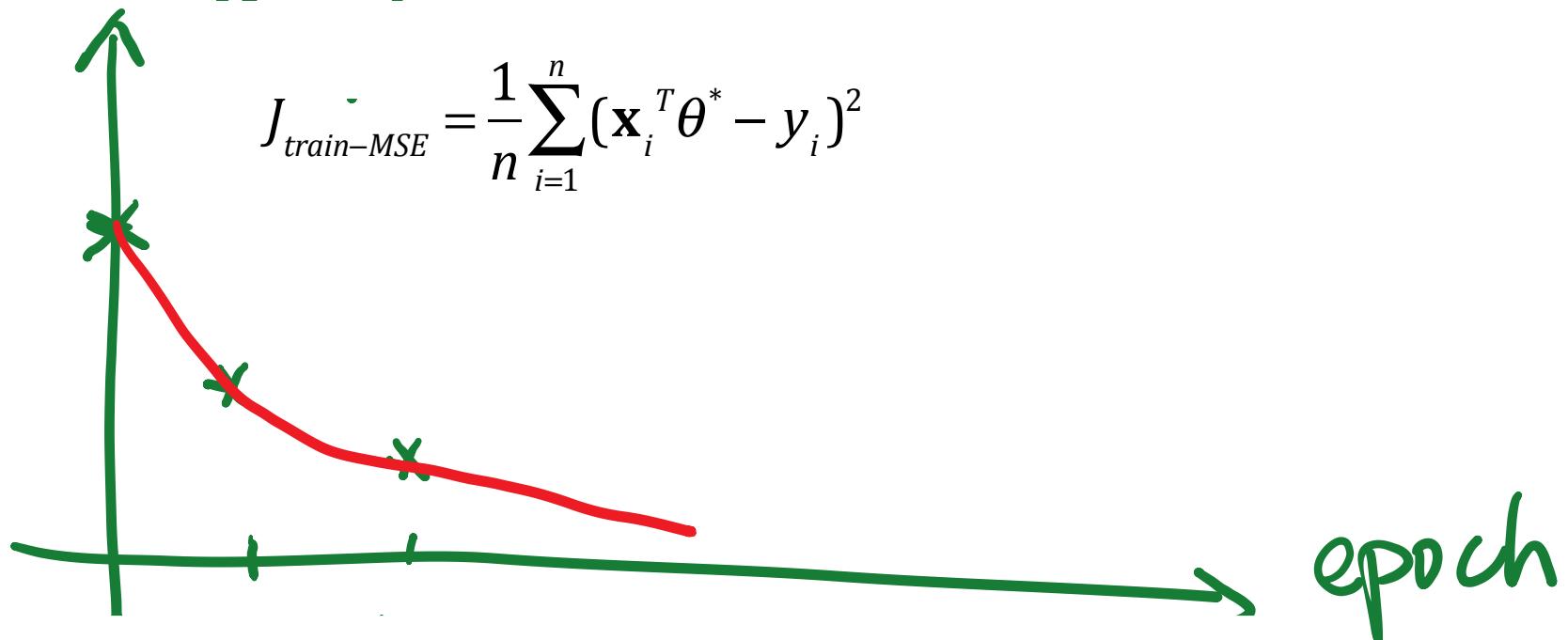
$$J_{train_MSE}^t = \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i^T \boldsymbol{\theta}^t - y_i)^2$$

In many situations, visualizing Train-MSE can be helpful to understand the behavior of your method, e.g., the influence of the hyper parameter you chose; 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



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
-

Summary so far: three 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

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

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

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:

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

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

Update Rule Per Feature
(Variable-Wise)

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: 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 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

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 α , or gradually decrease α .

$$\alpha_0 = 0.05$$

Extra: Computational Cost (Naïve..)

$$\hat{\theta}^* = \left(\begin{matrix} X^T \\ Y^T \end{matrix} \right)^{-1} \begin{matrix} X^T \\ Y^T \end{matrix}$$

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

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

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

References

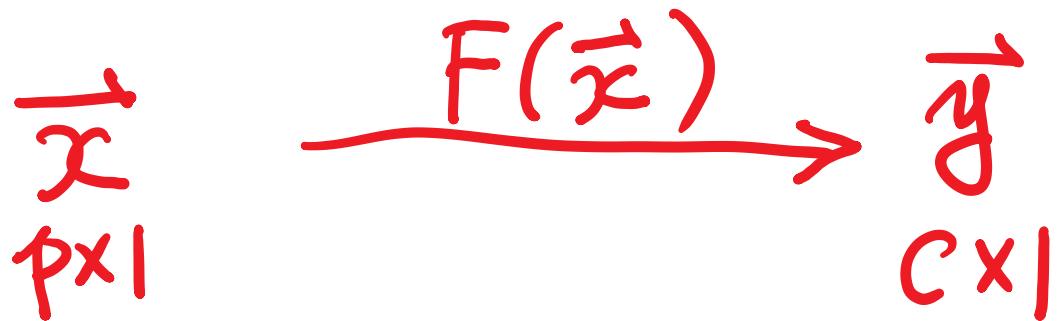
- Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- http://en.wikipedia.org/wiki/Matrix_calculus
- Prof. Nando de Freitas's tutorial slide
- An overview of gradient descent optimization algorithms,
<https://arxiv.org/abs/1609.04747>

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

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

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

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

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 !)

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

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

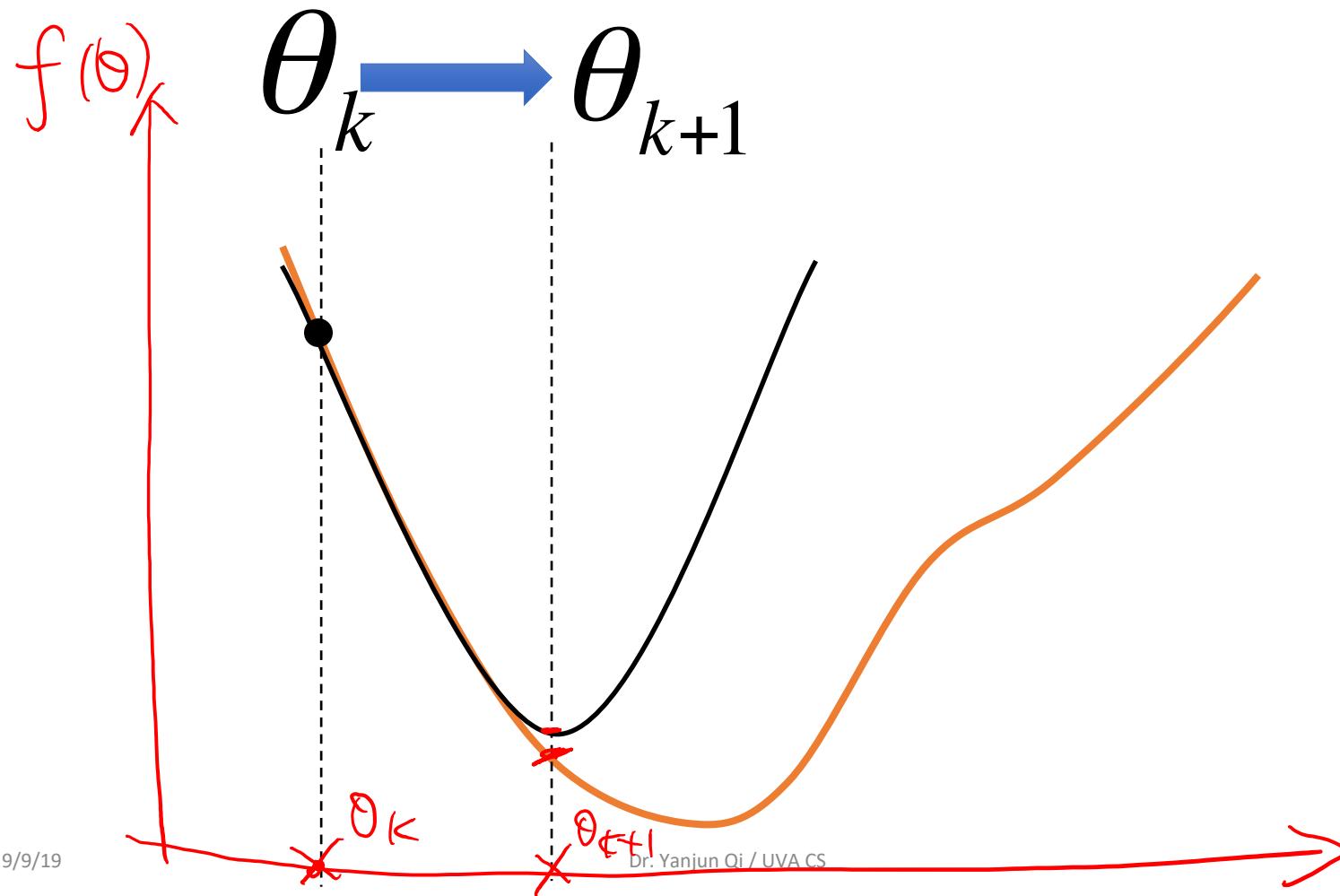
$$\frac{\partial \hat{f}(\theta)}{\partial \theta} = 0 + g_K + \underbrace{\frac{2}{2} H_K \theta - \frac{2}{2} H_K \theta_K}_{\text{See P24 handout}} := 0$$

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

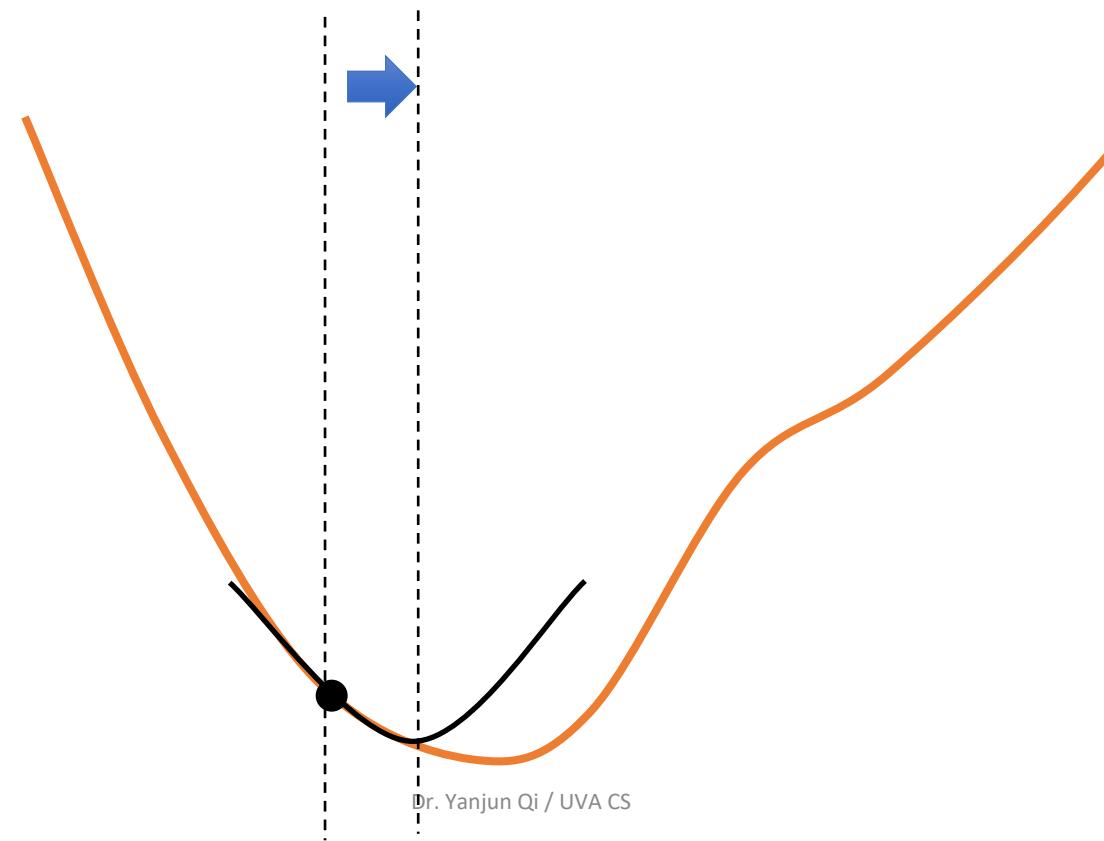
$$\Rightarrow \theta = \theta_K - H_K^{-1} g_K$$

where $\begin{cases} H_K \in \mathbb{R}^{P \times P} \\ g_K \in \mathbb{R}^P \end{cases}$
69

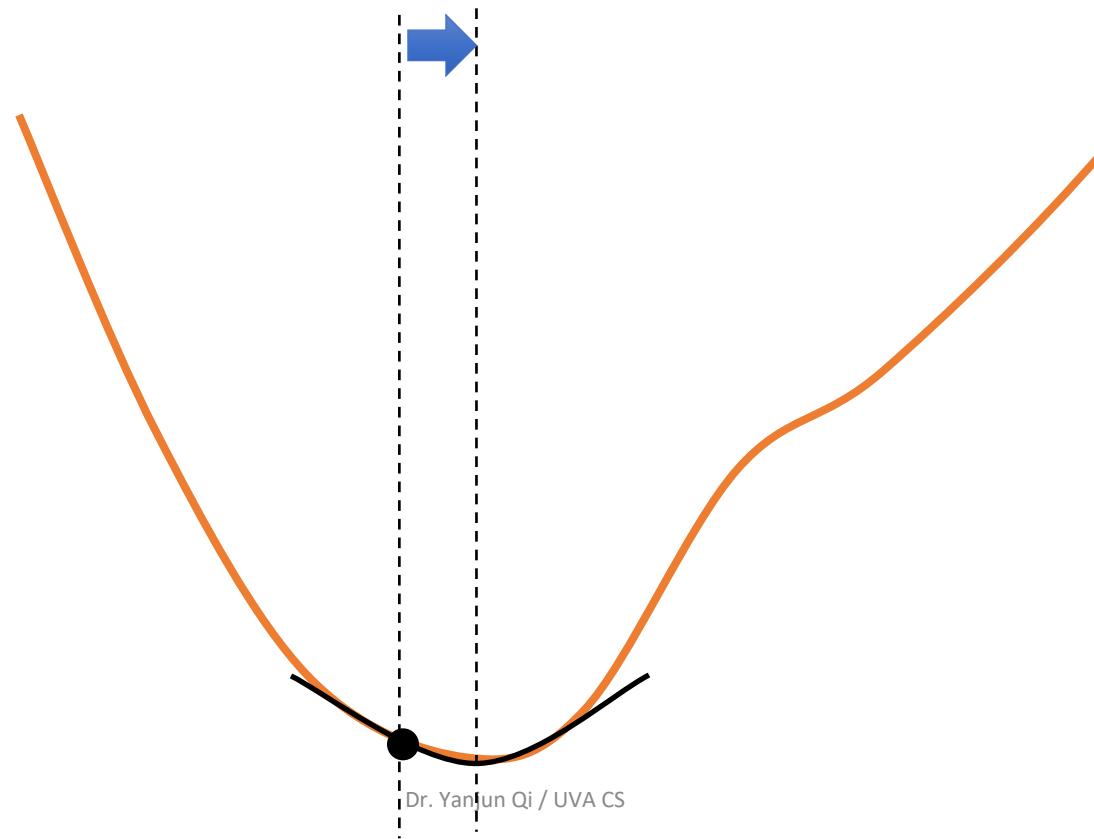
Newton's Method / second-order Taylor series approximation



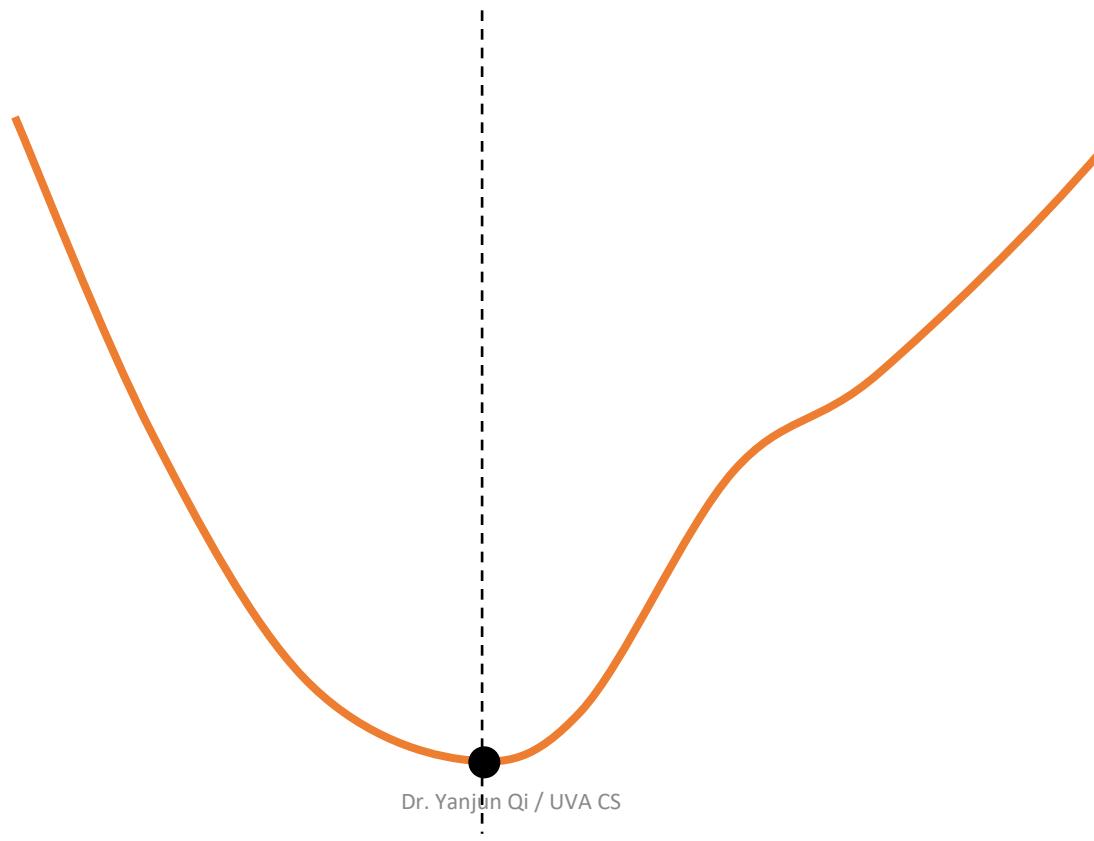
Newton's Method / second-order Taylor series approximation



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 !)

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



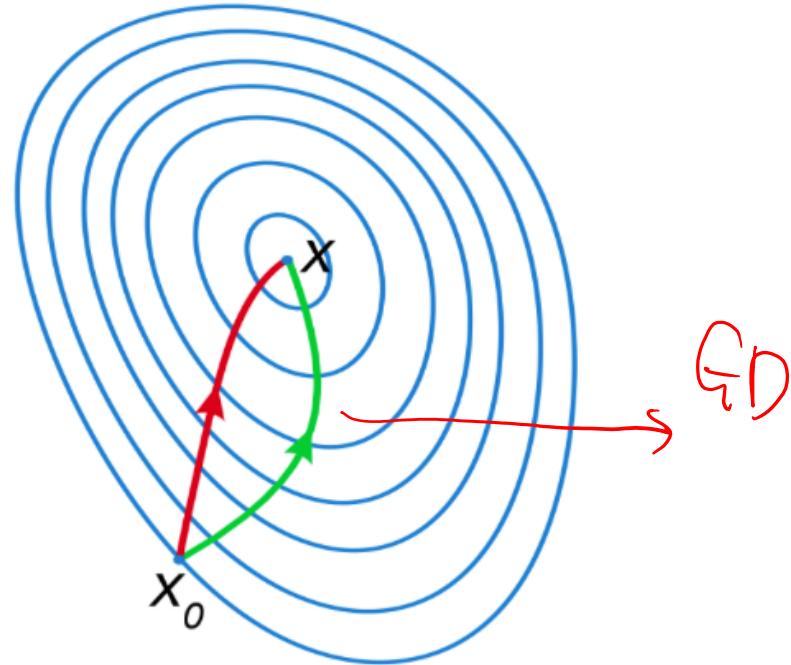
$$\Theta_{k+1} = \Theta_k - \frac{1}{\alpha} g(\Theta_k)$$

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 ...



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

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

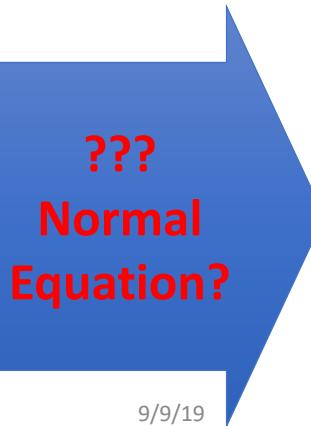
$$H = \nabla_{\theta}^2 J(\theta) = \mathbf{X}^T \mathbf{X}$$

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

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

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

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



Newton's method
for Linear Regression