

# UVA CS 4501 - 001 / 6501 – 007

## Introduction to Machine Learning and Data Mining

### Lecture 4: More optimization for Linear Regression

Yanjun Qi / Jane

University of Virginia  
Department of  
Computer Science

## Last Lecture Recap

- Linear regression (aka **least squares**)
- Learn to derive the least squares estimate by optimization
- Evaluation with Cross-validation

## e.g. SUPERVISED LEARNING

$$f: X \longrightarrow Y$$

- Find function to map **input** space  $X$  to **output** space  $Y$
- **Generalisation**: learn function / hypothesis from **past data** in order to “explain”, “predict”, “model” or “control” **new** data examples

KEY

9/4/14

3

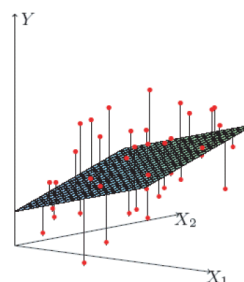
## Linear Regression Models

$$f: X \longrightarrow Y$$

→ e.g. Linear Regression Models

$$\hat{y} = f(x) = \theta_0 + \theta_1 x^1 + \theta_2 x^2$$

- **Features**:  
Living area, distance to campus, # bedroom ...
- **Target  $y$** :  
Rent → Continuous



9/4/14

4

## training / learning goal

- Using matrix form, we get the following general representation of the linear function on train set:

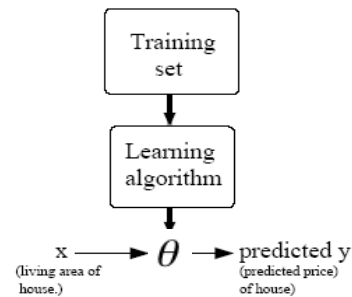
$$\hat{\mathbf{Y}} = \mathbf{X}\theta$$

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

- Our goal is to pick the optimal  $\theta$  that minimize the following cost function:

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

Our goal:



9/4/14

5

## Method I: normal equations

- Write the cost function in matrix form:

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

$$\mathbf{X} = \begin{bmatrix} \text{--} & \mathbf{x}_1^T & \text{--} \\ \text{--} & \mathbf{x}_2^T & \text{--} \\ \vdots & \vdots & \vdots \\ \text{--} & \mathbf{x}_n^T & \text{--} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

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

$$\Downarrow$$

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

9/4/14

6

## e.g. 10 fold Cross Validation

- Divide data into 10 equal pieces
- 9 pieces as training set, the rest 1 as test set
- Collect the scores from the diagonal

model	P1	P2	P3	P4	P5	P6	P7	P8	P9	P10
1	train	train	train	train	train	train	train	train	train	test
2	train	train	train	train	train	train	train	train	test	train
3	train	train	train	train	train	train	train	test	train	train
4	train	train	train	train	train	train	test	train	train	train
5	train	train	train	train	train	test	train	train	train	train
6	train	train	train	train	test	train	train	train	train	train
7	train	train	train	test	train	train	train	train	train	train
8	train	train	test	train	train	train	train	train	train	train
9	train	test	train	train	train	train	train	train	train	train
10	test	train	train	train	train	train	train	train	train	train

9/4/14

## Today

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

9/4/14

8

## Review: Definitions of gradient (from Stanford handout)

Suppose that  $f : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}$  is a function that takes as input a matrix  $A$  of size  $m \times n$  and returns a real value. Then the **gradient** of  $f$  (with respect to  $A \in \mathbb{R}^{m \times n}$ ) is the matrix of

$$\nabla_A f(A) \in \mathbb{R}^{m \times n} = \begin{bmatrix} \frac{\partial f(A)}{\partial A_{11}} & \frac{\partial f(A)}{\partial A_{12}} & \dots & \frac{\partial f(A)}{\partial A_{1n}} \\ \frac{\partial f(A)}{\partial A_{21}} & \frac{\partial f(A)}{\partial A_{22}} & \dots & \frac{\partial f(A)}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f(A)}{\partial A_{m1}} & \frac{\partial f(A)}{\partial A_{m2}} & \dots & \frac{\partial f(A)}{\partial A_{mn}} \end{bmatrix}$$

9/4/14

9

## Review: Definitions of gradient (from Stanford handout)

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

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \in \mathbb{R}^n \quad \text{if } x \in \mathbb{R}^n$$

9/4/14

10

Review: Definitions of gradient  
 (from [http://en.wikipedia.org/wiki/Matrix\\_calculus#Scalar-by-vector](http://en.wikipedia.org/wiki/Matrix_calculus#Scalar-by-vector))

The derivative of a scalar  $y$  function of a matrix  $\mathbf{X}$  of independent variables, with respect to the matrix  $\mathbf{X}$ , is given as

$$\frac{\partial y}{\partial \mathbf{X}} = \begin{bmatrix} \frac{\partial y}{\partial x_{11}} & \frac{\partial y}{\partial x_{21}} & \dots & \frac{\partial y}{\partial x_{p1}} \\ \frac{\partial y}{\partial x_{12}} & \frac{\partial y}{\partial x_{22}} & \dots & \frac{\partial y}{\partial x_{p2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial y}{\partial x_{1q}} & \frac{\partial y}{\partial x_{2q}} & \dots & \frac{\partial y}{\partial x_{pq}} \end{bmatrix}$$

Notice that the indexing of the gradient with respect to  $\mathbf{X}$  is transposed as compared with the indexing of  $\mathbf{X}$ .

9/4/14

11

(from [http://en.wikipedia.org/wiki/Matrix\\_calculus#Scalar-by-vector](http://en.wikipedia.org/wiki/Matrix_calculus#Scalar-by-vector))

The derivative of a scalar  $y$  by a vector  $\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$ , is

$$\frac{\partial y}{\partial \mathbf{x}} = \begin{bmatrix} \frac{\partial y}{\partial x_1} & \frac{\partial y}{\partial x_2} & \dots & \frac{\partial y}{\partial x_n} \end{bmatrix}$$

This gradient is a  $1 \times n$  row vector whose entries respectively contain the  $n$  partial derivatives

9/4/14

12

## Review: Derivative of a Function

$\lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}$  is called the derivative of  $f$  at  $a$ .

We write:  $f'(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$

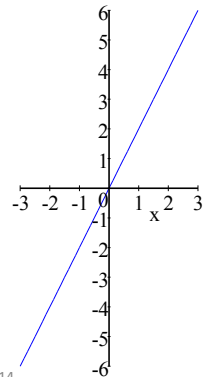
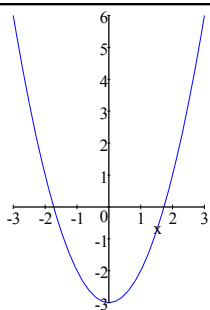
“The derivative of  $f$  with respect to  $x$  is ...”

There are many ways to write the derivative of  $y = f(x)$

→ e.g. define the slope of the curve  $y=f(x)$  at the point  $x$

9/4/14

13



9/4/14

## Review: Derivative of a Quadratic Function

$$y = x^2 - 3$$

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

$$y' = \lim_{h \rightarrow 0} \frac{\cancel{x^2} + 2x\cancel{h} + \cancel{h^2} - \cancel{x^2}}{\cancel{h}}$$

$$y' = \lim_{h \rightarrow 0} 2x + \cancel{h}^0$$

$$y' = 2x$$

14 →

## Today

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

## A little bit more about [ Optimization ]

- Objective function  $F(x)$
- Variables  $\mathcal{X}$
- Constraints

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



Yanjun Qi / UVA CS 4501-01-6501-07

## e.g. Gradient Descent ( Steepest Descent )

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.

17

9/4/14

## Gradient Descent (GD)

- Initialize  $k=0$ , choose  $x_0$
- While  $k < k_{\max}$ 

For the  $k$ -th epoch

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

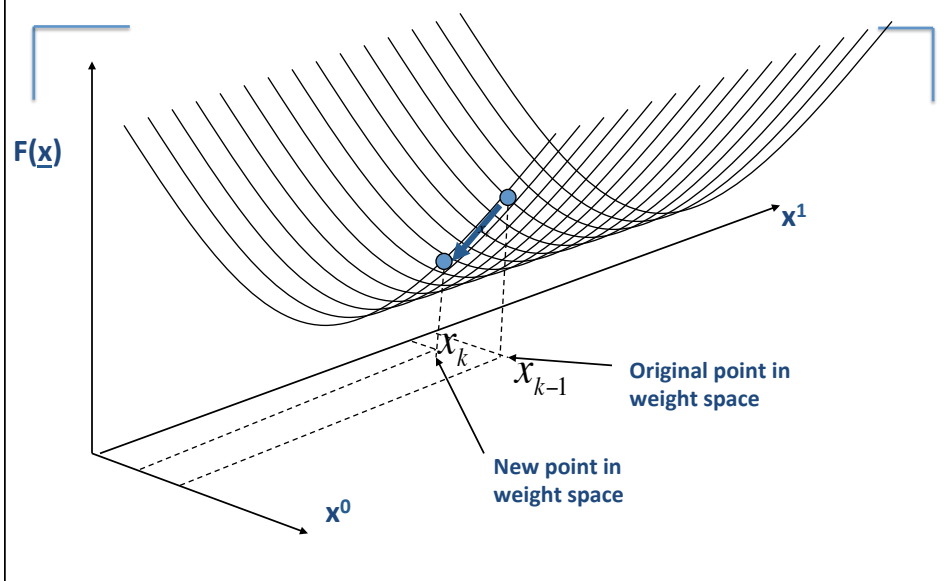
Please READ this note to clarify the confusion : <http://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/gradientDescent.pdf>

18

9/4/14

Yanjun Qi / UVA CS 4501-01-6501-07

## Illustration of Gradient Descent (2D case)

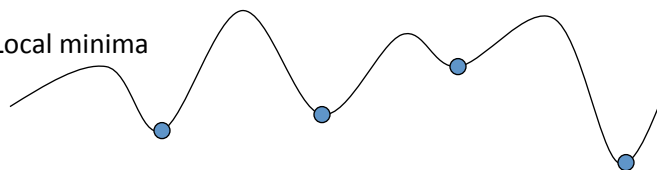


Yanjun Qi / UVA CS 4501-01-6501-07

## Comments on Gradient Descent Algorithm

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

- 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

9/4/14

20

## Method III: LR with batch GD

- The Cost Function:

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

- Consider a **gradient descent** algorithm:

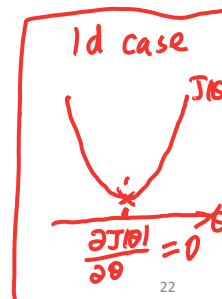
$$\theta_j^{t+1} = \theta_j^t - \alpha \left. \frac{\partial}{\partial \theta_j} J(\theta) \right|_t$$

For the (t+1)-th epoch

$$\begin{aligned} J(\theta) &= (\mathcal{X}\theta - \mathcal{y})^T (\mathcal{X}\theta - \mathcal{y}) \\ &= (\mathcal{X}\theta)^T - \mathcal{y}^T (\mathcal{X}\theta - \mathcal{y}) \\ &= (\theta^T \mathcal{X}^T - \mathcal{y}^T) (\mathcal{X}\theta - \mathcal{y}) \\ &= \theta^T \mathcal{X}^T \mathcal{X} \theta - \underbrace{\theta^T \mathcal{X}^T \mathcal{y} - \mathcal{y}^T \mathcal{X} \theta}_{\substack{\text{since } \theta^T \mathcal{X}^T \mathcal{y} = \mathcal{y}^T \mathcal{X} \theta \\ \langle \mathcal{X}\theta, \mathcal{y} \rangle = \langle \mathcal{y}, \mathcal{X}\theta \rangle}} + \mathcal{y}^T \mathcal{y} \end{aligned}$$

$$= \theta^T \mathcal{X}^T \mathcal{X} \theta - 2 \theta^T \mathcal{X}^T \mathcal{y} + \mathcal{y}^T \mathcal{y}$$

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



Yanjun Qi / UVA CS 4501-01-6501-07

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

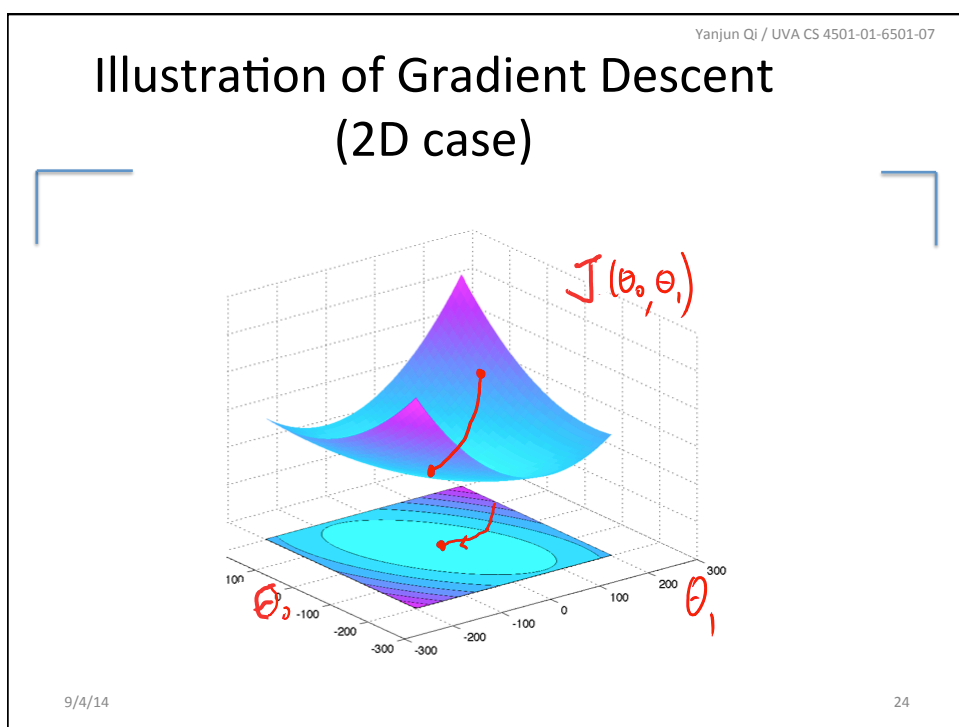
$$\nabla_{\theta} (\theta^T X^T X \theta) = 2 X^T X \theta \quad (P24)$$

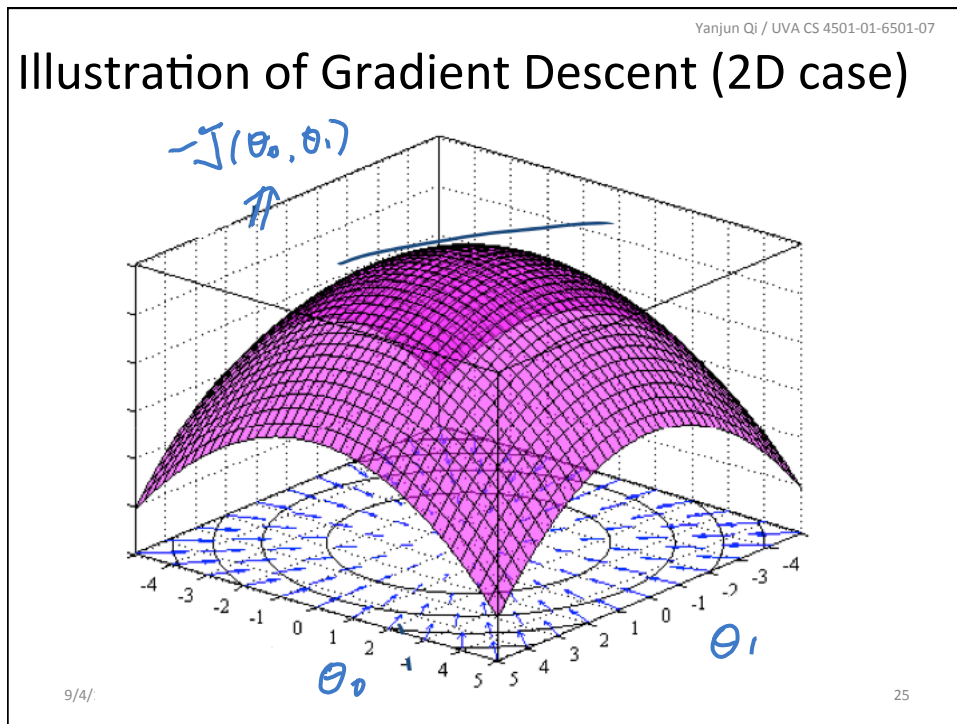
$$\nabla_{\theta} (-2 \theta^T X^T y) = -2 X^T y \quad (P24)$$

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

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

9/4/1423





Yanjun Qi / UVA CS 4501-01-6501-07

$$\begin{aligned}
 \nabla_{\theta} J(\theta) &= X^T X \theta - X^T Y \\
 &= X^T (X \theta - Y) \\
 &= X^T \left( \begin{matrix} -x_1^T \\ -x_2^T \\ \vdots \\ -x_n^T \end{matrix} \theta - \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \right) \\
 &= \sum_{p=1}^n \begin{bmatrix} x_1^T \theta - y_1 \\ x_2^T \theta - y_2 \\ \vdots \\ x_n^T \theta - y_n \end{bmatrix} = \begin{bmatrix} 1 & 1 & \dots & 1 \\ x_1 & x_2 & \dots & x_n \\ \vdots & \vdots & \dots & \vdots \\ 1 & 1 & \dots & 1 \end{bmatrix} \begin{bmatrix} x_1^T \theta - y_1 \\ \vdots \\ x_n^T \theta - y_n \end{bmatrix} \\
 &= \sum_{i=1}^n x_i \cdot (x_i^T \theta - y_i)
 \end{aligned}$$

9/4/14

26

## LR with batch GD

- Steepest descent / GD

– Note that:

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

Update Rule Per  
Feature Variable-  
Wise

$$\nabla_{\theta} J = \left[ \frac{\partial}{\partial \theta_1} J, \dots, \frac{\partial}{\partial \theta_k} J \right]^T = - \sum_{i=1}^n (y_n - \mathbf{x}_n^T \theta) \mathbf{x}_n$$

Based on Stanford  
Handout's Definition  
of Gradient

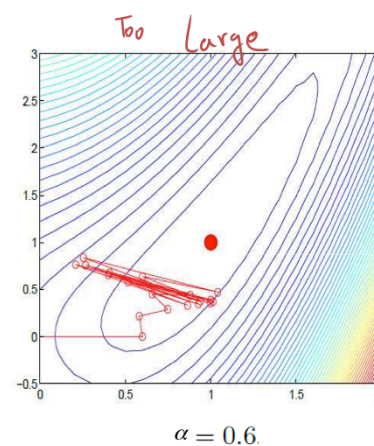
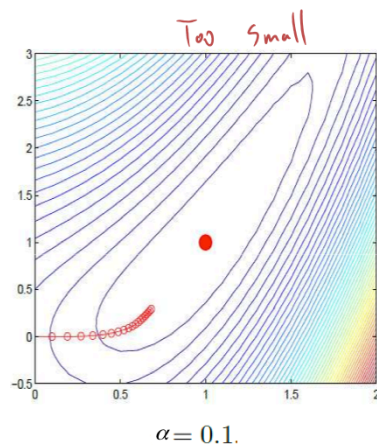
$$\theta^{t+1} = \theta^t + \alpha \sum_{i=1}^n (y_n - \mathbf{x}_n^T \theta^t) \mathbf{x}_n$$

– This is as a **batch** gradient descent algorithm

Please READ this note to clarify the confusion : <http://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/gradientDescent.pdf>

9/4/14

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



9/4/14

28

## Method III: LR with Stochastic GD →

- Now we have the following descent rule:

- For a single training point, we have:

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

- This is known as the Least-Mean-Square update rule, or the Widrow-Hoff learning rule
- This is actually a "stochastic", "coordinate" descent algorithm
- This can be used as a on-line algorithm

9/4/14

29

## Summary: three ways to learn LR

- Normal equations

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

- 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 or Steepest descent

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

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

- Stochastic LMS update rule

$$\theta_j^{t+1} = \theta_j^t + \alpha(y_n - \mathbf{x}_n^T \theta^t) x_{n,j}$$

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

9/4/14

30

## Direct (normal equation) vs. Iterative (GD) 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 or steepest gradient
  - Converging in a limiting sense
  - But more attractive in large practical problems
  - Caution is needed for deciding the learning rate  $\alpha$

## Convergence rate

- **Theorem:** the steepest descent 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  $\rightarrow$

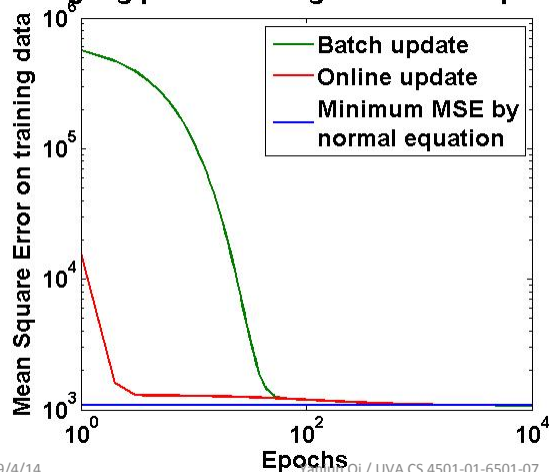
$$0 < \alpha < 2/\lambda_{\max}[X^T X]$$

- A formal analysis of LMS need more math; in practice, one can use a small  $\alpha$ , or gradually decrease  $\alpha$ .



## Convergence Curves, for an example

Log-log plot of training MSE versus epochs



- For the batch method, the training MSE is initially large due to uninformed initialization
- In the online update,  $N$  updates for every epoch reduces MSE to a much smaller value.

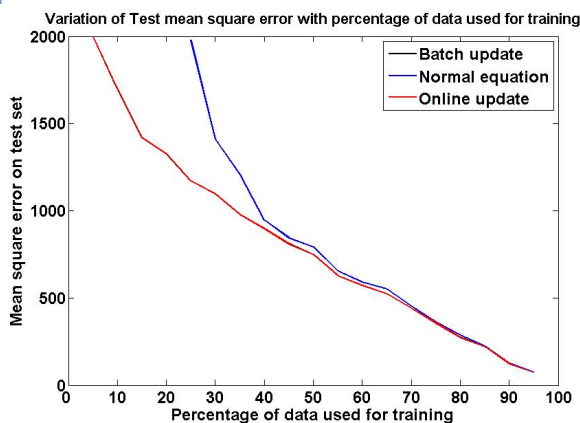
9/4/14

Yanjun Qi / UVA CS 4501-01-6501-07

33

## Performance vs. Training Size for an example

Yanjun Qi / UVA CS 4501-01-6501-07



- The results from B and O update are almost identical. So the plots coincide.
- The test MSE from the normal equation is more than that of B and O during small training. This is probably due to overfitting.
- In B and O, since only 2000 (for example) iterations are allowed at most. This roughly acts as a mechanism that avoids overfitting.

9/4/14

34

Yanjun Qi / UVA CS 4501-01-6501-07

## Geometric Interpretation of Least Mean Square Solution

- The predictions on the training data are:
- Note that  $\hat{\bar{y}} = X\theta^* = X(X^T X)^{-1} X^T \bar{y}$

and

$$\hat{\bar{y}} - \bar{y} = (X(X^T X)^{-1} X^T - I)\bar{y}$$

$$X^T(\hat{\bar{y}} - \bar{y}) = X^T(X(X^T X)^{-1} X^T - I)\bar{y}$$

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

$$= \mathbf{0} \quad !!$$

$\hat{\bar{y}}$  → the orthogonal projection of the true  $\bar{y}$  vector into the space spanned by the columns of  $X$

35

Yanjun Qi / UVA CS 4501-01-6501-07

## Today Recap

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

36

## References

- Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- **Notes about Gradient Descent from Toussaint:**  
(please read) <http://ipvs.informatik.uni-stuttgart.de/mlr/marc/notes/gradientDescent.pdf>
- [http://en.wikipedia.org/wiki/Matrix\\_calculus#Scalar-by-vector](http://en.wikipedia.org/wiki/Matrix_calculus#Scalar-by-vector)