

# UVA CS 6316 – Fall 2015 Graduate: Machine Learning

## Lecture 3: Linear Regression

Dr. Yanjun Qi

University of Virginia

Department of  
Computer Science

9/14/15

1

**HW1 OUT TOMORROW /  
DUE IN TWO WEEKS**

9/14/15

2

Where we are ? ➔

## Five major sections of this course

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

Today ➔

## Regression (supervised)

- ☐ Four ways to train / perform optimization for linear regression models
  - ☐ Normal Equation
  - ☐ Gradient Descent (GD)
  - ☐ Stochastic GD
  - ☐ Newton's method
- ☐ Supervised regression models
  - ☐ Linear regression (LR)
  - ☐ LR with non-linear basis functions
  - ☐ Locally weighted LR
  - ☐ LR with Regularizations

# Today

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

9/14/15

5

$X_1$	$X_2$	$X_3$	$Y$

A Dataset  
for **regression**

$$f : X \rightarrow Y$$

continuous  
valued  
variable

- **Data**/points/instances/examples/samples/records: [ rows ]
- **Features**/attributes/dimensions/independent variables/covariates/predictors/regressors: [ columns, except the last ]
- **Target**/outcome/response/label/dependent variable: special column to be predicted [ last column ]

9/14/15

6

# For Example, Machine learning for apartment hunting



- Now you've moved to Charlottesville !!  
And you want to find the **most reasonably priced** apartment satisfying your **needs**:  
square-ft., # of bedroom, distance to campus ...

Living area (ft <sup>2</sup> )	# bedroom	Rent (\$)
230	1	600
506	2	1000
433	2	1100
109	1	500
...		
150	1	?
270	1.5	?

9/14/15

7

# For Example, Machine learning for apartment hunting

Living area (ft <sup>2</sup> )	# bedroom	Rent (\$)
230	1	600
506	2	1000
433	2	1100
109	1	500
...		
150	1	?
270	1.5	?



*features labels*

	$X_1$	$X_2$	$Y$
$s_1$			
$s_2$			
$s_3$			
$s_4$			
$s_5$			
$s_6$			

9/14/15

8

# Linear SUPERVISED Regression

$$f: X \longrightarrow Y$$

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

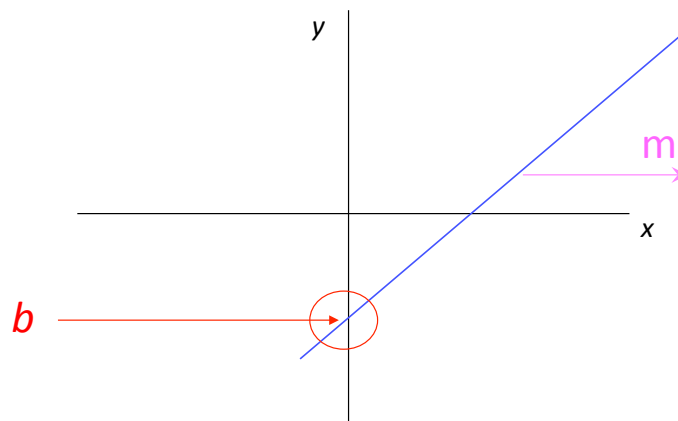
9/14/15

9

## Remember this: “Linear”? (1D case)

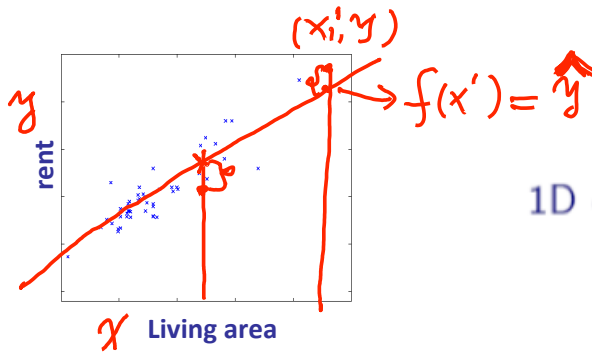
- $y = mx + b$ ?

A slope of 2 (i.e.  $m=2$ ) means that every 1-unit change in  $X$  yields a 2-unit change in  $Y$ .

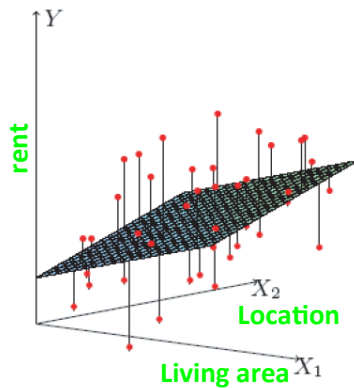


9/14/15

10



1D case ( $\mathcal{X} = \mathbb{R}$ ): a line



$\mathcal{X} = \mathbb{R}^2$ : a plane

11

## A new representation (for single sample)

- Assume that each sample  $\mathbf{x}$  is a column vector,

- Here we assume a pseudo "feature"  $x^0=1$  (this is the **intercept** term), and **RE-define** the feature vector to be:

$$\mathbf{x}^T = [x^0, x^1, x^2, \dots, x^{p-1}]$$

- the parameter vector  $\theta$  is also a column vector

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{p-1} \end{bmatrix}$$



$$\begin{aligned} \hat{y} &= f(\mathbf{x}) \\ &= \mathbf{x}^T \theta = \theta^T \mathbf{x} \end{aligned}$$

$$\mathbf{x} = \begin{bmatrix} x^1 \\ x^2 \\ \vdots \\ x^{p-1} \end{bmatrix} \quad \Downarrow \quad \mathbf{x} = \begin{bmatrix} 1 \\ x^1 \\ x^2 \\ \vdots \\ x^{p-1} \end{bmatrix}$$

$$\hat{y} = f(\mathbf{x}) = \theta_0 + \theta_1 x^1 + \theta_2 x^2 + \dots + \theta_{p-1} x^{p-1}$$

1D:  $y = mx + b = \begin{bmatrix} x & 1 \end{bmatrix} \begin{bmatrix} m \\ b \end{bmatrix} = \vec{x}^T \theta$

$p=2$

$\Rightarrow$

	$x^0$	$x^1$	$y$
$s_1$	1	$x'_1$	$y_1$
$s_2$	1	$x'_2$	$y_2$
$\vdots$			
$s_n$	1	$x'_n$	$y_n$

9/14/15

13

## Training / learning problem

- Now represent the whole Training set (with  $n$  samples) as matrix form :

$$\mathbf{X} = \begin{bmatrix} -- & \mathbf{x}_1^T & -- \\ -- & \mathbf{x}_2^T & -- \\ \vdots & \vdots & \vdots \\ -- & \mathbf{x}_n^T & -- \end{bmatrix} = \begin{bmatrix} x_1^0 & x_1^1 & \dots & x_1^{p-1} \\ x_2^0 & x_2^1 & \dots & x_2^{p-1} \\ \vdots & \vdots & \vdots & \vdots \\ x_n^0 & x_n^1 & \dots & x_n^{p-1} \end{bmatrix} \quad \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

$n \times p$        $n \times 1$

9/14/15

14

# REVIEW: Special Uses for Matrix Multiplication

## • Matrix-Vector Products (I)

Given a matrix  $A \in \mathbb{R}^{m \times n}$  and a vector  $x \in \mathbb{R}^n$ , their product is a vector  $y = Ax \in \mathbb{R}^m$ .

If we write  $A$  by rows, then we can express  $Ax$  as,

$$y = Ax = \begin{bmatrix} - & a_1^T & - \\ - & a_2^T & - \\ & \vdots & \\ - & a_m^T & - \end{bmatrix} x = \begin{bmatrix} a_1^T x \\ a_2^T x \\ \vdots \\ a_m^T x \end{bmatrix}.$$

# Training / learning problem

## • Represent as matrix form:

– Predicted output

$$\hat{Y} = \underset{n \times p}{X} \underset{p \times 1}{\theta} = \begin{bmatrix} f(\mathbf{x}_1) \\ f(\mathbf{x}_2) \\ \vdots \\ f(\mathbf{x}_n) \end{bmatrix} = \underset{n \times 1}{\begin{bmatrix} \mathbf{x}_1^T \theta \\ \mathbf{x}_2^T \theta \\ \vdots \\ \mathbf{x}_n^T \theta \end{bmatrix}}$$

– Labels (given output value)

$$Y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$



# Training / learning goal

- Using matrix form, we get the following general representation of the linear regression function:

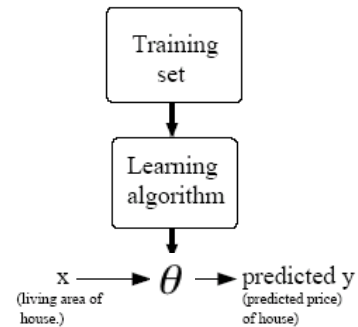
$$\hat{Y} = \mathbf{X}\theta$$

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

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

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

**Our goal:**



9/14/15

SSE: Sum of squared error

17

## Today

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

9/14/15

18

## Method I: normal equations

- Write the cost function in matrix form:

$$\begin{aligned}
 J(\theta) &= \frac{1}{2} \sum_{i=1}^n (\mathbf{x}_i^T \theta - y_i)^2 \\
 &= \frac{1}{2} (X\theta - \bar{y})^T (X\theta - \bar{y}) \\
 &= \frac{1}{2} (\theta^T X^T X \theta - \theta^T X^T \bar{y} - \bar{y}^T X \theta + \bar{y}^T \bar{y})
 \end{aligned}
 \quad
 \mathbf{X} = \begin{bmatrix} -- & \mathbf{x}_1^T & -- \\ -- & \mathbf{x}_2^T & -- \\ \vdots & \vdots & \vdots \\ -- & \mathbf{x}_n^T & -- \end{bmatrix}
 \quad
 \mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

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

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

The normal equations

$$\Downarrow$$

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

9/14/15

19

## Review: Special Uses for Matrix Multiplication

- Dot (or Inner) Product of two Vectors  $\langle \mathbf{x}, \mathbf{y} \rangle$

which is the sum of products of elements in similar positions for the two vectors

$$\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle \quad \mathbf{a}^T \mathbf{b} = \mathbf{b}^T \mathbf{a}$$

$$\text{Where } \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^T \mathbf{y} \in \mathbb{R} = \begin{bmatrix} x_1 & x_2 & \cdots & x_n \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} = \sum_{i=1}^n x_i y_i.$$

9/14/15

20

# Review: Matrix Calculus:

## Types of Matrix Derivatives

	Scalar	Vector	Matrix
Scalar	$\frac{dy}{dx}$	$\frac{d\mathbf{y}}{dx} = \left[ \frac{\partial y_i}{\partial x} \right]$	$\frac{d\mathbf{Y}}{dx} = \left[ \frac{\partial y_{ij}}{\partial x} \right]$
Vector	$\frac{dy}{d\mathbf{x}} = \left[ \frac{\partial y}{\partial x_j} \right]$	$\frac{d\mathbf{y}}{d\mathbf{x}} = \left[ \frac{\partial y_i}{\partial x_j} \right]$	
Matrix	$\frac{dy}{d\mathbf{X}} = \left[ \frac{\partial y}{\partial x_{ji}} \right]$		

By Thomas Minka. Old and New Matrix Algebra Useful for Statistics

# Review: Special Uses for

## Matrix Multiplication

- **Sum the Squared Elements of a Vector → L2 norm**

- Premultiply a column vector  $\mathbf{a}$  by its transpose – If

$$\mathbf{a} = \begin{bmatrix} 5 \\ 2 \\ 8 \end{bmatrix}$$

then premultiplication by a row vector  $\mathbf{a}^T$

$$\mathbf{a}^T = \begin{bmatrix} 5 & 2 & 8 \end{bmatrix}$$

will yield the sum of the squared values of elements for  $\mathbf{a}$ , i.e.

$$|\mathbf{a}|_2^2 = \mathbf{a}^T \mathbf{a} = \begin{bmatrix} 5 & 2 & 8 \end{bmatrix} \begin{bmatrix} 5 \\ 2 \\ 8 \end{bmatrix} = 5^2 + 2^2 + 8^2 = 93$$

Details for slide [19]:

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

$$= (\underbrace{\sum}_{n \times p} \theta - \underbrace{y}_{p \times 1})^T (\underbrace{\sum}_{n \times 1} \theta - \underbrace{y}_{n \times 1})$$

Since  $\underbrace{w^T w}_{\substack{\Downarrow \\ x\theta}} = \|w\|_2^2 = \sum_{i=1}^n w_i^2$

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

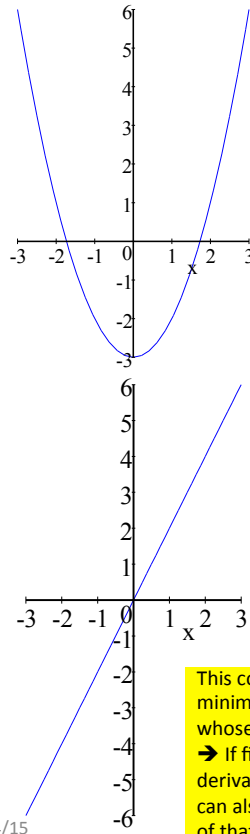
$$= \frac{1}{2} (\sum \theta - y)^T (\sum \theta - y)$$

$$= \frac{1}{2} (\theta^T \sum^T - y^T) (\sum \theta - y)$$

$$= \frac{1}{2} (\theta^T \sum^T \sum \theta - \underbrace{\theta^T \sum^T y - y^T \sum \theta}_{\text{same}} + y^T y)$$

Since  $\theta^T \sum^T y = y^T \sum \theta$   
 $\langle x\theta, y \rangle = \langle y, \sum \theta \rangle$

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



## Review: Derivative of a Quadratic Function

$$y = x^2 - 3$$

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

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

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

$$y' = 2x$$

$$y'' = 2$$

This convex function is minimized @ the unique point whose derivative (slope) is zero.  
 → If finding zeros of the derivative of this function, we can also find minima (or maxima) of that function.

25

## Review: Convex function

- Intuitively, a convex function (1D case) has a single point at which the derivative goes to zero, and this point is a minimum.
- Intuitively, a function  $f$  (1D case) is convex on the range  $[a,b]$  if a function's second derivative is positive every-where in that range.
- Intuitively, if a function's Hessians is psd (positive semi-definite!), this (multivariate) function is Convex
  - Intuitively, we can think "Positive definite" matrices as analogy to positive numbers in matrix case

# Review: Some important rules for taking derivatives

- Scalar multiplication:  $\partial_x[af(x)] = a[\partial_x f(x)]$
- Polynomials:  $\partial_x[x^k] = kx^{k-1}$
- Function addition:  $\partial_x[f(x) + g(x)] = [\partial_x f(x)] + [\partial_x g(x)]$
- Function multiplication:  $\partial_x[f(x)g(x)] = f(x)[\partial_x g(x)] + [\partial_x f(x)]g(x)$
- Function division:  $\partial_x \left[ \frac{f(x)}{g(x)} \right] = \frac{[\partial_x f(x)]g(x) - f(x)[\partial_x g(x)]}{[g(x)]^2}$
- Function composition:  $\partial_x[f(g(x))] = [\partial_x g(x)][\partial_x f](g(x))$
- Exponentiation:  $\partial_x[e^x] = e^x$  and  $\partial_x[a^x] = \log(a)e^x$
- Logarithms:  $\partial_x[\log x] = \frac{1}{x}$

# Review: Some important rules for taking gradient

- $\frac{\partial \mathbf{x}^T \mathbf{a}}{\partial \mathbf{x}} = \frac{\partial \mathbf{a}^T \mathbf{x}}{\partial \mathbf{x}} = \mathbf{a}$
- $\nabla_x x^T A x = 2Ax$  (if  $A$  symmetric)
- $\nabla_x^2 x^T A x = 2A$  (if  $A$  symmetric)

$\Rightarrow J(\theta)$  quadratic func of  $\theta$ ; if 1-d,  $\bigcup \frac{\partial J(\theta)}{\partial \theta} = 0$

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

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

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

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

$$\Rightarrow \nabla_{\theta} J(\theta) = X^T X \theta - X^T y \stackrel{\text{Set to}}{=} 0$$

$$\Rightarrow X^T X \theta = X^T y$$

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

under certain condition

## Comments on the normal equation

- In most situations of practical interest, the number of data points  $N$  is larger than the dimensionality  $p$  of the input space and the matrix  $\mathbf{X}$  is of full column rank. If this condition holds, then it is easy to verify that  $X^T X$  is necessarily invertible.
- The assumption that  $X^T X$  is invertible implies that it is positive definite, thus the critical point we have found is a minimum.
- What if  $\mathbf{X}$  has less than full column rank?  $\rightarrow$  regularization (later).

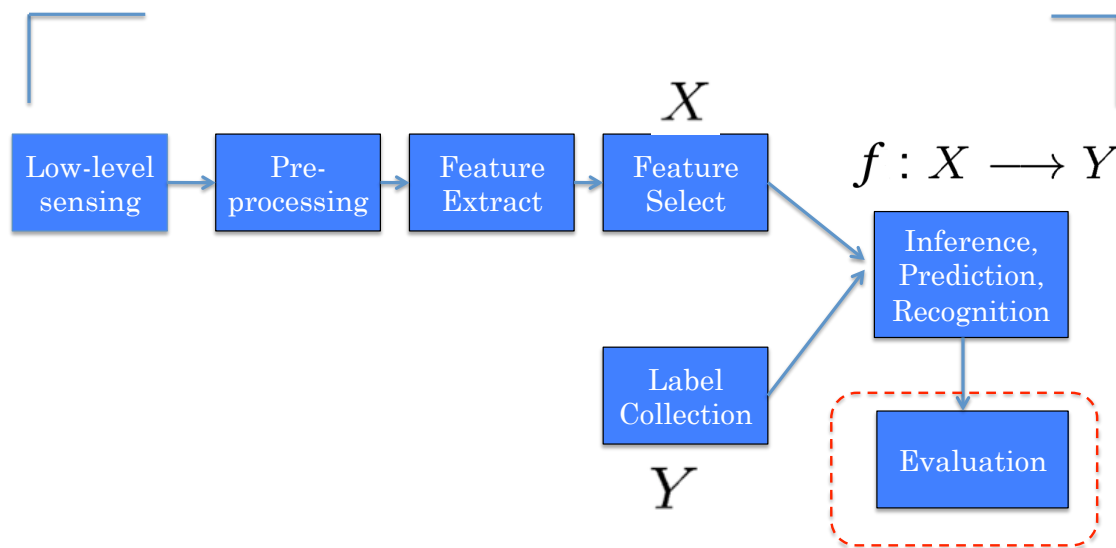
# Today

- ❑ Linear regression (aka **least squares**)
- ❑ Learn to derive the least squares estimate by optimization
- ❑ Evaluation with Train/Test OR k-folds Cross-validation

9/14/15

31

## TYPICAL MACHINE LEARNING SYSTEM

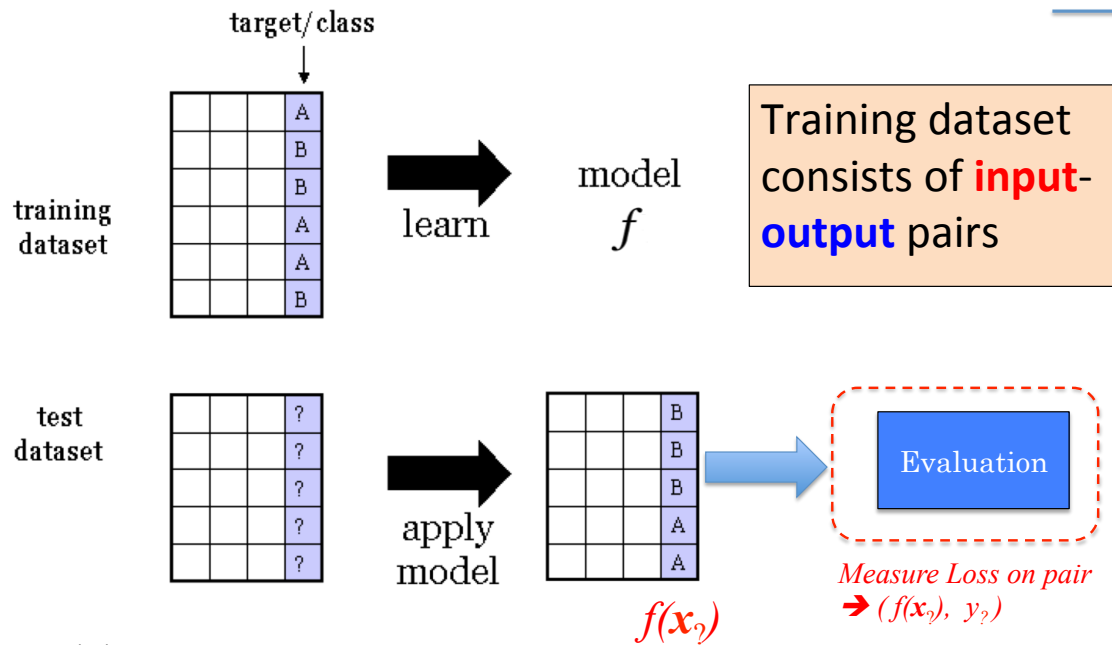


9/14/15

32



# Evaluation Choice-I: Train and Test

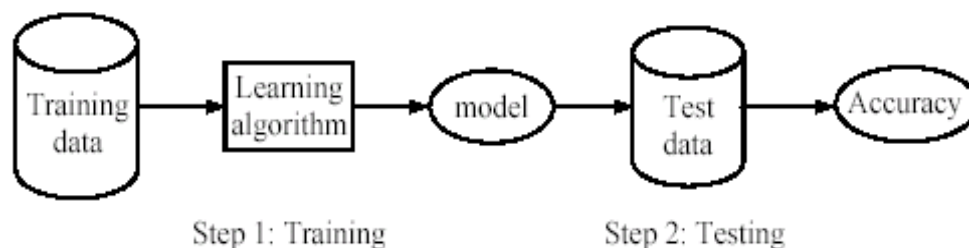


9/14/15

33

## Evaluation Choice-I: e.g. for supervised classification

- ✓ **Training (Learning):** Learn a model using the training data
- ✓ **Testing:** Test the model using **unseen test data** to assess the model accuracy

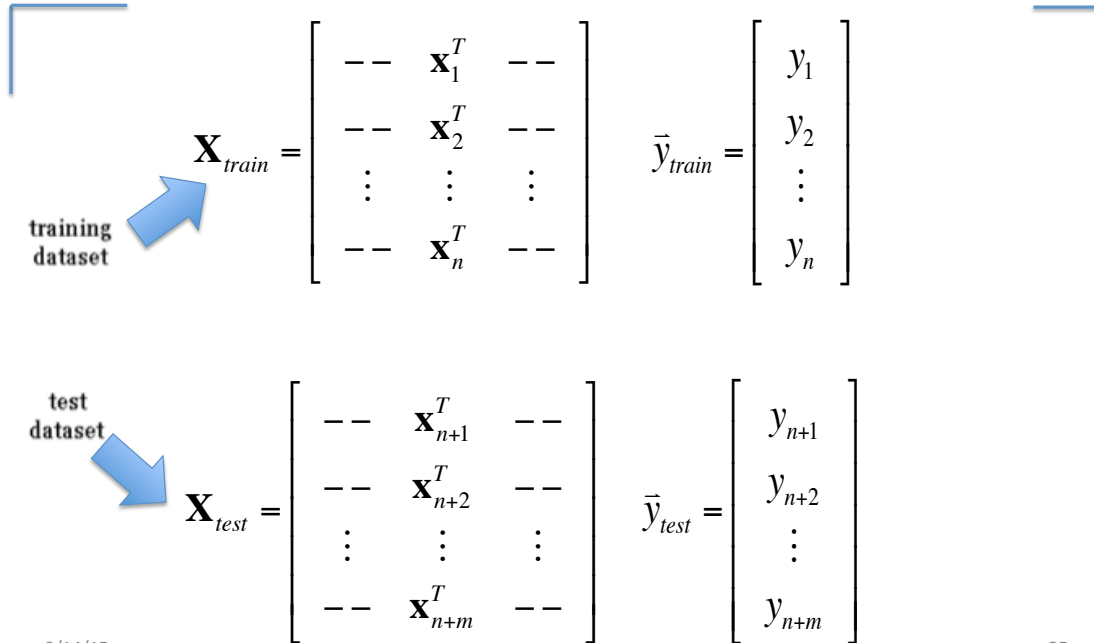


$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}},$$

9/14/15

## Evaluation Choice-I:

e.g. for linear regression models



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

$$\mathbf{X}_{test} = \begin{bmatrix} -- & \mathbf{x}_{n+1}^T & -- \\ -- & \mathbf{x}_{n+2}^T & -- \\ \vdots & \vdots & \vdots \\ -- & \mathbf{x}_{n+m}^T & -- \end{bmatrix} \quad \bar{\mathbf{y}}_{test} = \begin{bmatrix} y_{n+1} \\ y_{n+2} \\ \vdots \\ y_{n+m} \end{bmatrix}$$

9/14/15

35

## Evaluation Choice-I:

e.g. for linear regression models

- Training Error:

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

- Minimize  $J_{train}(\theta) \rightarrow$  Normal Equation to get

$$\theta^* = \operatorname{argmin} J_{train}(\theta) = \left( \mathbf{X}_{train}^T \mathbf{X}_{train} \right)^{-1} \mathbf{X}_{train}^T \bar{\mathbf{y}}_{train}$$

9/14/15

36

## Evaluation Choice-I: e.g. for Regression Models

- Testing MSE Error to report:

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

## Evaluation Choice-II: Cross Validation

- Problem: don't have enough data to set aside a test set
- Solution: Each data point is used both as train and test
- Common types:
  - K-fold cross-validation (e.g. K=5, K=10)
  - 2-fold cross-validation
  - Leave-one-out cross-validation (LOOCV, i.e., k=n\_reference)

# K-fold Cross Validation

- Basic idea:
  - Split the whole data to N pieces;
  - N-1 pieces for fit model; 1 for test;
  - Cycle through all N cases;
  - K=10 “folds” a common rule of thumb.
- The advantage:
  - all pieces are used for both training and validation;
  - each observation is used for validation exactly once.

9/14/15

39

## e.g. 10 fold Cross Validation

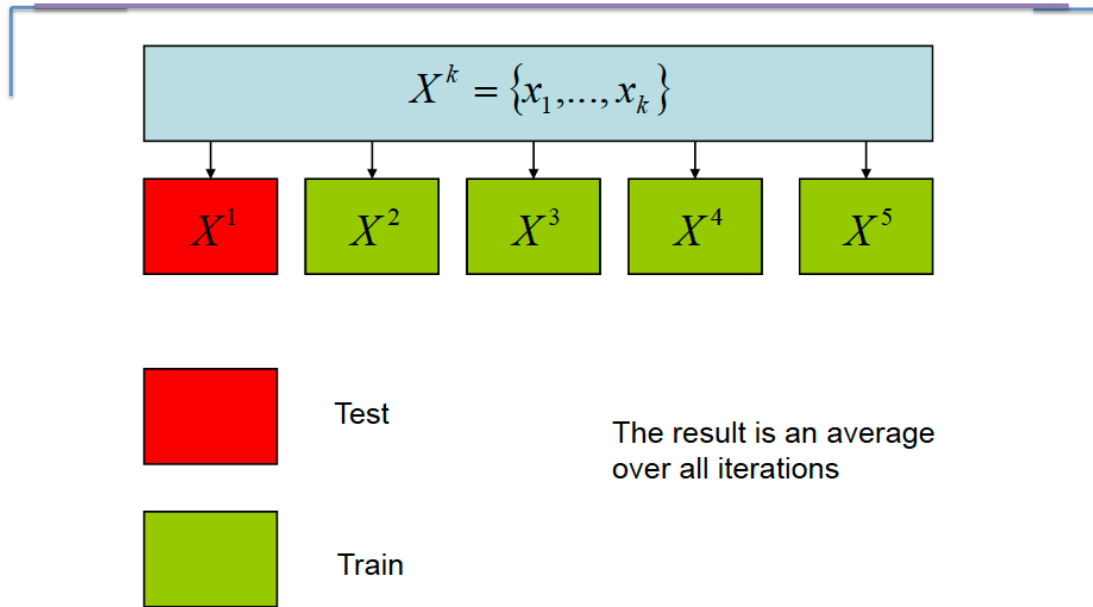
- Divide data into 10 equal pieces
- 9 pieces as training set, the rest 1 as test set
- Collect the scores from the diagonal
- We normally use the mean of the scores

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

9/14/15

40

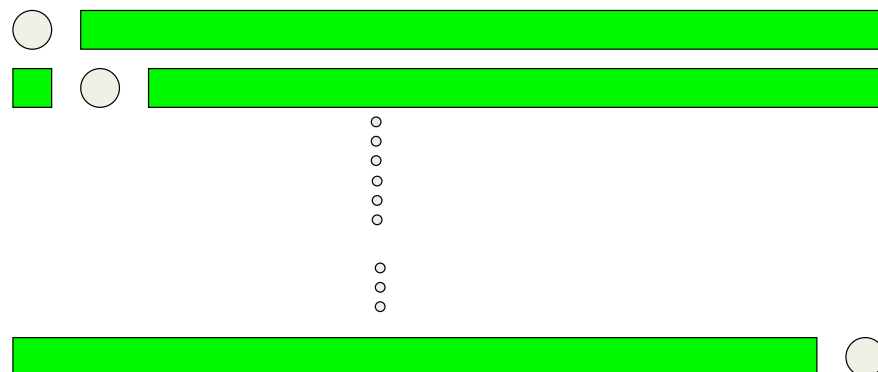
## e.g. 5 fold Cross Validation



9/14/15

41

## e.g. Leave-one-out / LOOCV (n-fold cross validation)



9/14/15

42

## Today Recap

- ❑ Linear regression (aka **least squares**)
- ❑ Learn to derive the least squares estimate by normal equation
- ❑ Evaluation with Train/Test OR k-folds Cross-validation

## References

- Big thanks to Prof. Eric Xing @ CMU for allowing me to reuse some of his slides
- ❑ <http://www.cs.cmu.edu/~zkolter/course/15-884/linalg-review.pdf> (please read)
- ❑ [http://www.cs.cmu.edu/~aarti/Class/10701/recitation/LinearAlgebra\\_Matlab\\_Review.ppt](http://www.cs.cmu.edu/~aarti/Class/10701/recitation/LinearAlgebra_Matlab_Review.ppt)
- ❑ Prof. Alexander Gray's slides