# UVA CS 4501 - 001 / 6501 – 007

# Introduction to Machine Learning and Data Mining

# Lecture 3: Linear Regression

Yanjun Qi / Jane

University of Virginia
Department of
Computer Science

9/2/14

1

# Last Lecture Recap

❑ **Data Representation**
❑ **Linear Algebra Review**

9/2/14

2

# e.g. SUPERVISED LEARNING

$$f : X \longrightarrow Y$$

- Find function to map input space X to output space Y

- So that the difference between *y* and *f(x)* of each example *x* is small.

9/2/14                                                                 3

---

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

$s_1$

$s_2$

$s_3$

$s_4$

$s_5$
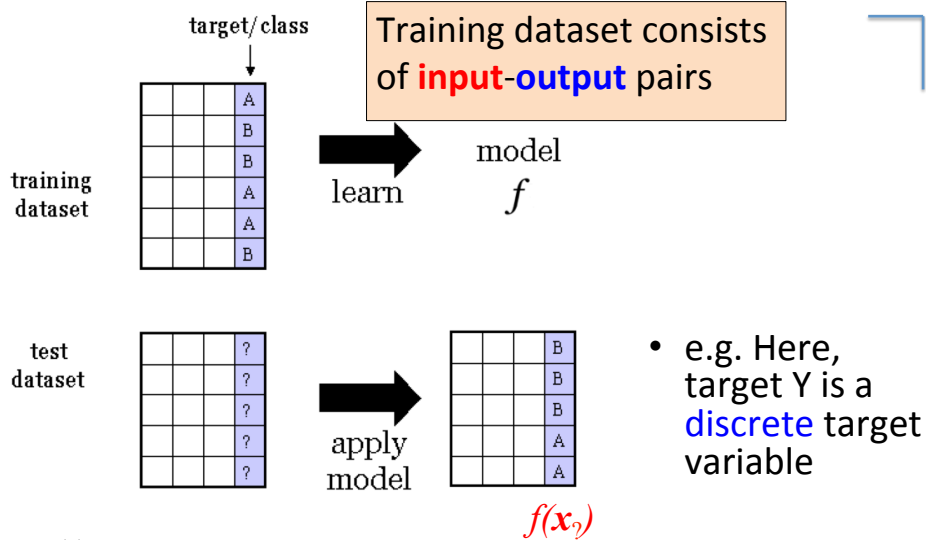
$s_6$

## A Dataset

$$f : X \longrightarrow Y$$

- **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 ]

9/2/14                                                                 4

# Main Types of Columns

|   | $X_1$ | $X_2$ | $X_3$ | $Y$ |
|---|---|---|---|---|
| $S_1$ | | | | |
| $S_2$ | | | | |
| $S_3$ | | | | |
| $S_4$ | | | | |
| $S_5$ | | | | |
| $S_6$ | | | | |

- *Continuous*: a real number, for example, age or height

- *Discrete*: a symbol, like "Good" or "Bad"

9/2/14                                                                 5

# e.g. SUPERVISED LEARNING

target/class

Training dataset consists of **input**-**output** pairs

training dataset

| | | | A |
| | | | B |
| | | | B |
| | | | A |
| | | | A |
| | | | B |

learn → model $f$

test dataset

| | | | ? |
| | | | ? |
| | | | ? |
| | | | ? |
| | | | ? |

apply model →

| | | | B |
| | | | B |
| | | | B |
| | | | A |
| | | | A |

$f(x_?)$

- e.g. Here, target Y is a discrete target variable

9/2/14                                                                 6

3

# MATRIX OPERATIONS

1) Transposition
2) Addition and Subtraction
3) Multiplication
4) Norm (of vector)
5) Matrix Inversion
6) Matrix Rank
7) Matrix calculus

9/2/14

# **Today**

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

9/2/14

8

---

# 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 (ft²) | # bedroom | Rent ($) |
|---|---|---|
| 230 | 1 | 600 |
| 506 | 2 | 1000 |
| 433 | 2 | 1100 |
| 109 | 1 | 500 |
| … | | |
| 150 | 1 | ? |
| 270 | 1.5 | ? |

9

---

# For Example,
## Machine learning for apartment hunting

| Living area (ft²) | # bedroom | Rent ($) |
|---|---|---|
| 230 | 1 | 600 |
| 506 | 2 | 1000 |
| 433 | 2 | 1100 |
| 109 | 1 | 500 |
| … | | |
| 150 | 1 | ? |
| 270 | 1.5 | ? |

features labels

$X_1$ $X_2$ $Y$

$S_1$
$S_2$
$S_3$
$S_4$
$S_5$
$S_6$

10

5

# Linear SUPERVISED LEARNING

$$f: X \longrightarrow Y$$

e.g.

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

Features:

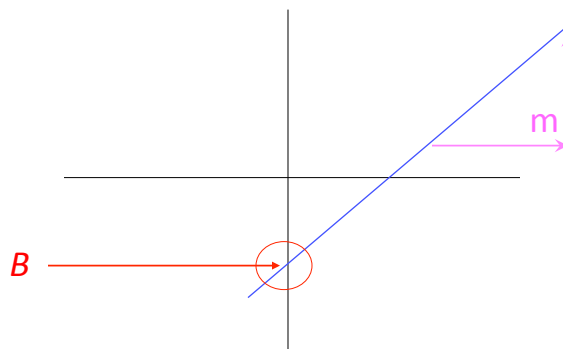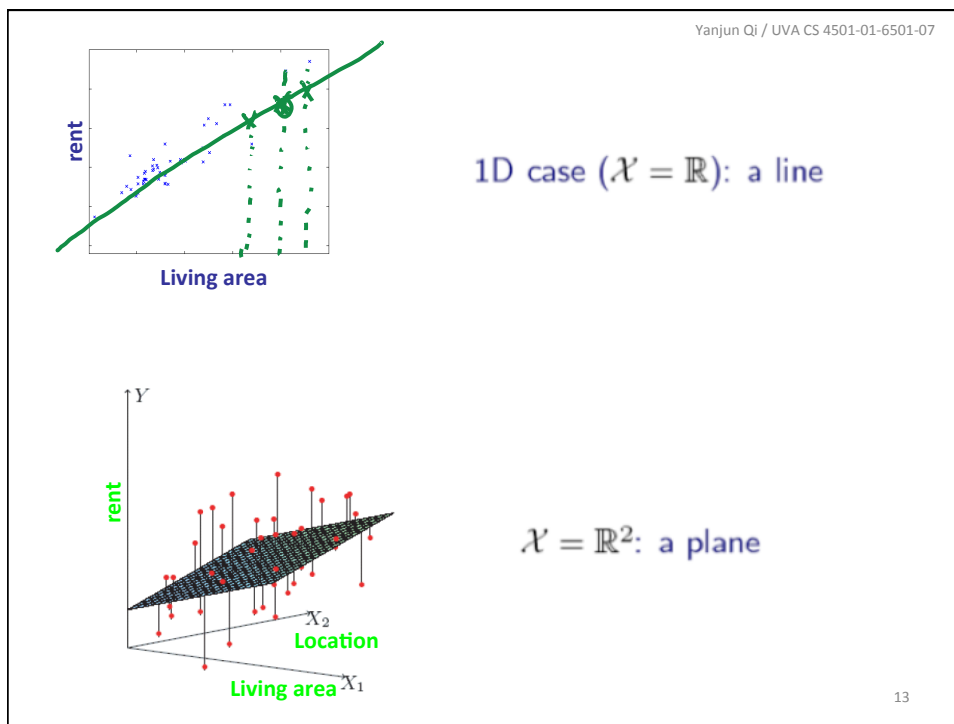Living area, distance to campus, # bedroom …

Target:

Rent

9/2/14

11

# Remember this: "Linear"?

- *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.

m

B

# A new representation

- Assume that each sample **x** is a column vector,
  - Here we assume a vacuous "feature" $X^0$=1 (this is the intercept term ), and define the feature vector to be:

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

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_{p-1} \end{bmatrix} \implies \hat{y} = f(x) = \mathbf{x}^T \theta$$

9/2/14

14

7

# Training / learning problem

- We can represent the whole Training set:

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

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

- Predicted output for each training sample:

$$\begin{bmatrix} f(\mathbf{x}_1^T) \\ f(\mathbf{x}_2^T) \\ \vdots \\ f(\mathbf{x}_n^T) \end{bmatrix} = \begin{bmatrix} \mathbf{x}_1^T \theta \\ \mathbf{x}_2^T \theta \\ \vdots \\ \mathbf{x}_n^T \theta \end{bmatrix} = X\theta$$

9/2/14

15

# training / learning goal

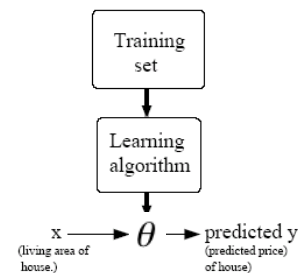- Using matrix form, we get the following general representation of the linear function:

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

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

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

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

9/2/14

16

8

# Today

❑ Linear regression (aka **least squares**)

❑ Learn to derive the least squares estimate by optimization

❑ Evaluation with Cross-validation

---

# Method I: normal equations

• Write the cost function in matrix form:

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

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

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

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

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

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

**The normal equations**

$$\Downarrow$$

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

# Review: Special Uses for Matrix Multiplication

- Dot (or Inner) Product of two Vectors <x, y>

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

    <x, y> = <y, x>

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

9/2/14                                                                                                          19

# Review: Special Uses for Matrix Multiplication

- Sum the Squared Elements of a Vector
    - Premultiply a column vector **a** by its transpose
        - If
            $$\mathbf{a} = \begin{bmatrix} 5 \\ 2 \\ 8 \end{bmatrix}$$

        then premultiplication by a row vector **a**$^T$

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

        will yield the sum of the squared values of elements for **a**, i.e.

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

9/2/14                                                                                                          20

# Review: Matrix Calculus:
# Types of Matrix Derivatives

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

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

9/2/14                                                                                              21

---

Details for slide [18] :

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

$$= (X\theta - y)^T (X\theta - y)$$

$$\underset{n\times p}{} \quad \underset{p\times 1}{} \quad \underset{n\times 1}{}$$

Since    $W^T W = \|W\|_2^2 = \sum_{i=1}^{n} W_i^2$

$$\Updownarrow$$

$$X\theta$$

9/2/14                                                                                              22

11

$$J(\theta) = (X\theta - y)^T (X\theta - y)$$

$$= ((X\theta)^T - y^T)(X\theta - y)$$

$$= (\theta^T X^T - y^T)(X\theta - y)$$

$$= \theta^T X^T X \theta - \underbrace{\theta^T X^T y - y^T X \theta}_{} - y^T y$$

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

$\langle X\theta, y \rangle \qquad \langle y, X\theta \rangle$

$$= \underbrace{\theta^T X^T X \theta}_{} - 2 \underbrace{\theta^T X^T}_{} y - y^T y$$

$\Rightarrow \quad J(\theta) \quad$ quadratic func of $\theta$; if 1-d,

$\dfrac{\partial J|\theta|}{\partial \theta} = 0$

23

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

$$\nabla_\theta (\theta^T X^T X \theta) = 2 X^T X \theta \qquad (P_{24})$$

$$\nabla_\theta (-2 \theta^T X^T y) = -2 X^T Y \qquad (P_{24})$$

$$\nabla_\theta (-y^T y) = 0$$

$$\Rightarrow \quad \nabla_\theta J(\theta) = 2 X^T X \theta - 2 X^T Y \xupparrow{\text{Set to}} 0$$

$$\Rightarrow \quad X^T X \theta = X^T Y$$

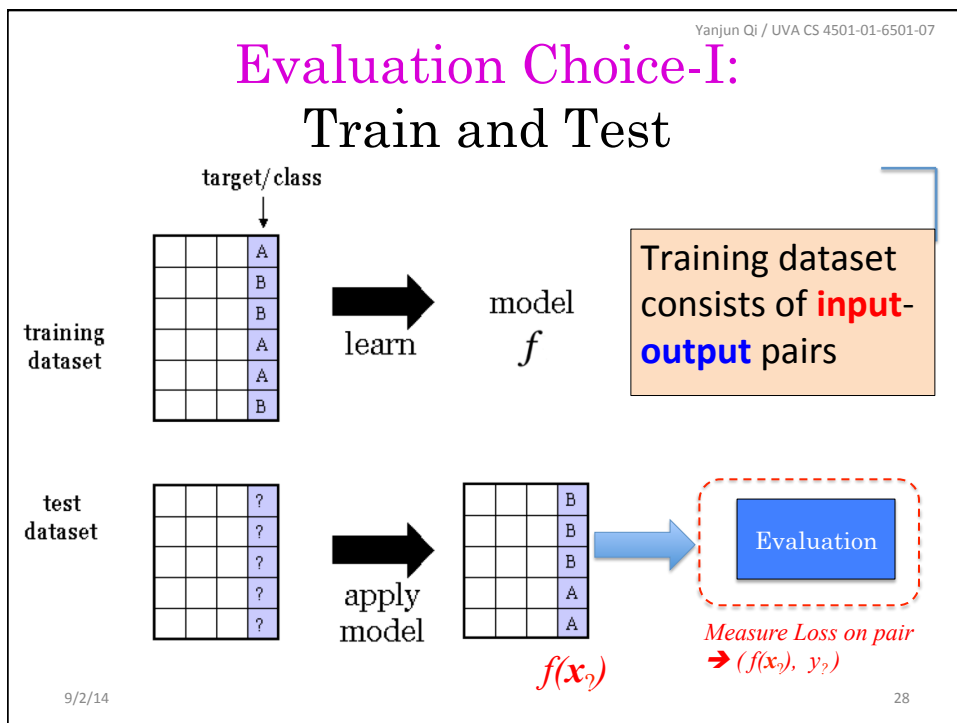$$\Rightarrow \quad \theta = (X^T X)^{-1} X^T Y$$

under certain condition

24

# 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 $X^TX$ is necessarily invertible.

- The assumption that $X^TX$ 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? → regularization (later).

9/2/14

25

# Today

❑ Linear regression (aka **least squares**)

❑ Learn to derive the least squares estimate by optimization

❑ Evaluation with Train/Test OR k-folds Cross-validation

9/2/14

26

# TYPICAL MACHINE LEARNING SYSTEM



$$X$$

| Low-level sensing | → | Pre-processing | → | Feature Extract | → | Feature Select |

$$f : X \longrightarrow Y$$

Inference, Prediction, Recognition

Label Collection

$$Y$$

Evaluation

9/2/14

27

---

# Evaluation Choice-I:
# Train and Test



target/class

training dataset

| | | A |
| | | B |
| | | B |
| | | A |
| | | A |
| | | B |

learn → model $f$

Training dataset consists of **input-output** pairs

test dataset

| | | ? |
| | | ? |
| | | ? |
| | | ? |
| | | ? |

apply model

| | | B |
| | | B |
| | | B |
| | | A |
| | | A |

$f(x_?)$

Evaluation

*Measure Loss on pair*
➔ ( f(x_?), y_? )

9/2/14

28

14

# 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



Step 1: Training          Step 2: Testing

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

9/2/14

---

# 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

9/2/14                                                                    30

---

# 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.

9/2/14        31

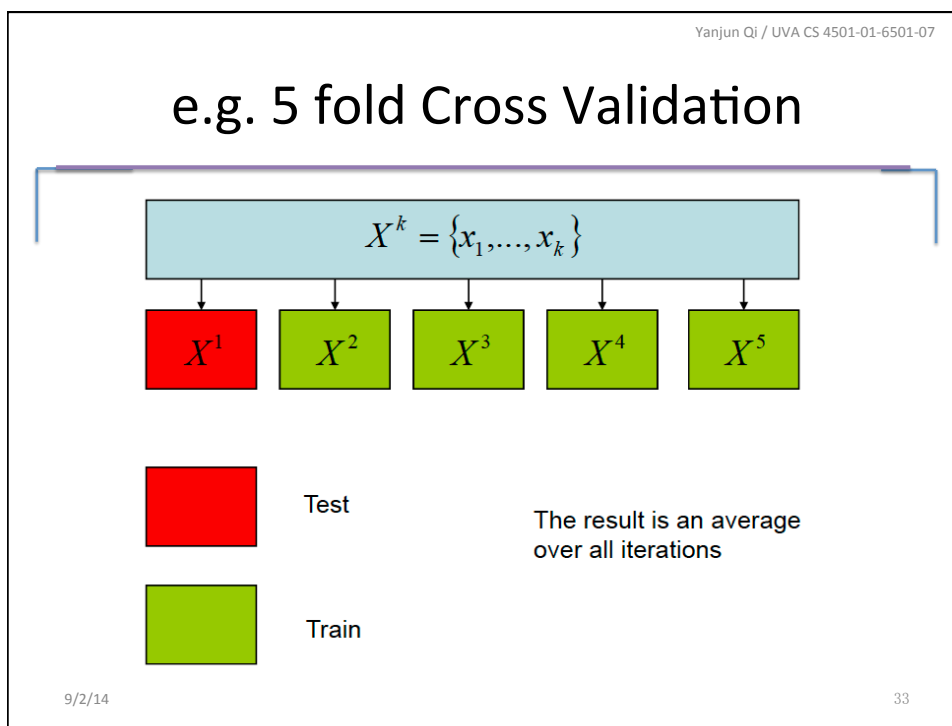---

# 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 |

9/2/14        32

# e.g. 5 fold Cross Validation

$$X^k = \{x_1, \ldots, x_k\}$$

$X^1$ $X^2$ $X^3$ $X^4$ $X^5$

Test

The result is an average over all iterations

Train

# Today Recap

❑ Linear regression (aka **least squares**)

❑ Learn to derive the least squares estimate by optimization

❑ 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_Rev iew.ppt
- ❑ Prof. Alexander Gray's slides

9/2/14                                                                                   35

18