Georgia Tech

# Linear Regression

Mahdi Roozbahani
Georgia Tech

These slides are inspired based on slides from Le Song, Chao Zhang, Yaser Abu-Mostafa, Andrew Zisserman.

# Outline

- Supervised Learning ⬅
- Linear Regression
- Extension

$X_{n \times d}$  $Y_{n \times 1} \rightarrow$ actual

# Supervised Learning: Overview

Functions $\mathcal{F}$

$$f : \mathcal{X} \to \mathcal{Y}$$

Training data

$$\{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}\}$$

$\rightarrow$ predicted

**LEARNING**

find $\hat{f} \in \mathcal{F}$

s.t. $y_i \approx \hat{f}(x_i) = \hat{y}$

**Learning machine**

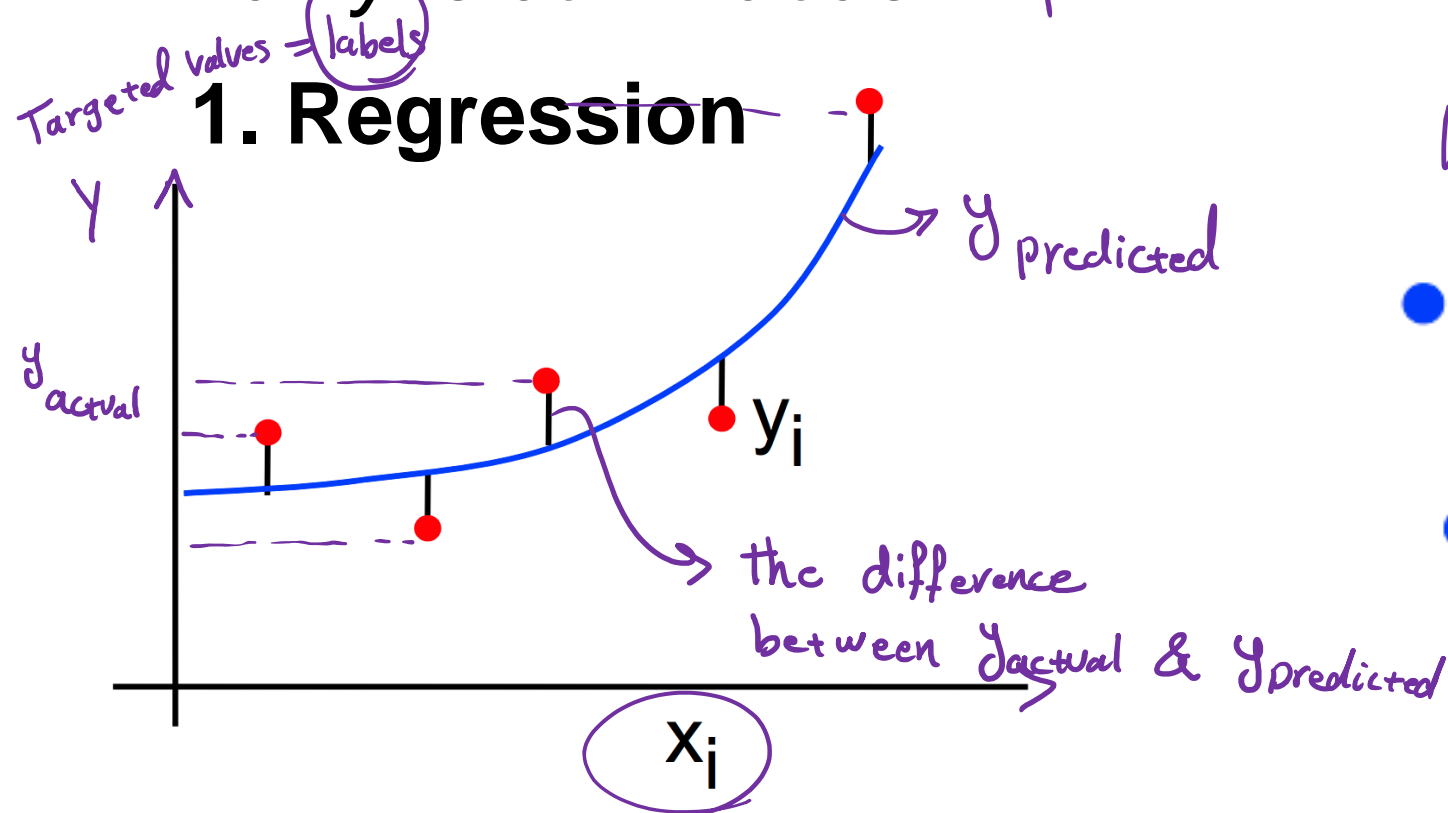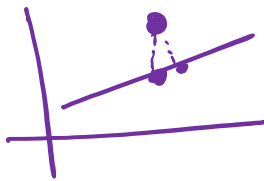**PREDICTION**

$$y = \hat{f}(x)$$

New data

$$x$$

4

# Supervised Learning: Two Types of Tasks

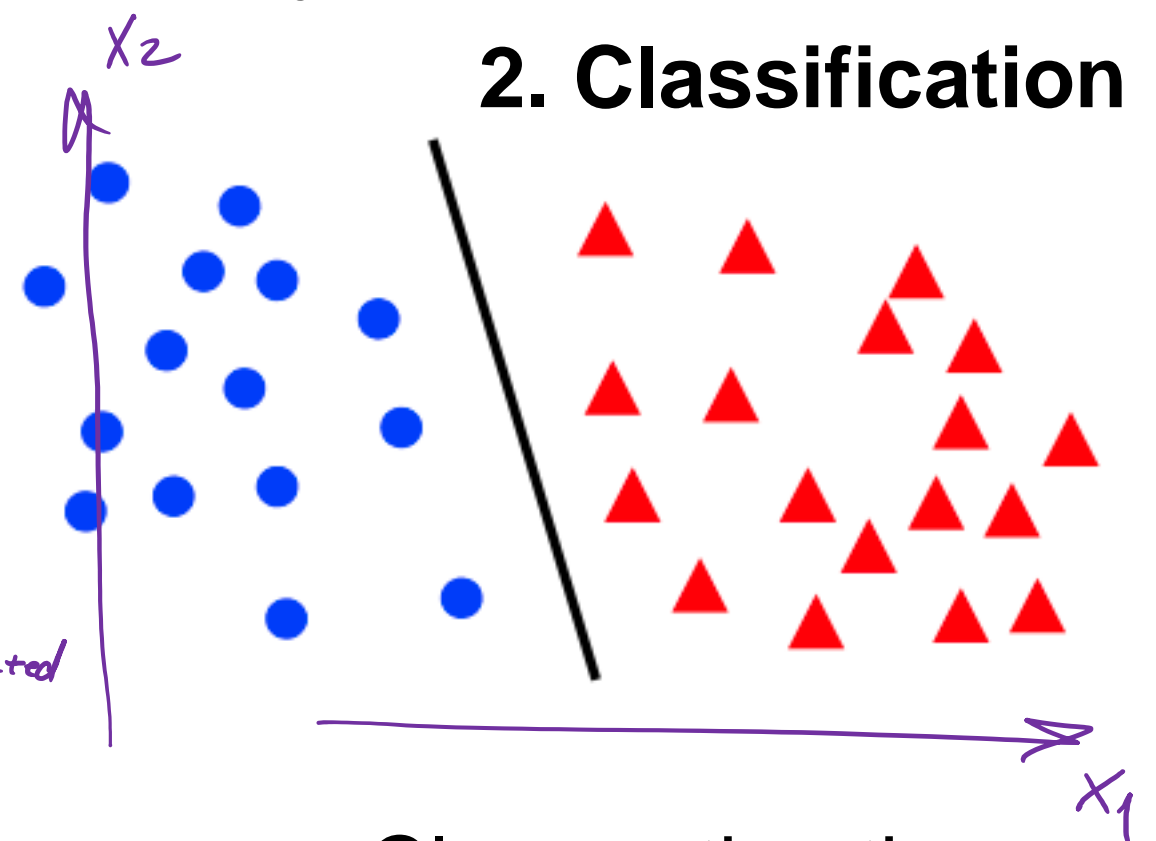**Given**: training data $\{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \ldots, (\mathbf{x}_n, y_n)\}$

**Learn**: a function $f(\mathbf{x}) : y = f(\mathbf{x})$

*When y is continuous:*

Targeted values = labels

**1. Regression**

$y_{actual}$

$y_{predicted}$

$y_i$

the difference between $y_{actual}$ & $y_{predicted}$

$x_i$

Curve fitting

*When y is discrete:*
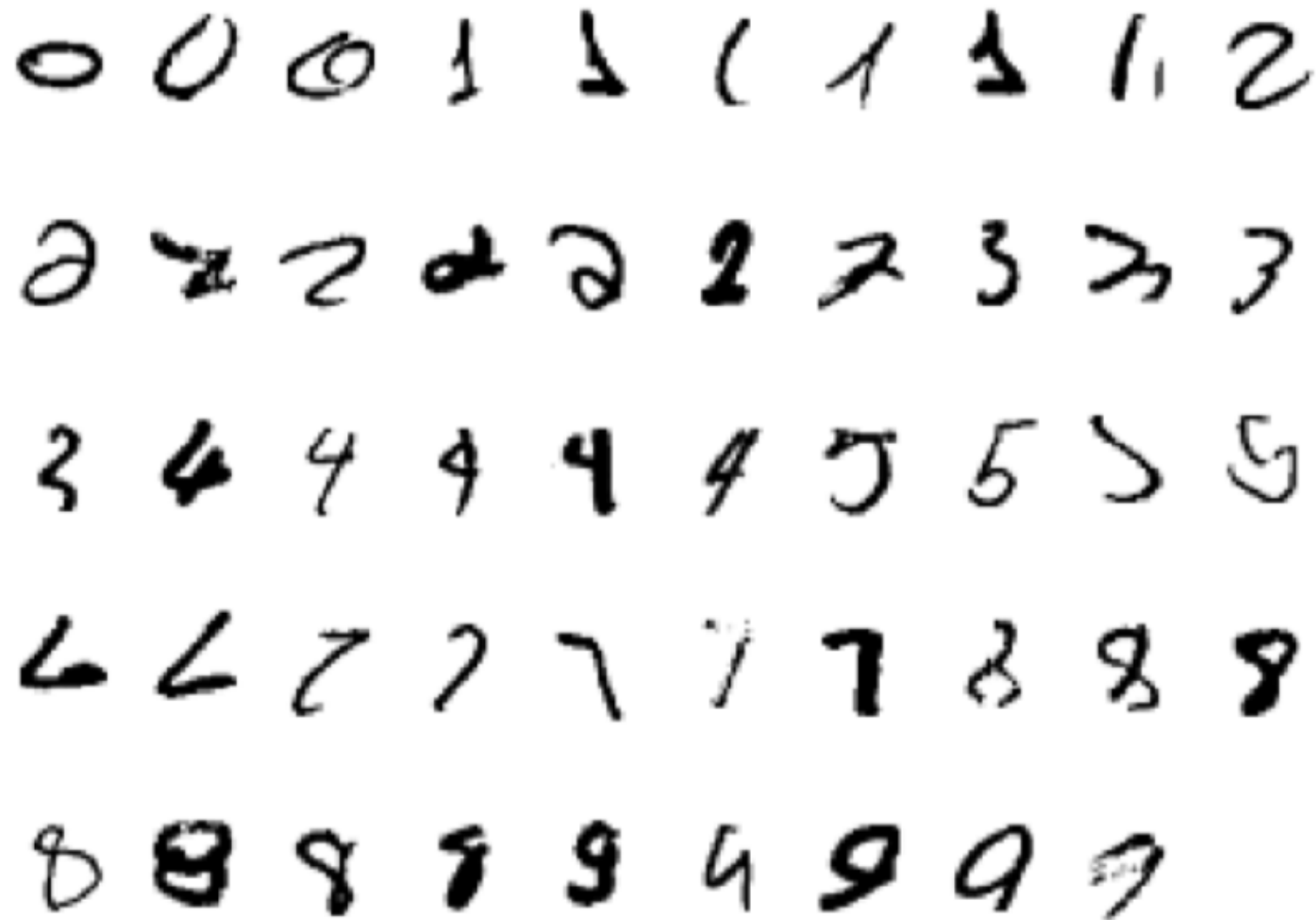
$x_2$

**2. Classification**

$x_1$

Class estimation

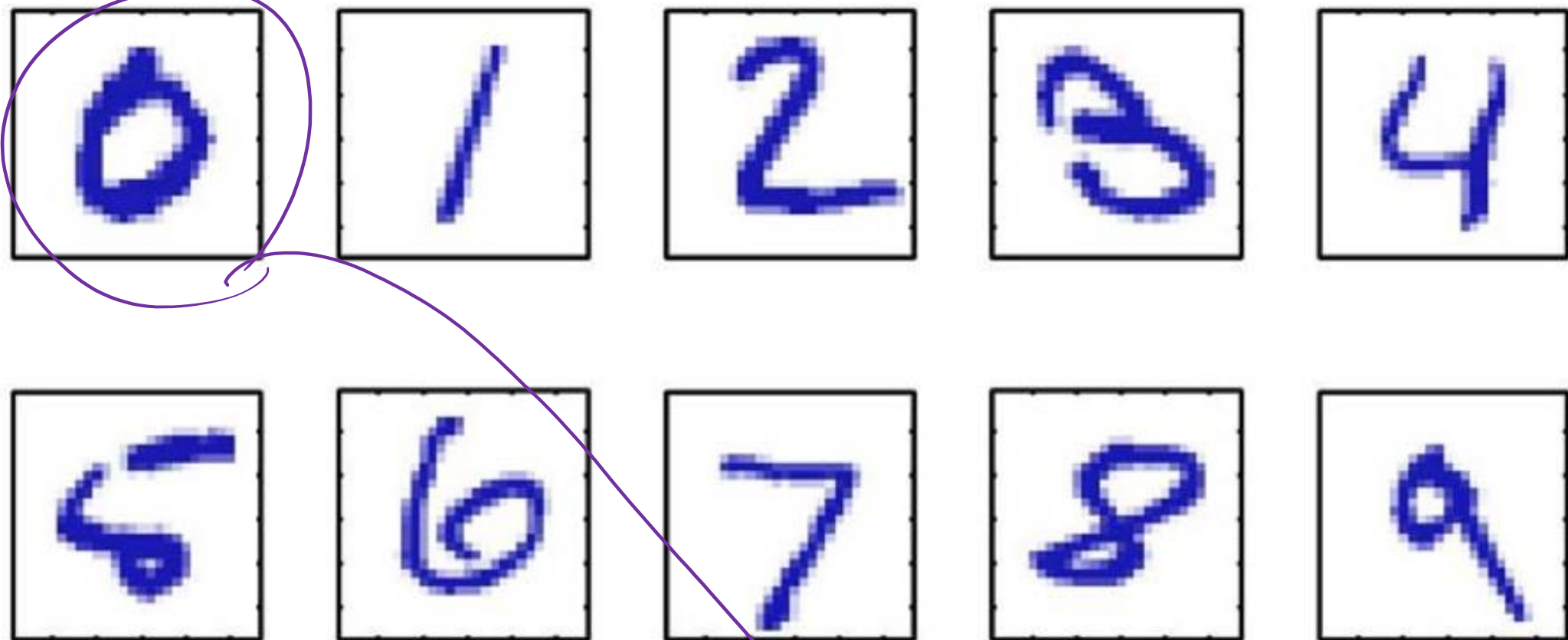# Classification Example 1: Handwritten digit recognition

As a supervised classification problem

Start with training data, e.g. 6000 examples of each digit



- Can achieve testing error of 0.4%

- One of first commercial and widely used ML systems (for zip codes & checks)

# Classification Example 1: Hand-Written Digit Recognition
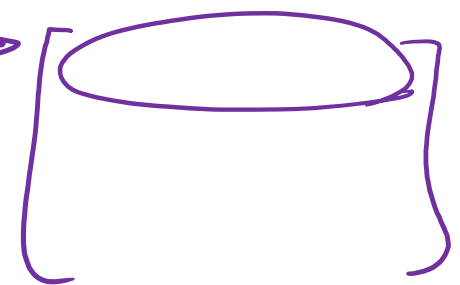
Images are 28 x 28 pixels

**A classification problem**

Represent input image as a vector $\mathbf{x} \in \mathbb{R}^{784}$

Learn a classifier $f(\mathbf{x})$ such that,

$$f : \mathbf{x} \to \{0, 1, 2, 3, 4, 5, 6, 7, 8, 9\}$$
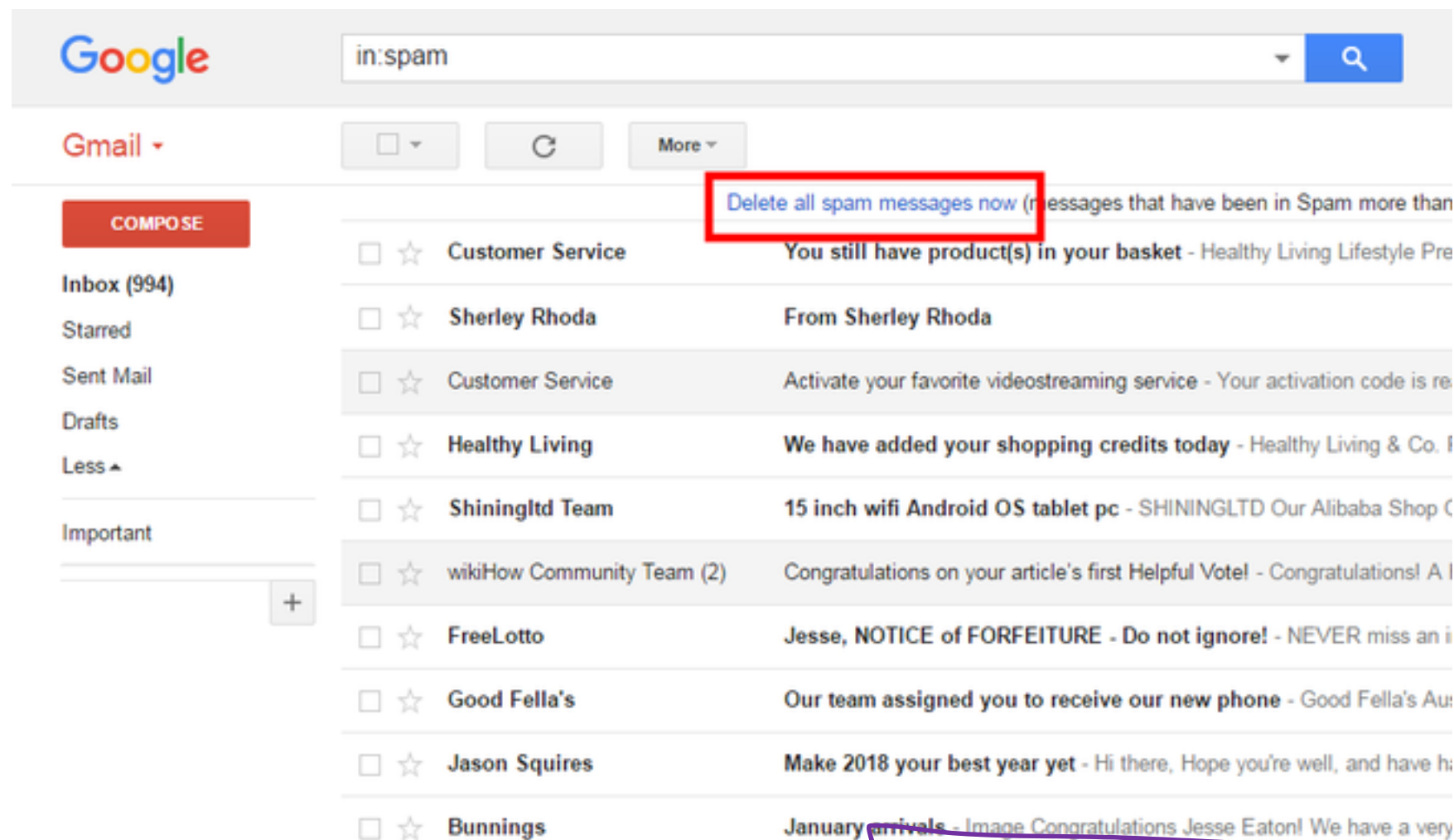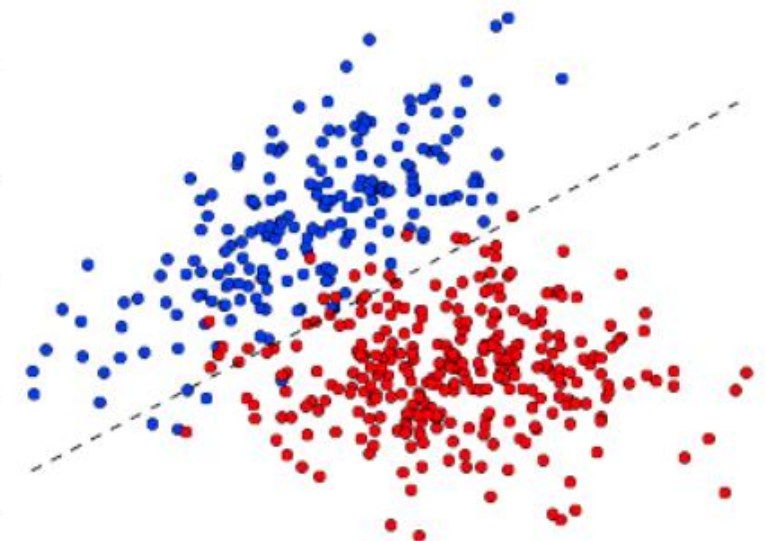
$x =$

$10 \times d$

# Classification Example 2: Spam Detection

*TF-IDF*



NOT SPAM

SPAM

**A classification problem**

- This is a classification problem

- Task is to classify email into spam/non-spam

- Data $x_i$ is word count

- Requires a learning system as "enemy" keeps innovating

"Bag of words encoding"

$$X = \begin{array}{c} \text{Email 1} \\ \text{Email 2} \\ \vdots \\ \text{Email}_n \end{array} \begin{bmatrix} 2 & 10 & 0 & 1 & 0 \\ 0 & 0 & 2 & 3 & 4 \\ & & & & \\ & & & & \end{bmatrix}$$

# Regression Example 1: Apartment Rent Prediction

- Suppose you are to move to Atlanta
- And you want to find the **most reasonably priced** apartment satisfying your **needs:**

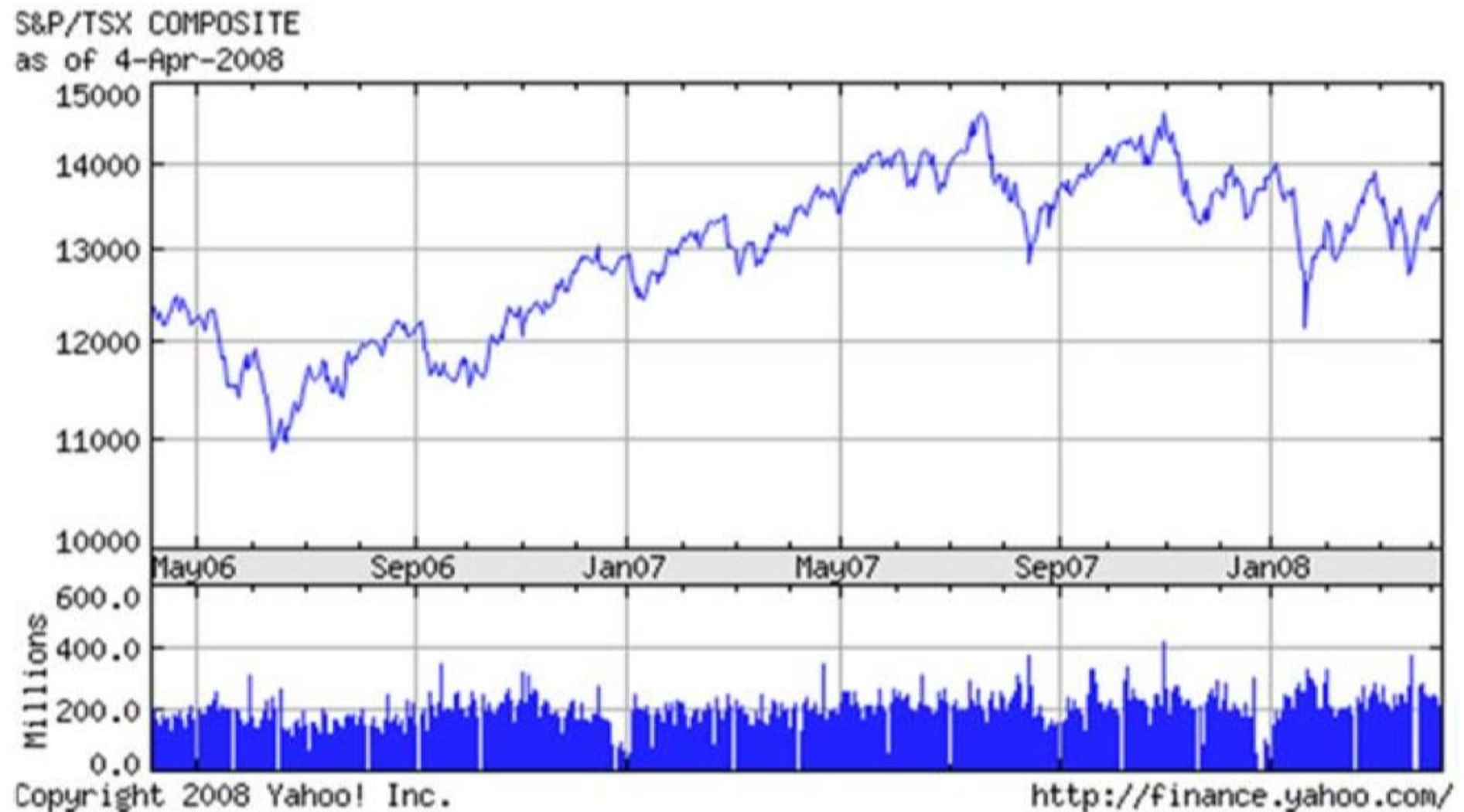    square-ft., # of bedroom, distance to campus ...

**A regression problem**

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

# Regression Example 2: Stock Price Prediction



S&P/TSX COMPOSITE
as of 4-Apr-2008

Copyright 2008 Yahoo! Inc.                    http://finance.yahoo.com/

- Task is to predict stock price at future date

**A regression problem**

slope $\hat{y} = \underbrace{m}\,x + \underbrace{b} \rightarrow$ Intercept

Living area (feature)

$\hat{y} = \Theta_0 + \Theta_1 X_1$

bias term

- **Features:**
  - Living area, distance to campus, # bedroom …
  - Denote as $x = (x_1, x_2, \dots, x_d)$

Linear combination of FEATURES

$\hat{y} \nearrow^{\mathbb{R}} = m_1 X_{Loc} + m_2 X_{Liv} + b$

- **Target:**
  - Rent

  $\hat{y} = \Theta_0 + \Theta_1 X_1 + \Theta_2 X_2 = \boxed{X\Theta}$

  - Denoted as $y$

  $X_i = [1 \quad X_1 \quad X_2 \cdots X_d]$
  
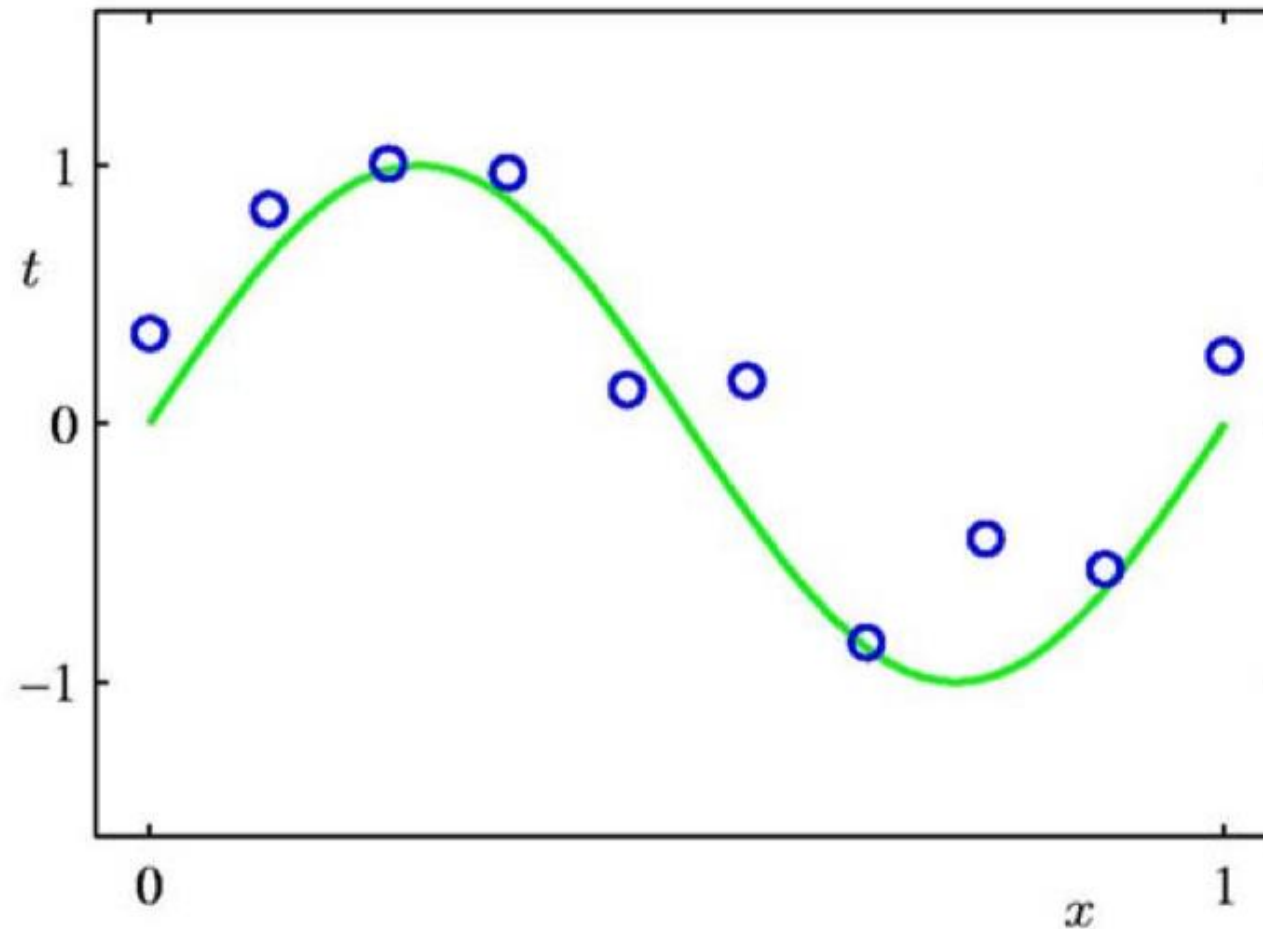  $1 \times (d+1)$

  $\Theta_{(d+1)\times 1} = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \vdots \\ \Theta_d \end{bmatrix}$

- **Training set:**
  - $x = \{x_1, x_2, \dots, x_n\} \in R^d$
  - $y = \{y_1, y_2, \dots, y_n\}$



11

# Regression: Problem Setup



- Suppose we are given a training set of N observations

$$(x_1, \ldots, x_N) \text{ and } (y_1, \ldots, y_N), x_i, y_i \in \mathbb{R}$$

- Regression problem is to estimate y(x) from this data

# Outline

- Supervised Learning
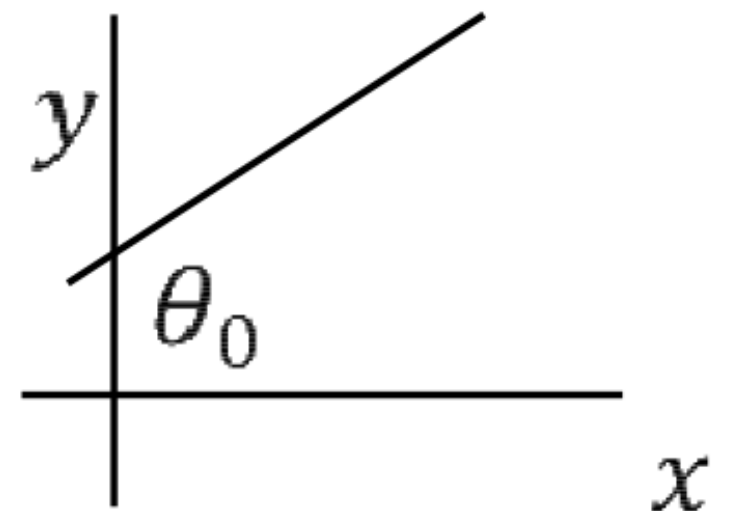
- Linear Regression ⬅

- Extension

# Linear Regression

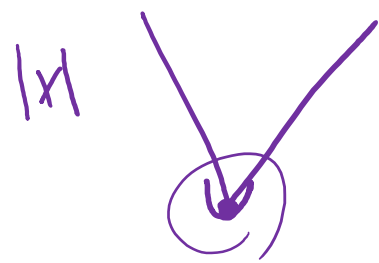- Assume $y$ is a linear function of $x$ (features) plus noise $\epsilon$

$$y = \theta_0 + \theta_1 x_1 + \cdots + \theta_d x_d + \epsilon$$

- where $\epsilon$ is an error term of unmodeled effects or <u>random noise</u>

- Let $\theta = (\theta_0, \theta_1, \ldots, \theta_d)^\top$, and augment data by one dimension

*linear combination of features*

- Then $\hat{y} = x\theta + \epsilon$  *noise*

# Least Mean Square Method

$|x|$ $x^2$ $|y_i - x_i\theta|$

- Given n data points, find $\theta$ that minimizes the mean square error

Training

$$\hat{\theta} = argmin_\theta \, L(\theta) = \frac{1}{n}\sum_{i=1}^{n}(y_i - x_i\theta)^2$$

$E(\theta)$  $y_{actual}$  $y_{predicted}$  $x_i$  $\theta$

$1\times(d+1)$  $(d+1)\times1$

$|x1|$  $|x1|$

- Our usual trick: set gradient to 0 and find parameter $\dfrac{\partial L(\theta)}{\partial \theta} = 0$

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n}\sum_{i=1}^{n} x_i^T(y_i - x_i\theta) = 0$$

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n}\sum_{i=1}^{n} x_i^T y_i + \frac{2}{n}\sum_{i=1}^{n} x_i^T x_i\theta = 0$$

# Matrix form

$$x = \begin{bmatrix} 1 & x_1^{\{1\}} & \dots & x_1^{\{d\}} \\ 1 & x_2^{\{1\}} & \ddots & x_2^{\{d\}} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_n^{\{1\}} & \dots & x_n^{\{d\}} \end{bmatrix}_{n \times (d+1)} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}_{(d+1) \times 1}$$

$$MSE(\theta) = argmin_\theta \, L(\theta) = \frac{1}{n} \overset{n \times 1}{(y - \widetilde{x\theta})}^{\mathrm{T}} (y - x\theta)$$

$$x\theta = \begin{bmatrix} \theta_0 + \theta_1 x_1^{\{1\}} + \theta_2 x_1^{\{2\}} + \dots + \theta_d x_1^{\{d\}} \\ \theta_0 + \theta_1 x_2^{\{1\}} + \theta_2 x_2^{\{2\}} + \dots + \theta_d x_2^{\{d\}} \\ \vdots \\ \theta_0 + \theta_1 x_n^{\{1\}} + \theta_2 x_n^{\{2\}} + \dots + \theta_d x_n^{\{d\}} \end{bmatrix}_{n \times 1}$$

# Matrix Version and Optimization

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n}\sum_{i=1}^{n} x_i^T y_i + \frac{2}{n}\sum_{i=1}^{n} x_i^T x_i \theta = 0$$

Let's rewrite it as:

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n}(x_1, \dots, x_n)^T (y_1, \dots, y_n) + \frac{2}{n}(x_1, \dots, x_n)^T (x_1, \dots, x_n)\theta = 0$$

$n \times (d+1)$

Define X $= (x_1, \dots, x_n)$ and y $= (y_1, \dots, y_n)$

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n}X^T y + \frac{2}{n}X^T X \theta = 0$$

$(d+1) \times 1$
$(d+1) \times n$
$n \times 1$
$(d+1) \times (d+1) \times (d+1) \times 1$
$(d+1) \times 1$

$$\Rightarrow \theta = (X^T X)^{-1} X^T y = X^+ y \qquad X^+ \text{ is the \textbf{pseudo-inverse} of } X$$

$(d+1) \times 1$
$(d+1) \times (d+1)$
$(d+1) \times n$
$n \times 1$

$$X^T X X^+ = X^T$$

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

$$X_{n \times d} \qquad n = \text{instances} \quad d = \text{dimension}$$

$$X^T X = \begin{bmatrix} d \times n \end{bmatrix} \begin{bmatrix} n \times d \end{bmatrix} = \begin{bmatrix} d \times d \end{bmatrix}$$

Not a big matrix because $n \gg d$ This matrix is invertible most of the times. If we are VERY unlucky and columns of $\mathbf{X^T X}$ are not linearly independent (it's not a full rank matrix), then it is not invertible.

# Alternative Way to Optimize

- The matrix inversion in $\bar{\theta} = (X^TX)^{-1}X^Ty$ can be very expensive to compute

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

$$\theta^{t+1} \leftarrow \theta^t - \alpha \frac{\partial L(\theta)}{\partial \theta}$$

learning step

- Gradient descent

$$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \frac{\alpha}{n}\sum_{i=1}^{n} x_i^T(y_i - x_i\theta)$$

- Stochastic gradient descent (use one data point at a time)

SGD $\leftarrow$ $$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \beta_t \times x_i^T(y_i - x_i\theta)$$

Batch $\leftarrow$ BGD

Training data

Actual labels

$\left(X\right)_{n \times (d+1)}$  $Y_{n \times 1}$  $\Theta_{(d+1) \times 1}$

$f(x) = \hat{y} \rightarrow$ predicted labels

$y_p$ is as close as possible to $y_a$

$f(x) = \hat{y} = \Theta_0 + \Theta_1 X_1 + \cdots + \Theta X = X\Theta$  Linear combination of features

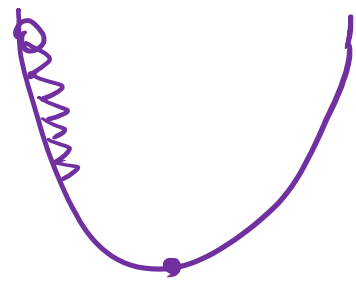$$L(\Theta) = E(\Theta) = \frac{1}{N} \sum_{i=1}^{N} (\underset{\text{actual}}{\widetilde{y_i}} - \underset{\text{predicted}}{\hat{y_i}})^2 = \frac{1}{N} \sum_{i=1}^{N} (y_i - X_i\Theta)^2$$

i.e. Convex

$$= E\left[(y_i - X_i\Theta)^2\right] = bias^2 + Variance$$

$$\underset{(d+1) \times 1}{\Theta} = \underbrace{\left(X^T X\right)^{-1}}_{(d+1) \times (d+1)} \underset{(d+1) \times n}{X^T} \underset{n \times 1}{Y}$$

$(d+1) \times 1$

① GD

$$L(\theta) = \frac{1}{N} \sum_{i=1}^{N} (y_i - x_i \theta)^2$$

i.e.

Initialize $\theta$ with random numbers or zero ($t=0$)

$$\theta^{t+1} \leftarrow \theta^t - \alpha \frac{\delta L(\theta)}{\delta(\theta)}$$

$$\frac{\delta L(\theta)}{\delta \theta} = -\frac{2}{N} \sum_{i=1}^{N} x_i^T (y_i - x_i \theta)$$

$$\theta^{t+1} \leftarrow \theta^t + \frac{\alpha}{N} \sum_{i=1}^{N} x_i^T (y_i - x_i \theta)$$

$$\theta^t \longrightarrow \theta^{t+1}$$

② SGD

$$\theta^{t+1} \leftarrow \theta^t + \beta x_i^T (y_i - x_i \theta)$$

③ BGD or Batch GD    Using a Batch of datapoints

It takes 5 iteration $\leftarrow$ to reach to one epoch

In Total, I have 100 datapoints and I use 20 datapoints in each iteration

# Recap

- Stochastic gradient update rule

$$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \beta_t \times x_i^T(y_i - x_i\theta)$$

  - Pros: on-line, low per-step cost
  - Cons: coordinate, maybe slow-converging
- Gradient descent

$$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \frac{\alpha}{n}\sum_{i=1}^{n} x_i^T(y_i - x_i\theta)$$

  - Pros: fast-converging, easy to implement
  - Cons: need to read all data
- Solve normal equations

$$\theta = (X^TX)^{-1}X^Ty$$

  - Pros: a single-shot algorithm! Easiest to implement.
  - Cons: need to compute inverse $(X^TX)^{-1}$, expensive, numerical issues (e.g., matrix is singular ..)

# Linear regression for classification

Raw Input $x = (x_0, x_1, \ldots, x_{256})$

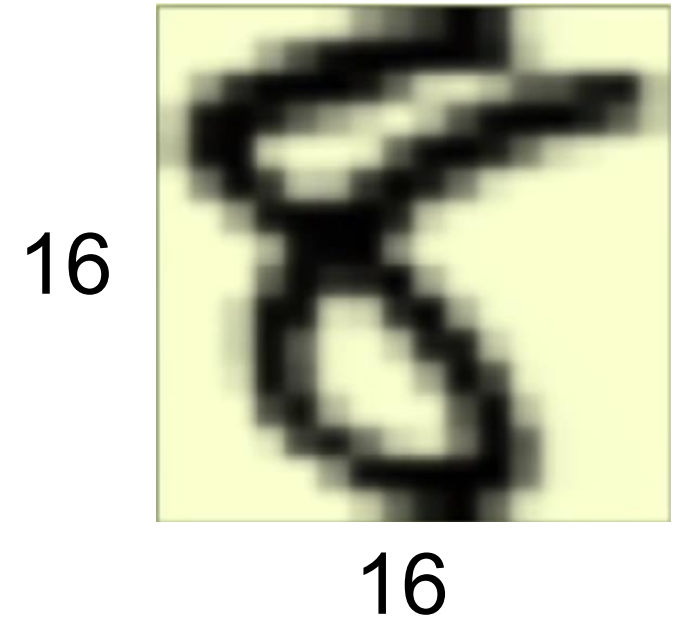Linear model $(\theta_0, \theta_1, \ldots, \theta_{256})$

<span style="color:purple">Extract useful information</span>

$intensity\ and\ symmetry\ x = (1, x_1, x_2)$

16

16

$Sum\ up\ all\ the\ pixels = intensity$

$Symmetry = \text{-}(difference\ between\ flip\ version)$

$$\hat{y} = \Theta_0 + \Theta_1 \underset{\downarrow}{X_1} + \Theta_2 \underset{\downarrow}{X_2} \rightsquigarrow \mathbb{R}$$

intensity   symmetry

$x = (1, x_1, x_2)$

$x_1 = intensity \ x_2 = symmetry$

$$\begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \end{array} \left[ \begin{array}{cc} \text{int} & \text{sym} \\ \cdots & \cdots \\ & \\ & \\ & \end{array} \right]$$

It is almost linearly separable

*symmetry*



*intensity*

# Linear regression for classification

Binary-valued functions are also real-valued $\pm 1 \in R$
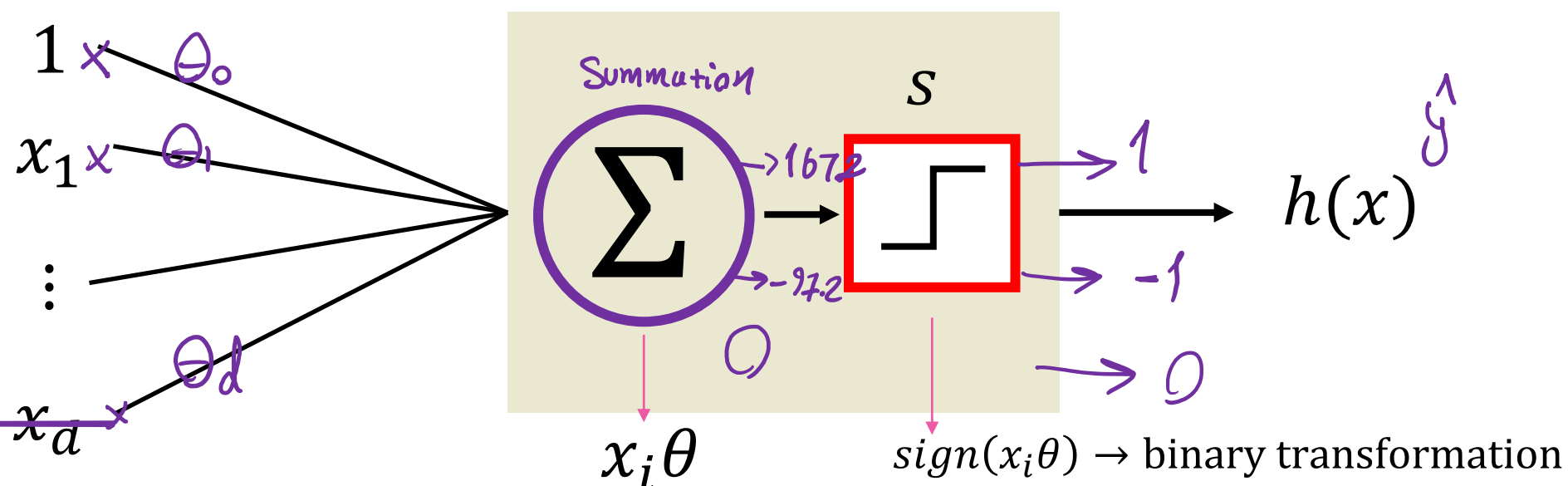
Use linear regression $x_i\theta \approx y_n = \pm 1$    $i$ = index of a data-point
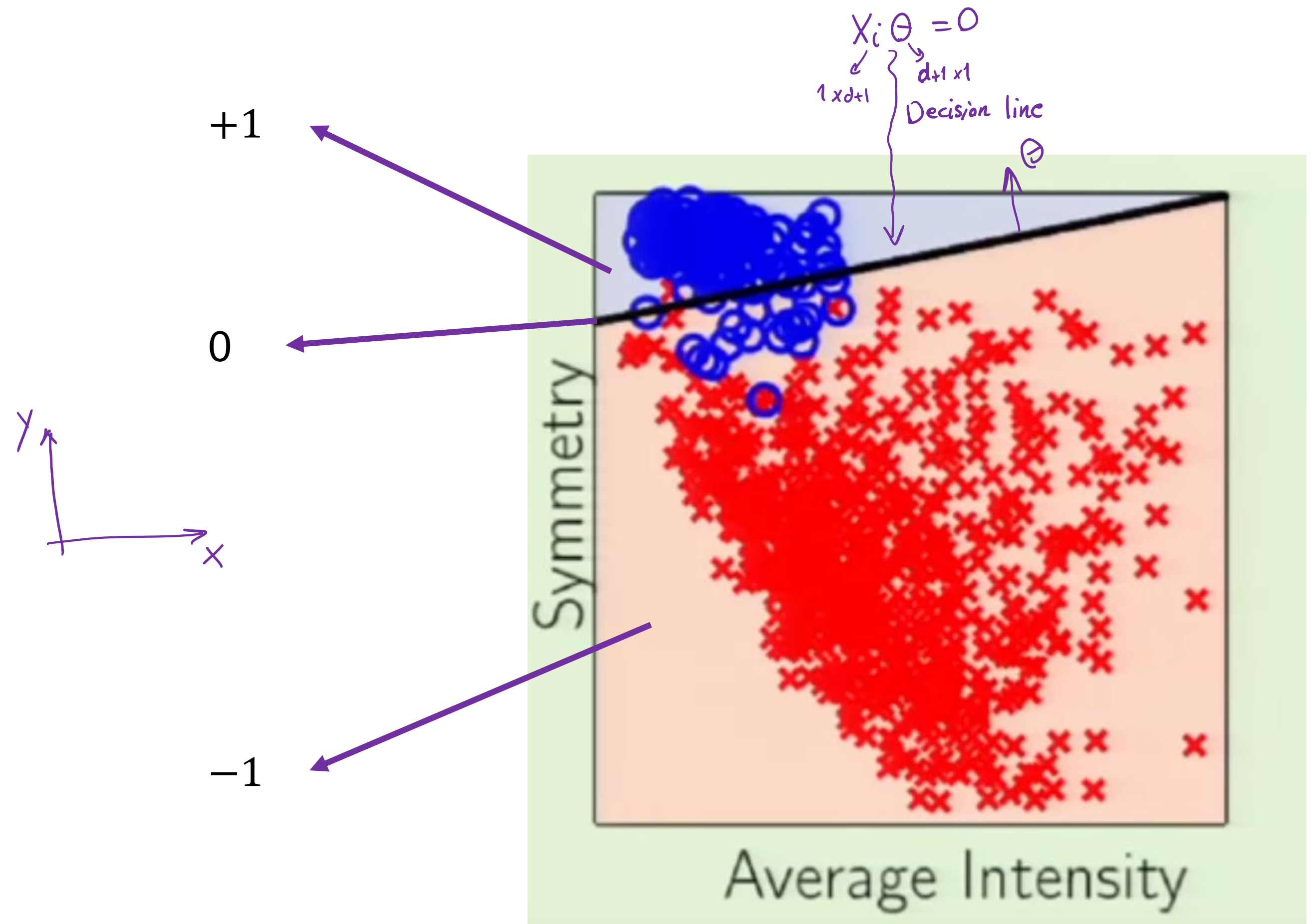
Let's calculate, $sign(x_i\theta) = \begin{cases} -1 & x_i\theta < 0 \\ 0 & x_i\theta = 0 \\ 1 & x_i\theta > 0 \end{cases}$

$\Sigma = \Theta_0 + \Theta_1 X_1 + \cdots + \Theta_d X_d$   A linear combination of features

One learning ← block

For one data point (data-point $i$) with **d** dimensions (instance):



Summation

$S$

$1 \times \quad \Theta_0$
$x_1 \times \quad \Theta_1$
$\vdots$
$x_d \times \quad \Theta_d$

$\Sigma \rightarrow 167.2 \rightarrow 1$
$\rightarrow -97.2 \rightarrow -1$
$\rightarrow 0$
$h(x)$ $\hat{y}$

$x_i\theta$         $sign(x_i\theta) \rightarrow$ binary transformation

$X_i \, \Theta = 0$

$1 \times d+1$    $d+1 \times 1$

Decision line

$\Theta$

+1

0

−1

$y$

$x$

Symmetry

Average Intensity
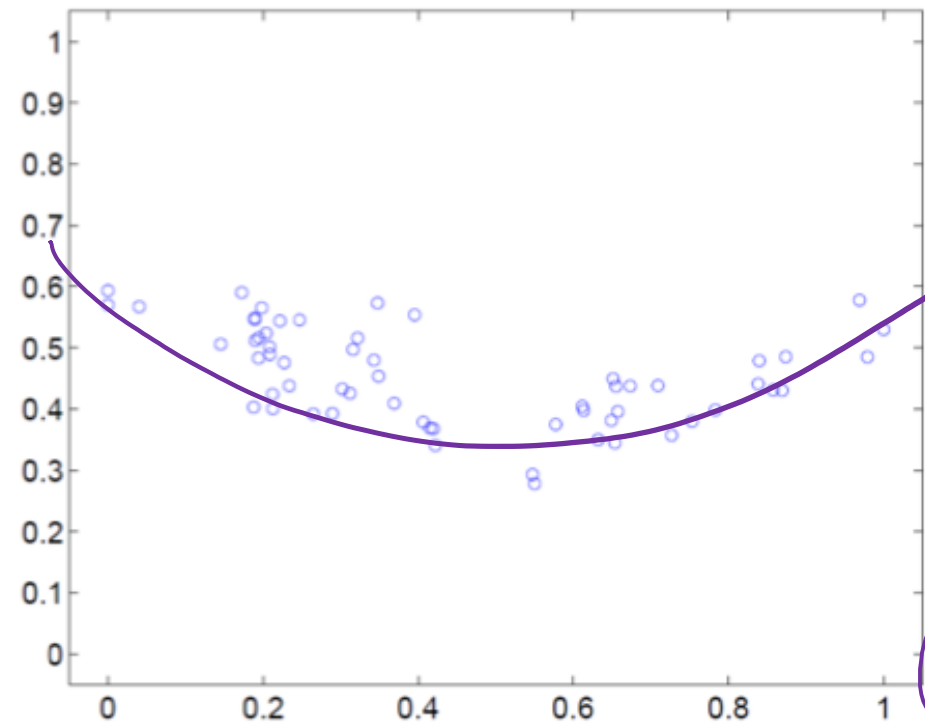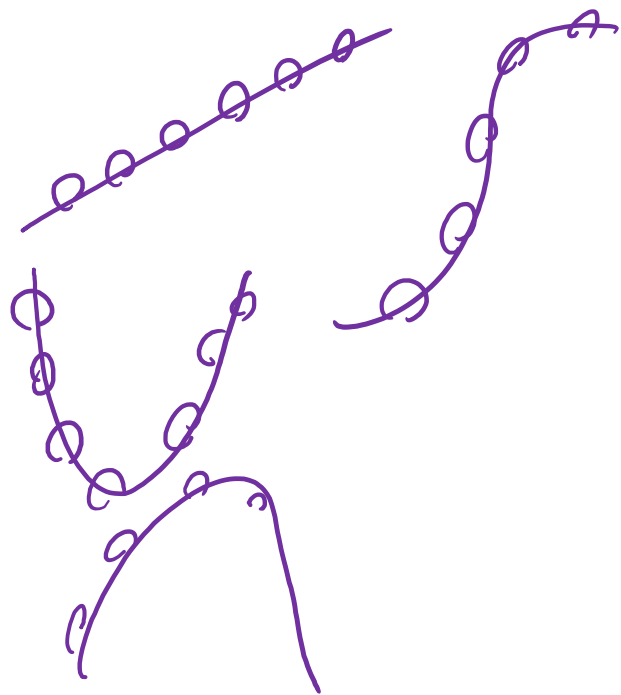
Not really the best for classification, but t's a good start

# Outline

- Supervised Learning

- Linear Regression

- Extension ←

# Extension to Higher-Order Regression



$$X = \begin{bmatrix} \\ \\ \\ \end{bmatrix} \qquad Z = \begin{bmatrix} h & h^2 & h^3 & \cdots h^d \\ \\ \\ \end{bmatrix}$$

$$X = \begin{bmatrix} h & \omega \\ \\ \\ \end{bmatrix} \qquad Z = \begin{bmatrix} h & \omega & h^2 & \omega^2 & h\omega \\ \\ h^2\omega & \omega^2h \end{bmatrix}$$

- Want to fit a polynomial regression model

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \cdots + \theta_d x^{\mathrm{d}} + \epsilon$$

$$\theta_0 + \theta_1 z_1 + \theta_2 z_2 + \cdots + \theta_d z_d$$

- $z = \{1, x, x^2, \ldots, x^d\} \in R^d$ and $\theta = (\theta_0, \theta_1, \theta_2, \ldots, \theta_{\mathrm{d}})^{\mathsf{T}}$

$$y = z\theta$$

28

# Least Mean Square Still Works the Same

- Given $n$ data points, find $\theta$ that minimizes the mean square error

$$\hat{\theta} = argmin_{\theta} \, L(\theta) = \frac{1}{n} \sum_{i=1}^{n} (y_i - z_i \theta)^2$$

- Our usual trick: set gradient to 0 and find parameter

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n} \sum_{i=1}^{n} z_i^T (y_i - z_i \theta) = 0$$

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n} \sum_{i=1}^{n} z_i^T y_i + \frac{2}{n} \sum_{i=1}^{n} z_i^T z_i \theta = 0$$

# Matrix Version of the Gradient

$$z = \{1, x, x^2, \ldots, x^d\} \in R^d \qquad\qquad y = \{y_1, y_2, \ldots, y_n\}$$

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n} z^{\mathrm{T}} y + \frac{2}{n} z^T z \theta = 0$$

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

- If we choose a different maximal degree **d** for the polynomial, the solution will be different.
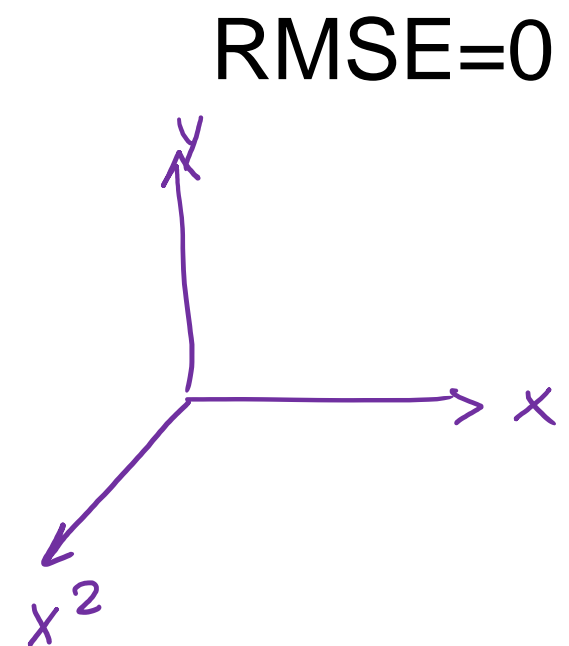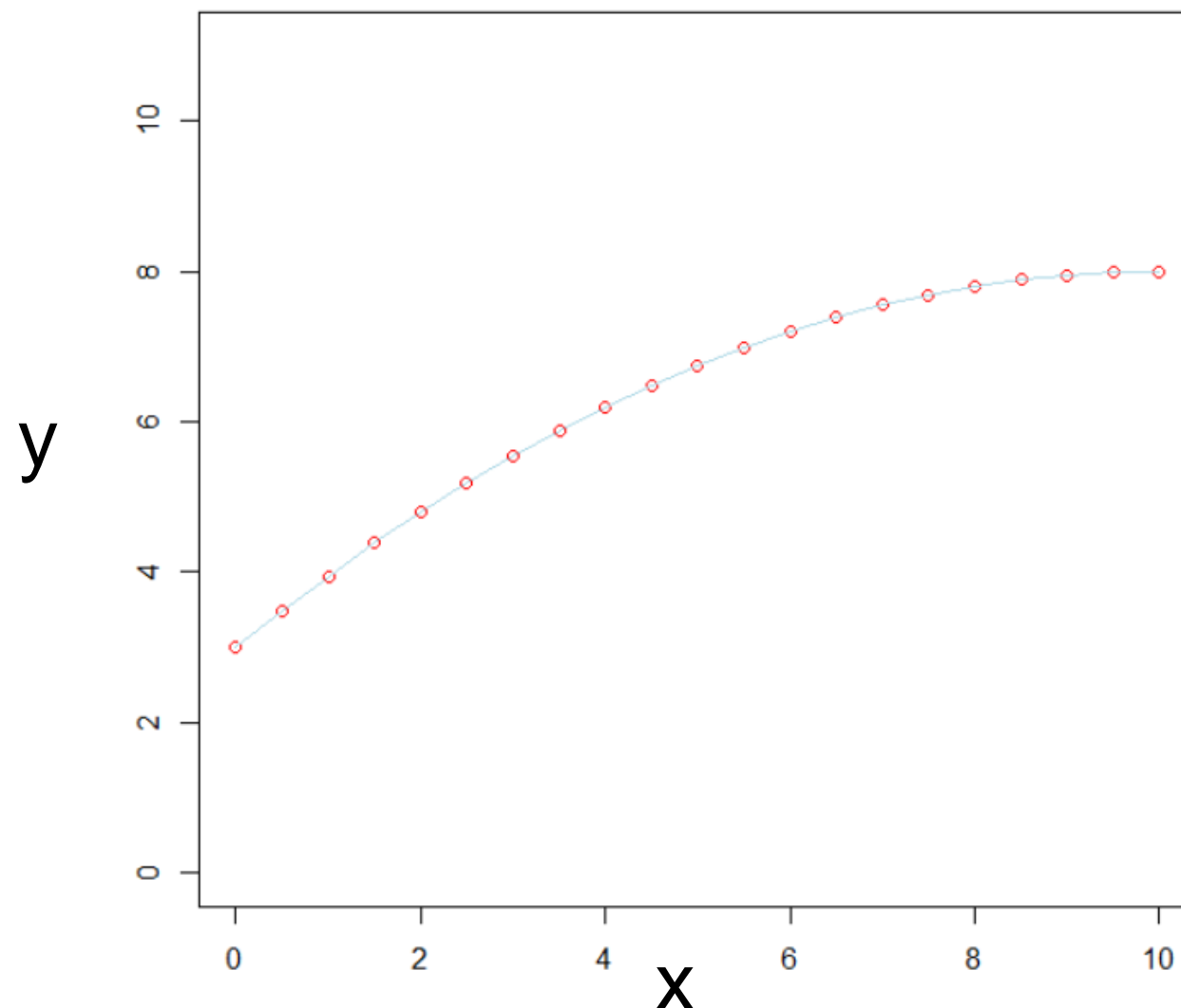
# What is happening in polynomial regression?
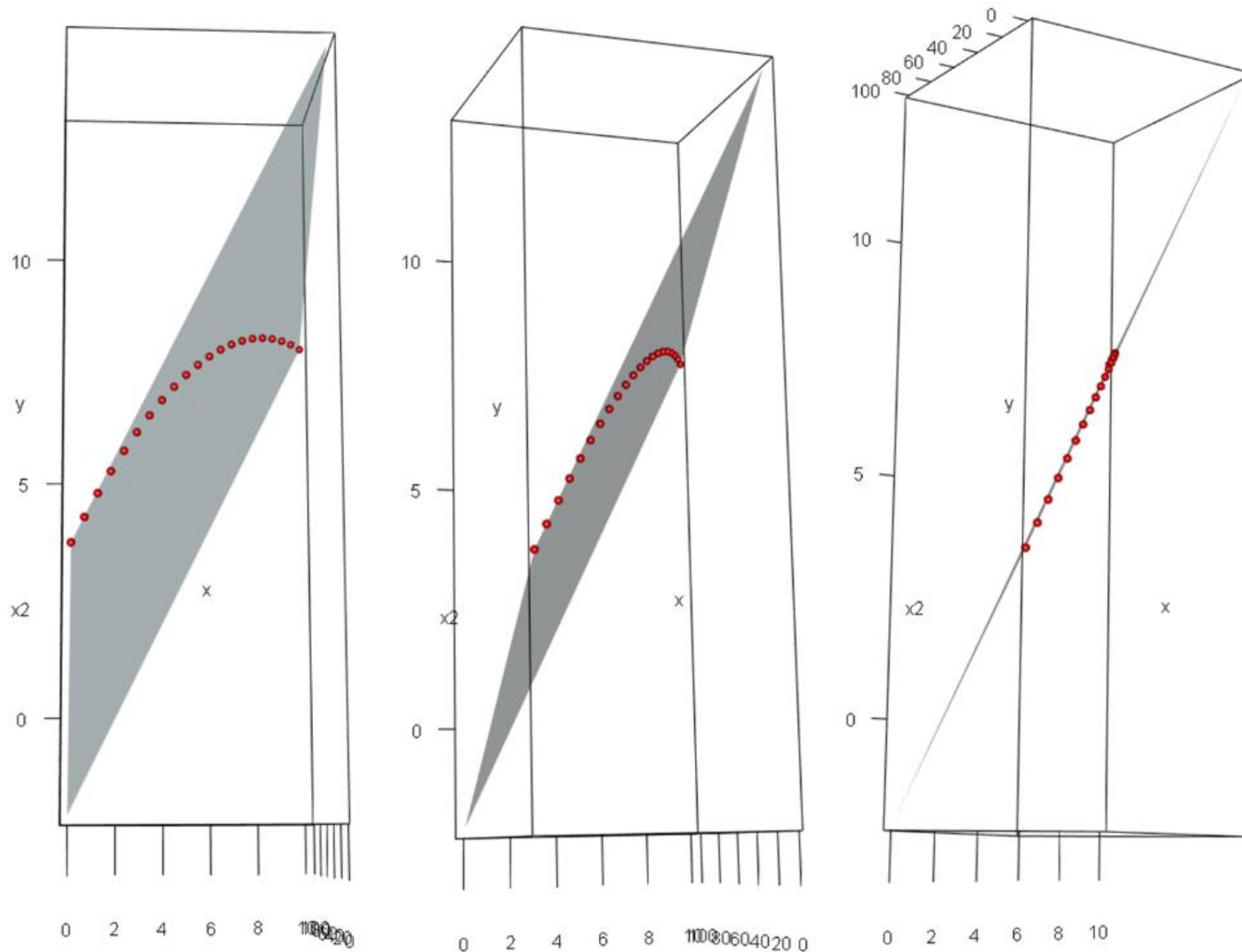
$$x = [0, 0.5, 1, \ldots, 9.5, 10]$$

$$y = [3, 3.4875, 3.95, \ldots, 7.98, 8]$$

$$f = \theta_0 + \theta_1 x + \theta_2 x^2$$

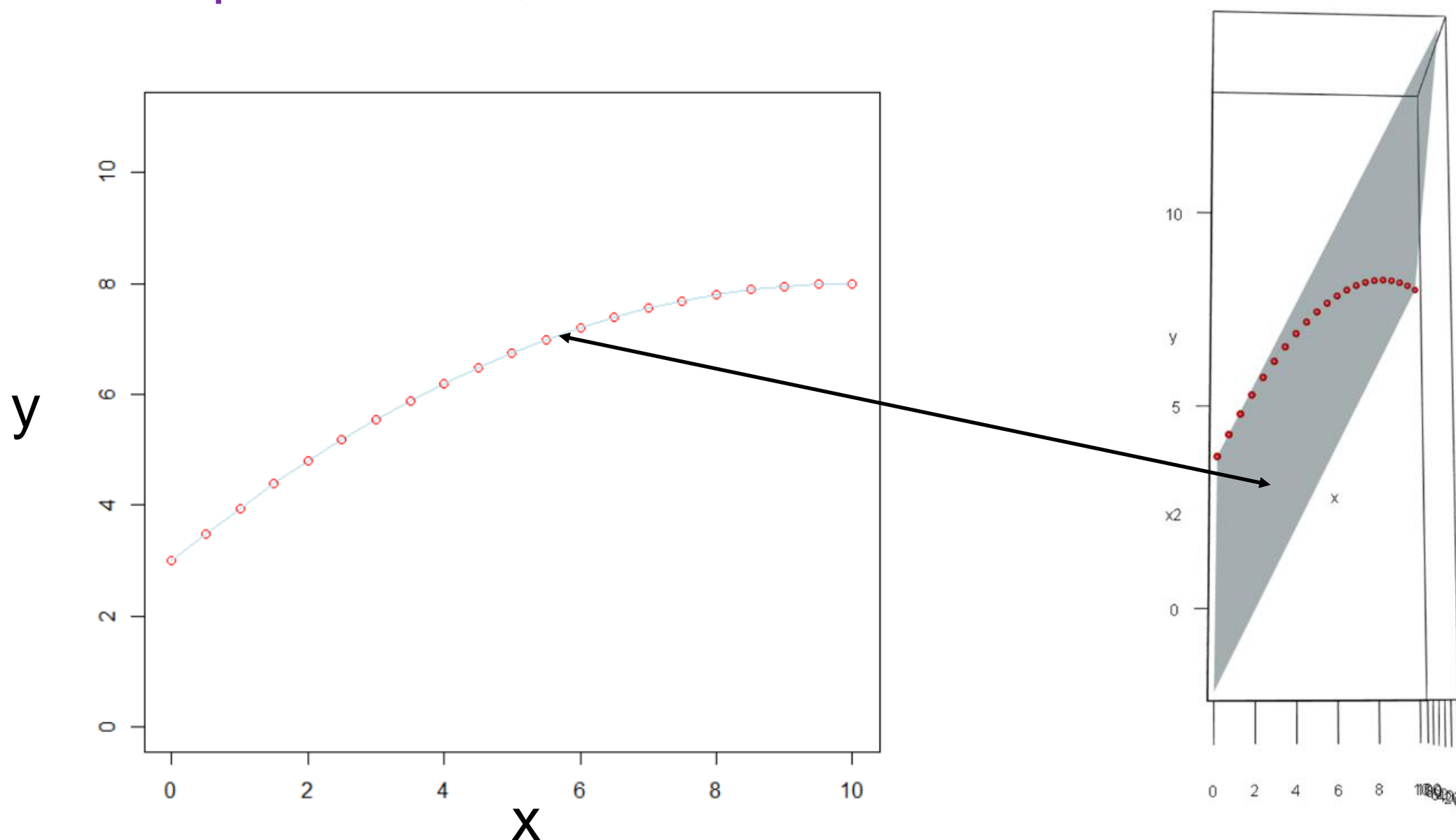$$\theta_0 = 3; \theta_1 = 1; \theta_2 = -0.5$$



RMSE=0

# Let's add to the feature space

$$x_1 = [0, 0.5, 1, \ldots, 9.5, 10] \qquad x_2^2 = [0, 0.25, 1, \ldots, 90.25, 100]$$
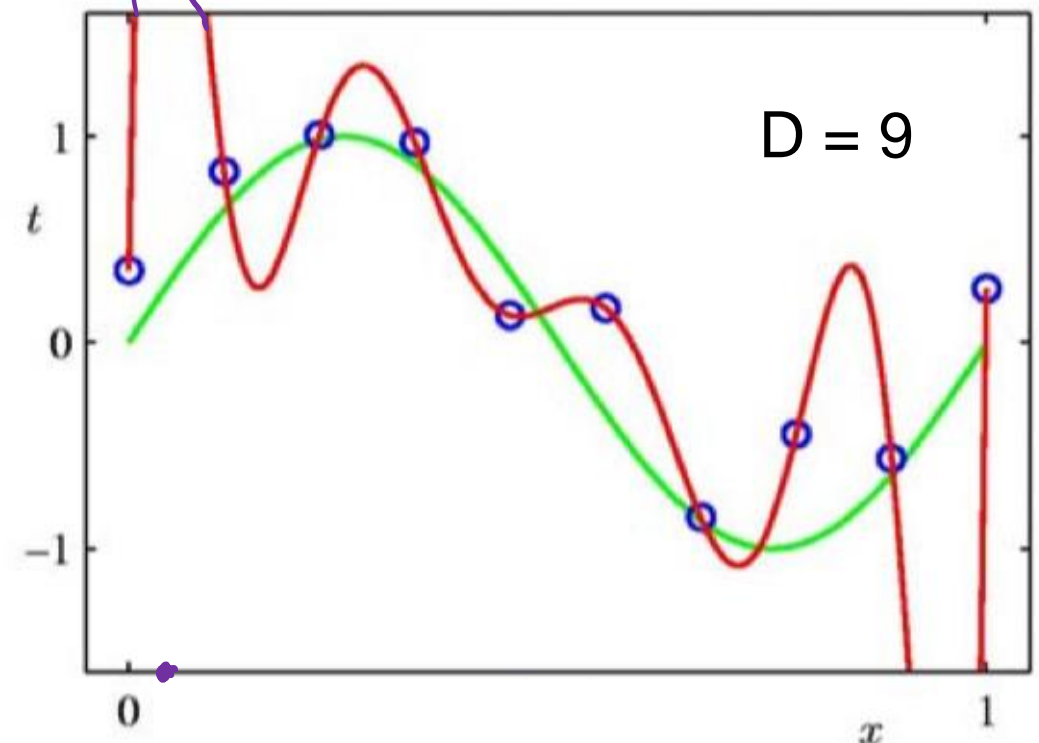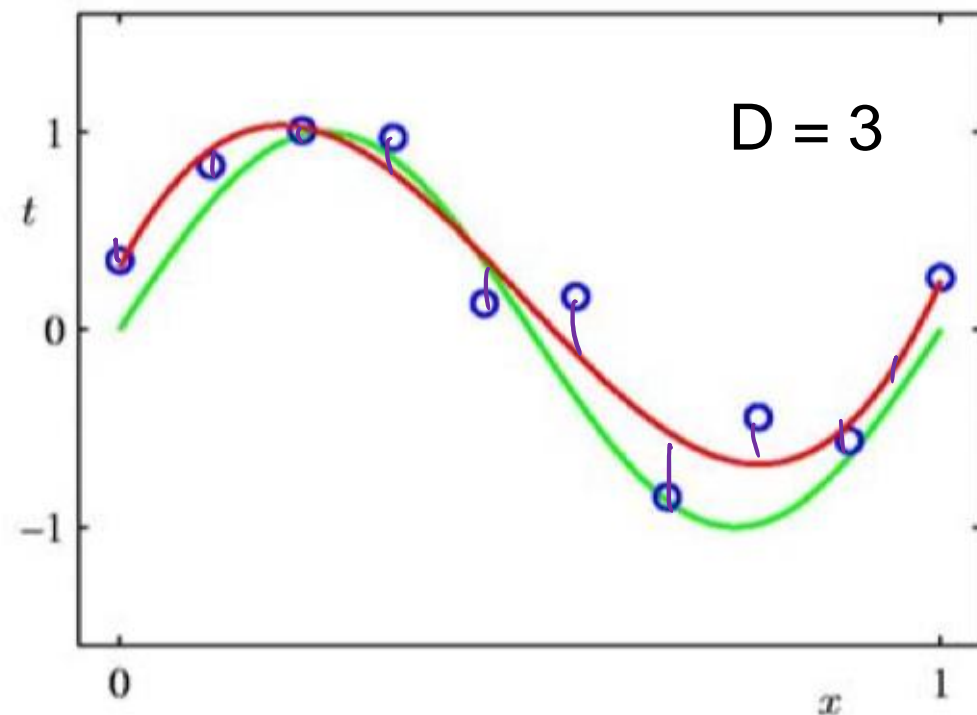
$$y = [3, 3.4875, 3.95, \ldots, 7.98, 8]$$
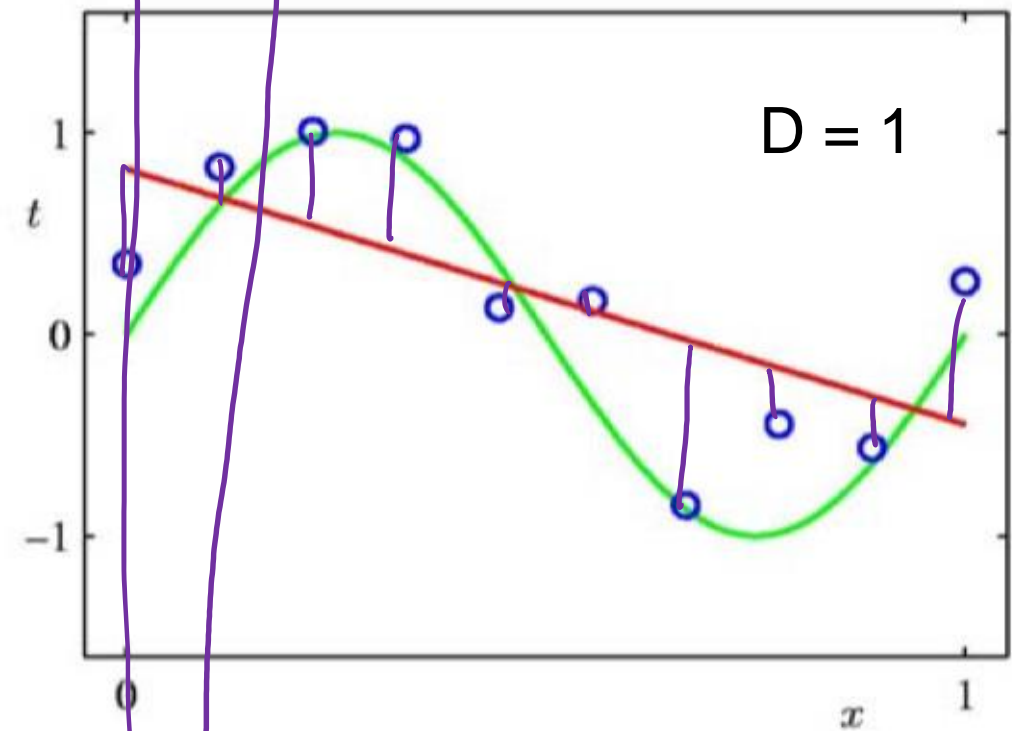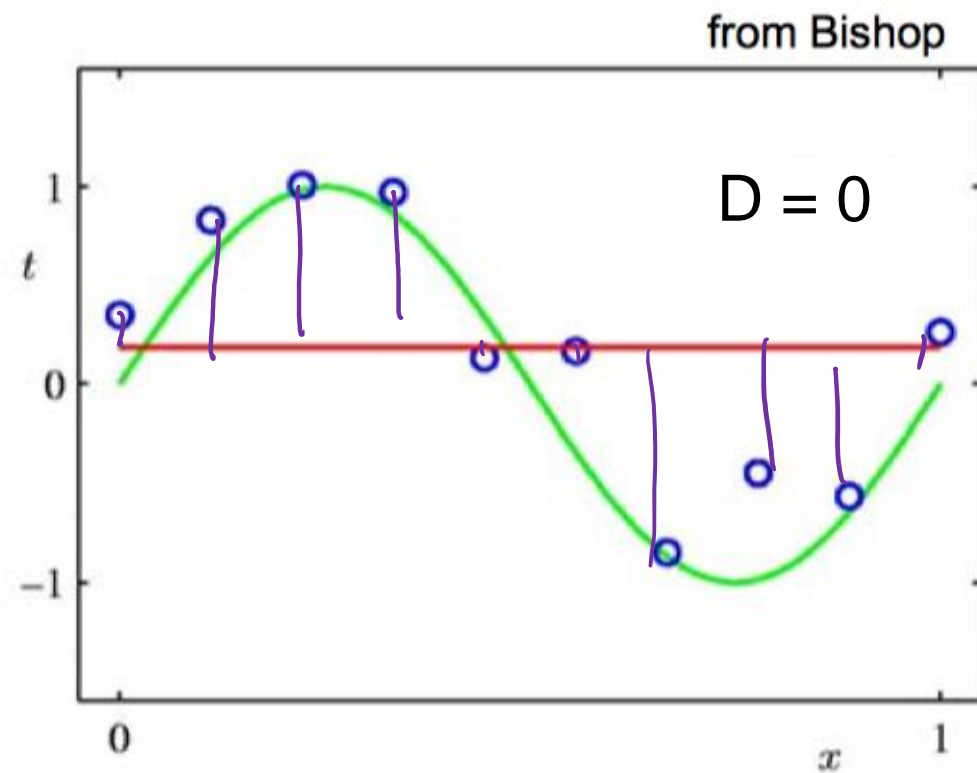
We are fitting a D-dimensional hyperplane in a D+1 dimensional hyperspace (in above example a 2D plane in a 3D space). That hyperplane really is 'flat' / 'linear' in 3D. It can be seen a non-linear regression (a curvy line) in our 2D example in fact it is a flat surface in 3D. So the fact that it is mentioned that the model is linear in parameters, it is shown here.

$y = \theta_0 + \theta_1 x_1 + \cdots + \theta_d x_d$

# Increasing the Maximal Degree

from Bishop



D = 0

D = 1

D = 3

D = 9

# Bias-Variance Trade off

We will have multiple prediction values (i.e. through Cross validation) $E[y_p]$

$$L(\theta) = \frac{1}{n}\sum_{i=1}^{n}(y_i - x_i\theta)^2 = E\left[(y_a - y_p)^2\right]$$

$$(y_a - y_p)^2 = (y_a - E[y_p] + E[y_p] - y_p)^2$$

$E[E(y_p)] = E[y_p]$

$$E$$

$$= (y_a - E[y_p])^2 + (E[y_p] - y_p)^2 + 2(y_a - E[y_p])(E[y_p] - y_p)$$

$E[y_p]$

$E[y_p]$

$0$

$$\mathbf{E}\left[(\boldsymbol{y_a} - \boldsymbol{y_p})^2\right] = (y_a - E[y_p])^2 + E\left[(E[y_p] - y_p)^2\right]$$
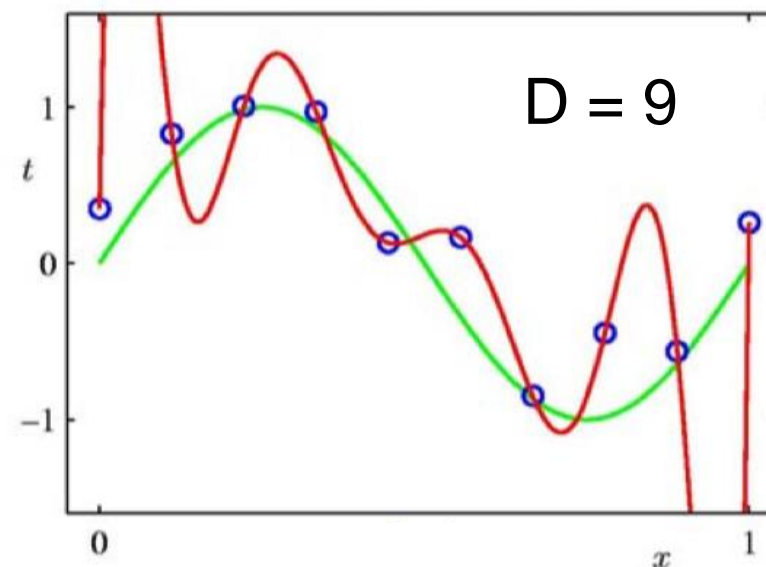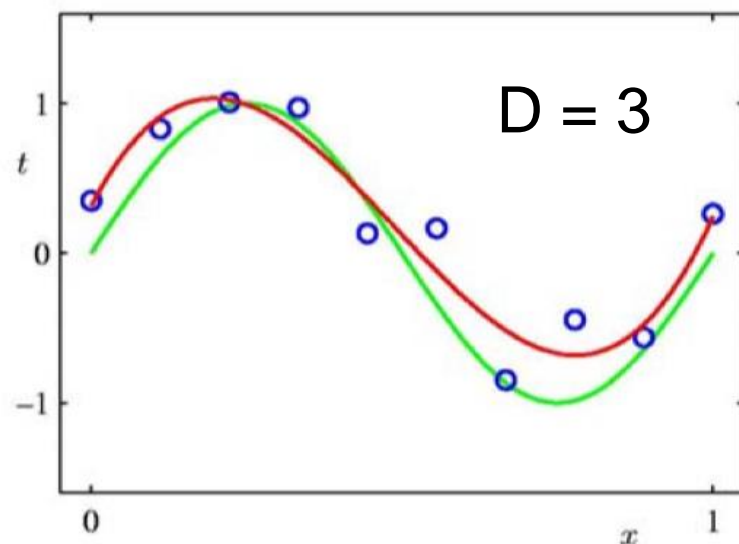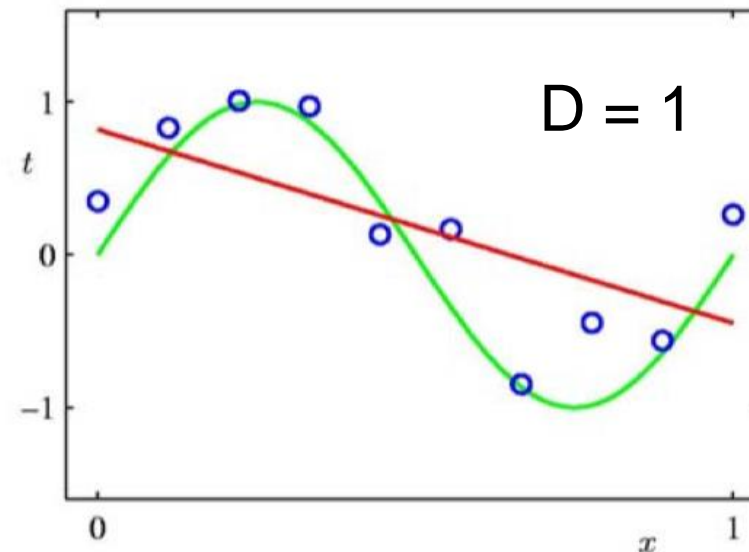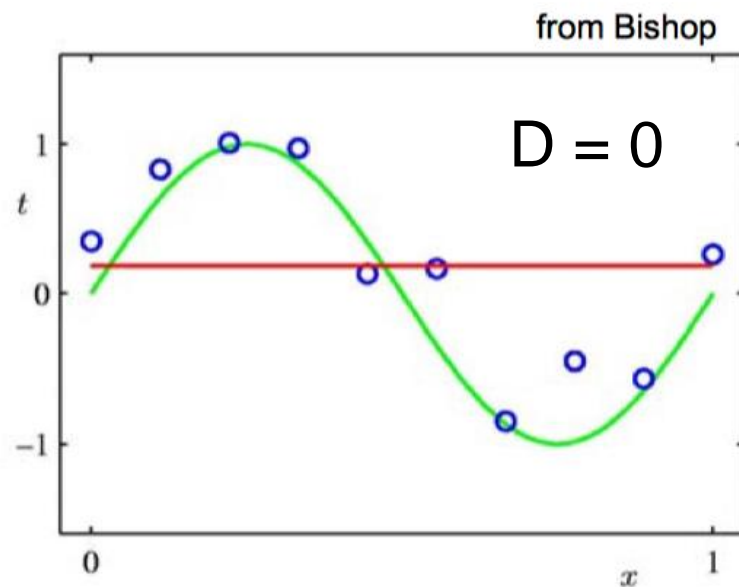
$$= [Bias]^2 + Variance$$

$$= [true\ value - mean(predictions)]^2 - mean[(mean(prediction) - prediction)^2]$$

Why $E[2(y_a - E[y_p])(E[y_p] - y_p)] = 0$ ?

$y_a - E[y_p]$ is a scalar, therefore $E\left[y_a - E[y_p]\right] = y_a - E[y_p]$

$$E[2(y_a - E[y_p])(E[y_p] - y_p)]$$

$$= 2(y_a - E[y_p])E[E[y_p] - y_p]$$

$$= 2(y_a - E[y_p])\left(E\left[E[y_p]\right] - E[y_p]\right)$$

$$= 2(y_a - E[y_p])(E[y_p] - E[y_p]) = 0$$

# Which One is Better?



from Bishop

D = 0

D = 1

D = 3

D = 9

- Can we increase the maximal polynomial degree to very large, such that the curve passes through all training points?

  - We will know the answer in next lecture.

# Take-Home Messages

- Supervised learning paradigm

- Linear regression and least mean square

- Extension to high-order polynomials