

$$X = \begin{bmatrix} h & w & a \end{bmatrix} \rightarrow Z = \begin{bmatrix} h \end{bmatrix} \rightsquigarrow \text{feed it in an ML algorithm (NB, SVM, ...)} \\ \hookrightarrow \text{acc} = 80\%$$

$$Z = \begin{bmatrix} w \end{bmatrix} \rightsquigarrow \text{acc} = 70\% \quad Z = \begin{bmatrix} a \end{bmatrix} \rightsquigarrow \text{acc} = 65\%$$

$$Z = \begin{bmatrix} h & w \end{bmatrix} \rightsquigarrow \text{acc} = 85\% \quad Z = \begin{bmatrix} h & a \end{bmatrix} \rightsquigarrow \text{acc} = 58\%$$

Forward Feature Selection

Backward Feature Selection

$$X = \begin{bmatrix} h & \omega & a \end{bmatrix}^T$$

→ LDA

Linear Regression

Mahdi Roozbahani
Georgia Tech

The Moment



Incoming ML HW3

You; after HW2

Outline

- Supervised Learning 
- Linear Regression
- Extension

Supervised Learning: Overview

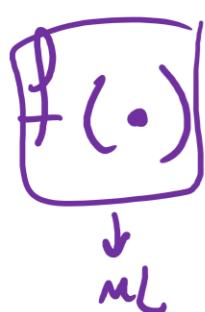
Functions \mathcal{F}

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

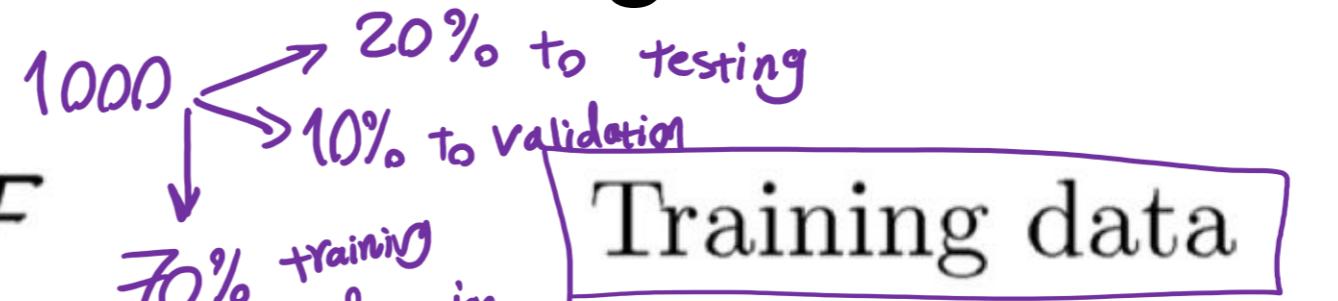
$$X = \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}_{n \times d} \quad y_a = \begin{bmatrix} 0 \\ \vdots \\ 0 \end{bmatrix}_{n \times 1}$$

LEARNING

$$f(x^i) = \hat{y}_p \quad \hat{y}_p = \begin{bmatrix} 0 \end{bmatrix}$$



PREDICTION



$$\{(x^{(i)}, y^{(i)}) \in \mathcal{X} \times \mathcal{Y}\}$$

find $\hat{f} \in \mathcal{F}$
s.t. $y_a \approx f(x^{(i)})$

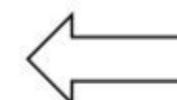


Learning machine



New data

$$\hat{y}_p = f(x^{(i)})$$



$$x$$

Supervised Learning: Two Types of Tasks

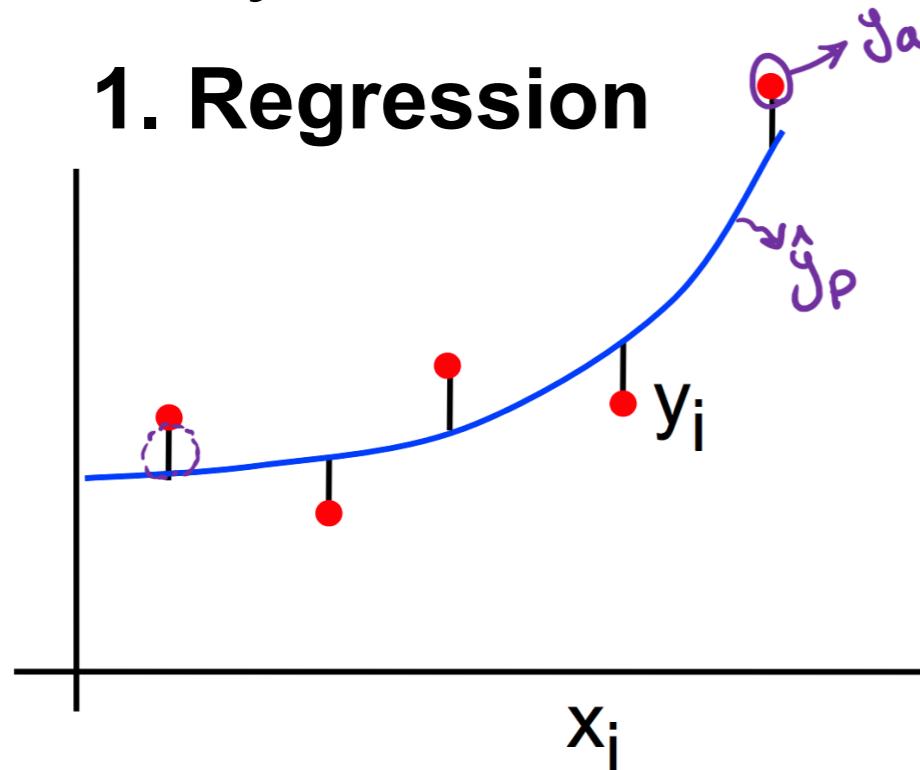
Given: training data $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(n)}, y^{(n)})\}$

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

①

When y is continuous:

1. Regression



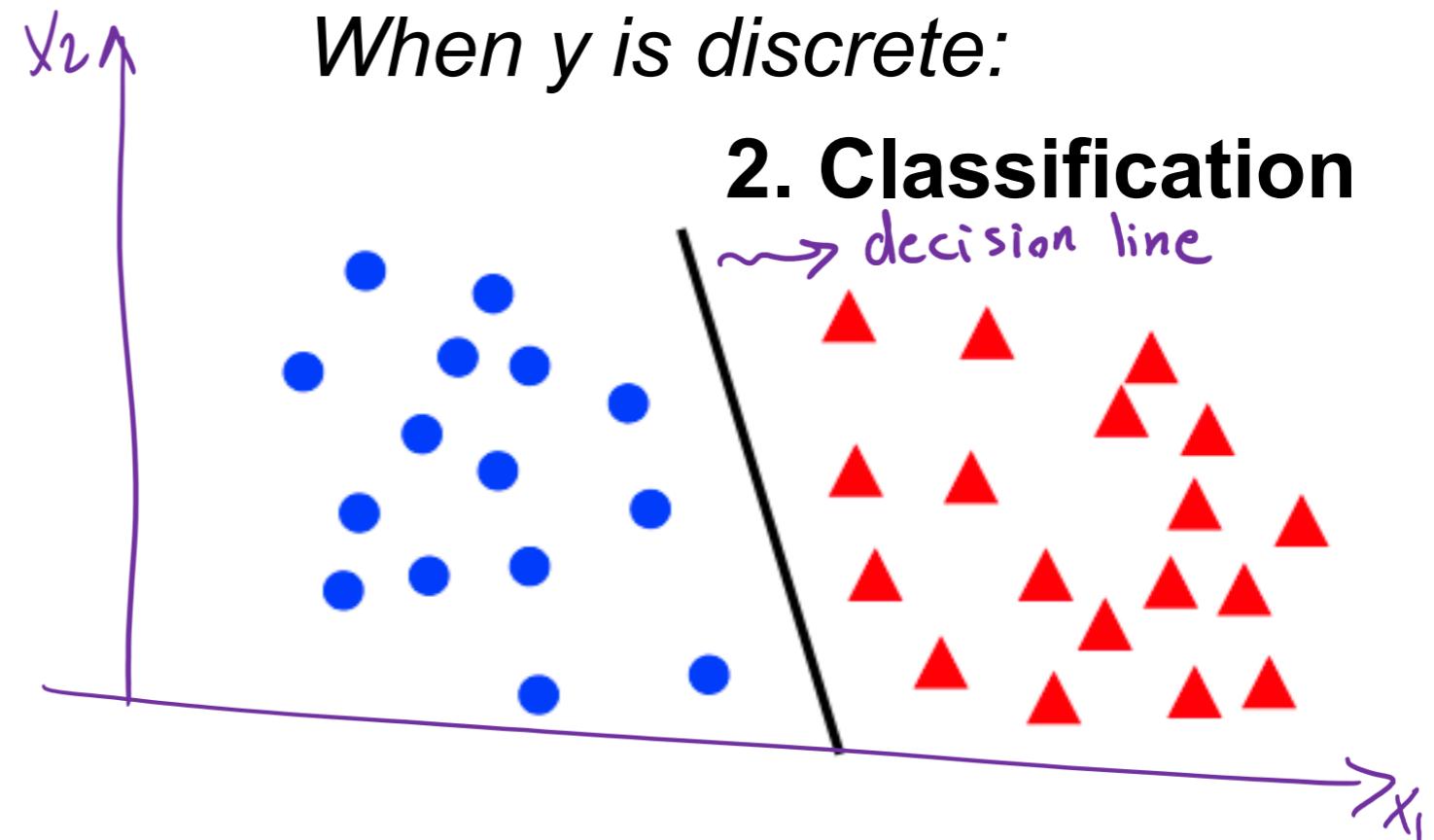
Curve fitting

difference between \hat{y}_p & y_a

②

When y is discrete:

2. Classification



Class estimation

Classification Example 1: Handwritten digit recognition

As a supervised classification problem

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

0 0 0 1 1 1 1 1 1 2

2 2 2 2 2 2 2 3 3 3

Multi class Classification

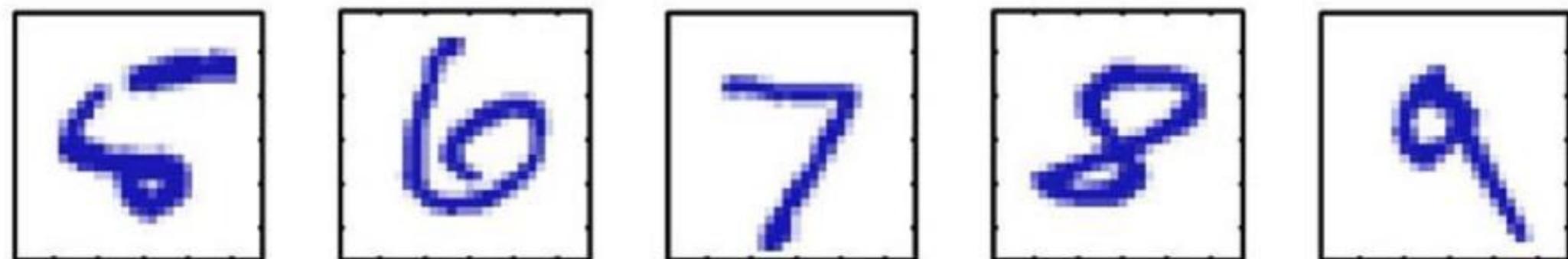
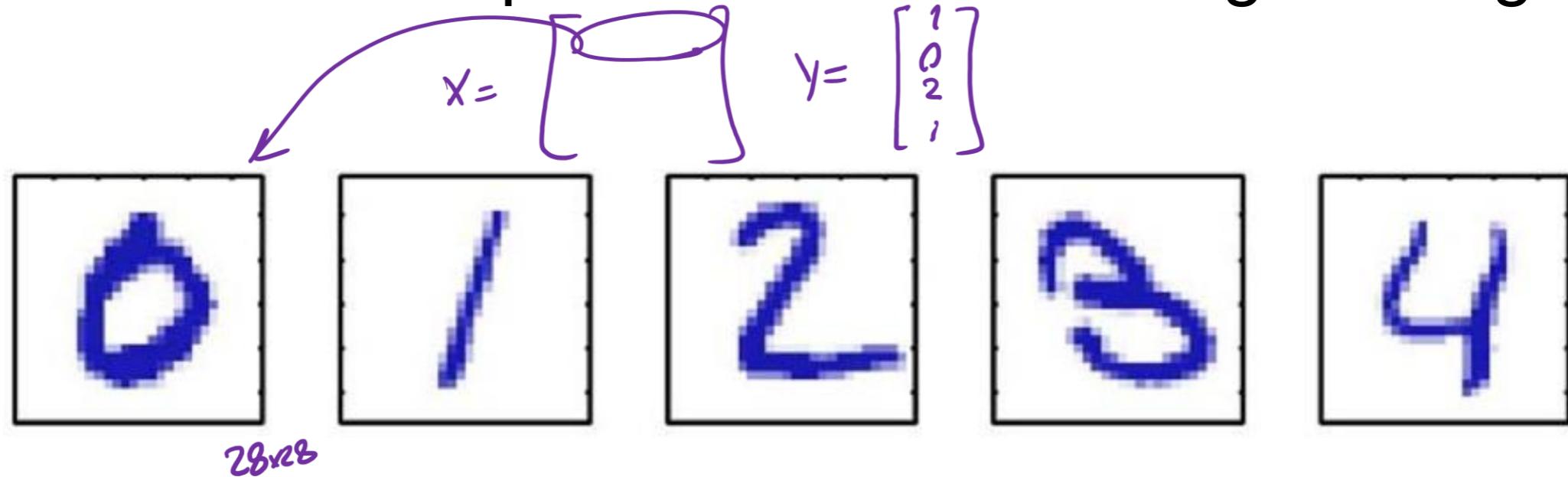
3 4 4 4 4 4 5 5 5 5

6 6 7 7 7 7 7 8 8 8

9 9 9 9 9 9 9 9 9 9

- 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 28×28 pixels

A classification problem

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

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

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

Classification Example 2: Spam Detection

Google in:spam

Gmail ▾

COMPOSE

Inbox (994)

Starred

Sent Mail

Drafts

Less ▾

Important

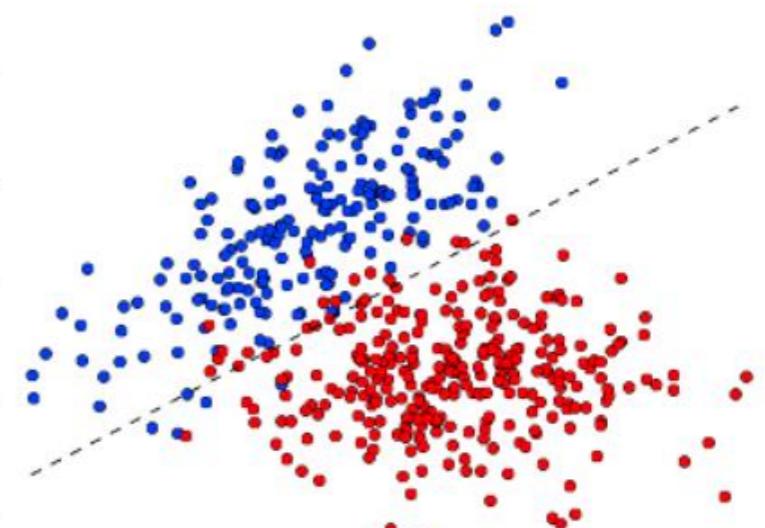
+

Customer Service You still have product(s) in your basket - Healthy Living Lifestyle Pre
Sherley Rhoda From Sherley Rhoda
Customer Service Activate your favorite videotracking service - Your activation code is re
Healthy Living We have added your shopping credits today - Healthy Living & Co. I
Shiningltd Team 15 inch wifi Android OS tablet pc - SHININGLTD Our Alibaba Shop C
wikiHow Community Team (2) Congratulations on your article's first Helpful Vote! - Congratulations! A I
FreeLotto Jesse, NOTICE of FORFEITURE - Do not ignore! - NEVER miss an i
Good Fella's Our team assigned you to receive our new phone - Good Fella's Au
Jason Squires Make 2018 your best year yet - Hi there, Hope you're well, and have h
Bunnings January arrivals - Image Congratulations Jesse Eaton! We have a very

Delete all spam messages now (messages that have been in Spam more than

Binary Classification

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

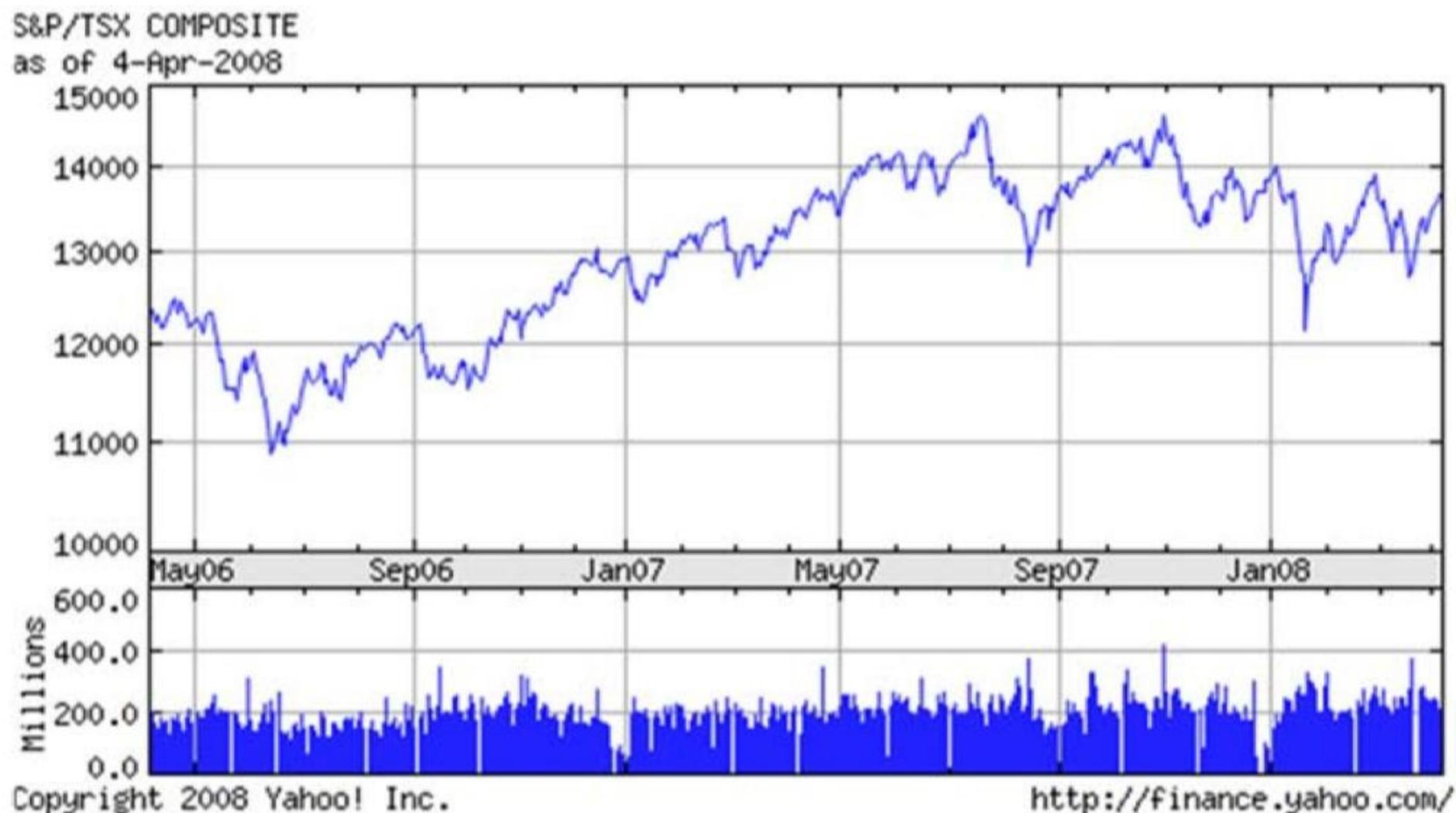
Regression Example 1: Apartment Rent Prediction

- 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

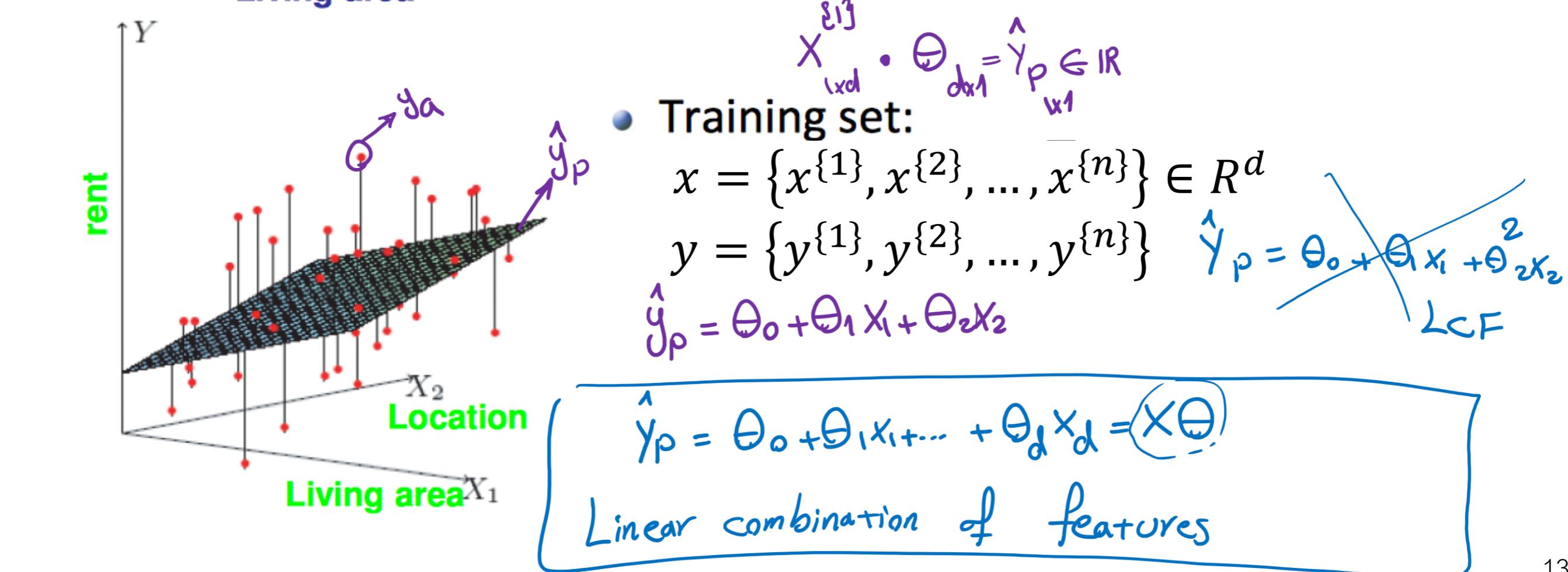
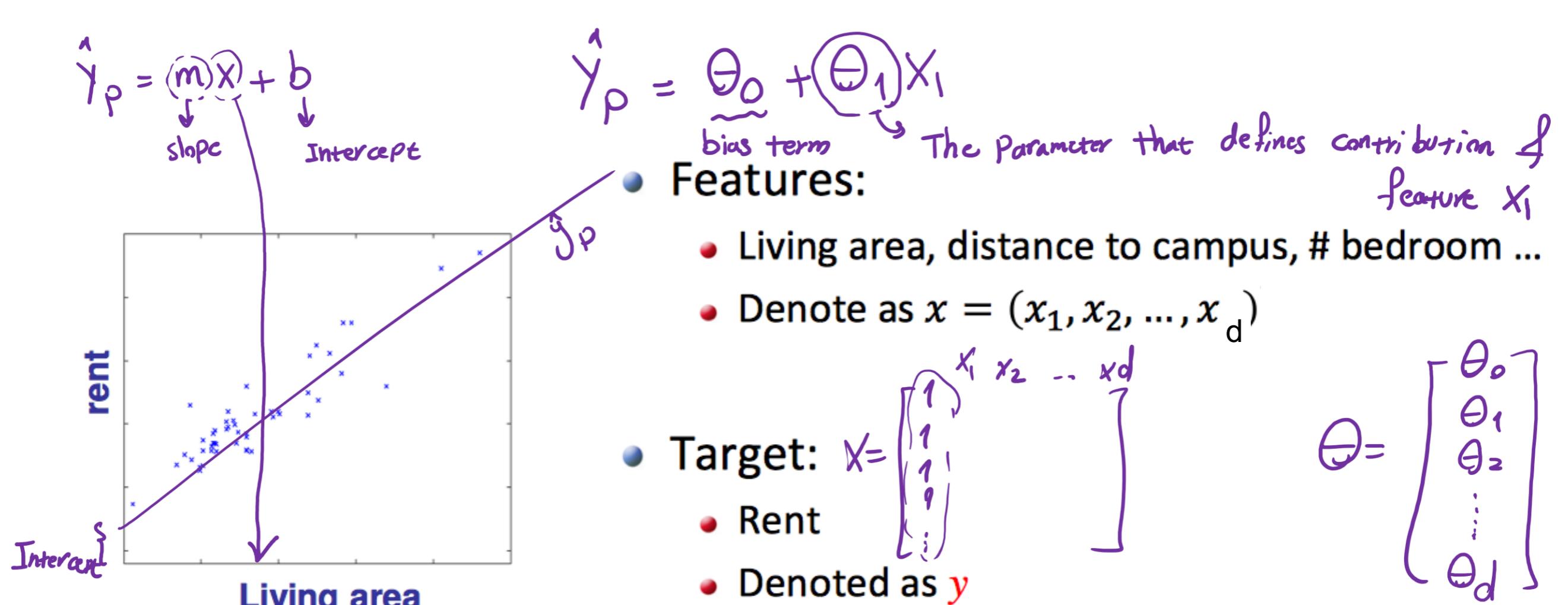
Living area (ft ²)	# bedroom	Rent (\$)
230	1	600
506	2	1000
433	2	1100
109	1	500
...		
150	1	?
270	1.5	?

Regression Example 2: Stock Price Prediction

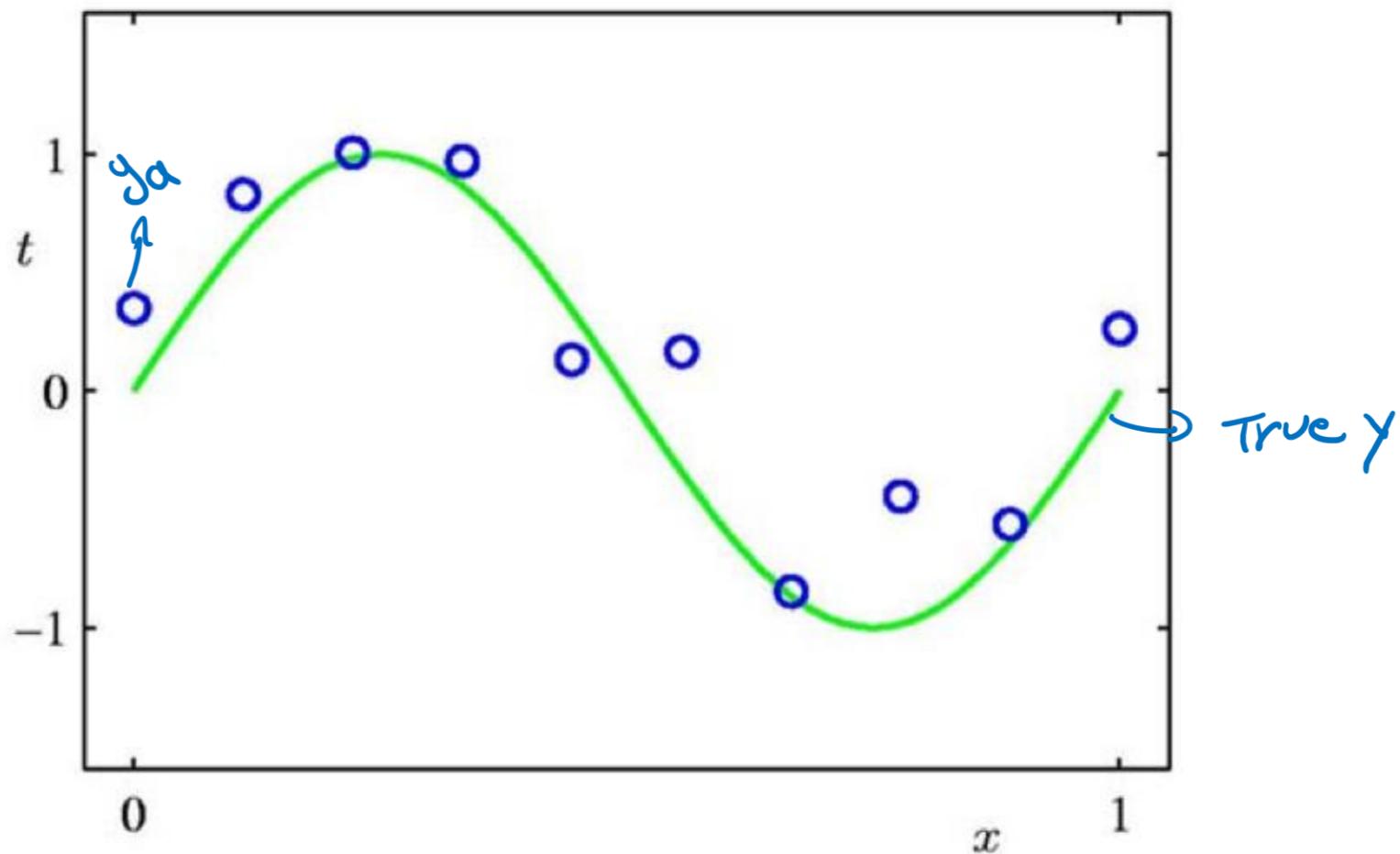


- Task is to predict stock price at future date

A regression problem



Regression: Problem Setup



- Suppose we are given a training set of N observations
- Regression problem is to estimate $y(x)$ from this data

Outline

- Supervised Learning
- Linear Regression ←
- Extension

Linear Regression

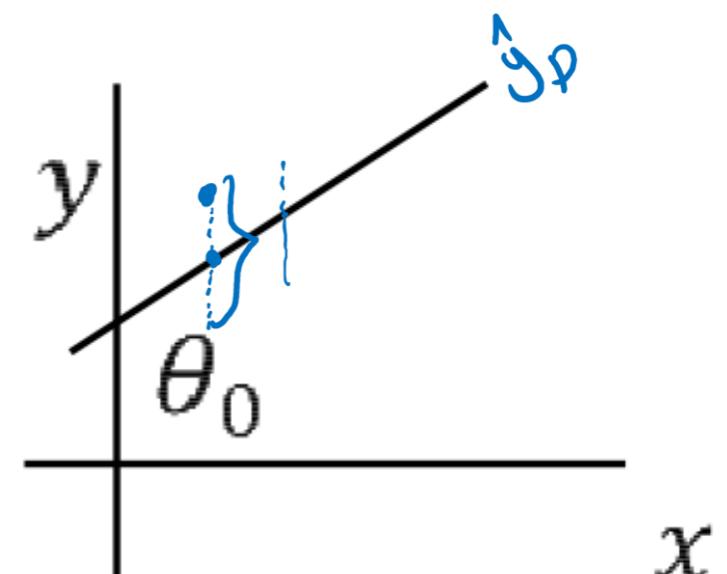
- Assume y is a linear function of x (features) plus noise ϵ

$$y = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d + \epsilon = X\theta$$

LCF

- where ϵ is an error term of unmodeled effects or random noise
- Let $\theta = (\theta_0, \theta_1, \dots, \theta_d)^T$, and augment data by one dimension

- Then $y = x\theta + \epsilon$



$E[\cdot] \rightarrow$ expectation

$$y = |x| \quad \vee \quad y = x^2$$

$E(\cdot) \rightarrow$ error

Least Mean Square Method

- Given n data points, find θ that minimizes the mean square error

Training $\hat{\theta} = \operatorname{argmin}_{\theta} L(\theta)$

$\xrightarrow{\text{Convex}}$

$$L(\theta) = \frac{1}{n} \sum_{i=1}^N (y^{(i)} - \hat{y}_P)^2 = E_{\theta}[(y^{(i)} - \hat{y}_P)^2]$$

$\in \mathbb{R}$

\downarrow error

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

$$\frac{\partial L(\theta)}{\partial \theta} = 0$$

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

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

For Loop

$$X\theta = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

Matrix form

↑
Column vector

$$x = \begin{bmatrix} 1 & x_1^{(1)} & \dots & x_d^{(1)} \\ 1 & x_1^{(2)} & \ddots & x_d^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{(n)} & \dots & x_d^{(n)} \end{bmatrix}_{n \times (d+1)}$$

$$y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}_{n \times 1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}_{(d+1) \times 1}$$

$$MSE(\theta) = \underset{\theta}{argmin} L(\theta) = \frac{1}{n} (y - x\theta)^T (y - x\theta)$$

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

$\hat{y}_p^{(1)}$
 $\hat{y}_p^{(2)}$
 $\approx \hat{y}_p$

Matrix Version and Optimization

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

$x^T x \theta$

Let's rewrite it as:

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n} (x^{(1)}, \dots, x^{(n)})^T (y^{(1)}, \dots, y^{(n)}) + \frac{2}{n} (x^{(1)}, \dots, x^{(n)})^T (x^{(1)}, \dots, x^{(n)}) \theta = 0$$

Define $X = (x^{(1)}, \dots, x^{(n)})$ and $y = (y^{(1)}, \dots, y^{(n)})$

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

$$\Rightarrow \theta = \underbrace{(X^T X)^{-1}}_{d \times d} \underbrace{X^T y}_{d \times 1} = X^+ y$$

$$\hat{y}_p = X \cdot \theta \Rightarrow \theta = X^{-1} y$$

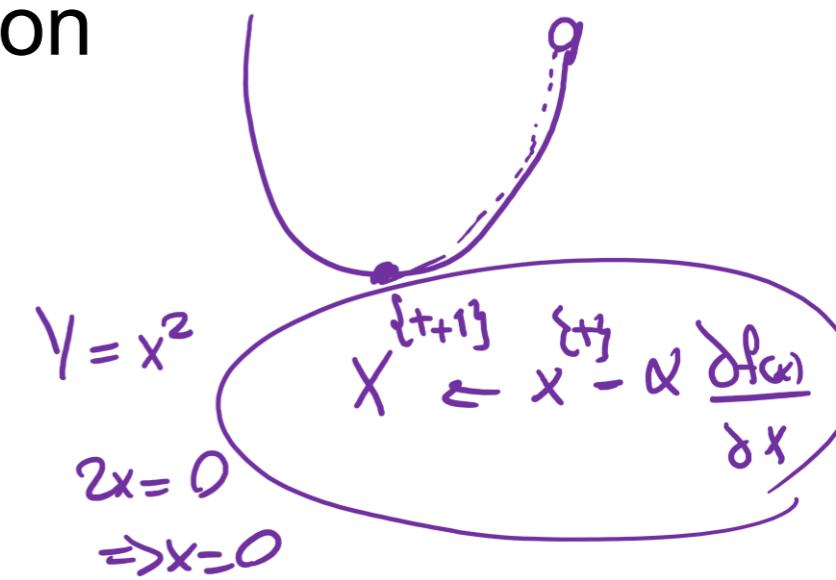
X^+ is the **pseudo-inverse** of X
 $X^T X X^+ = X^T$

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

$X_{n \times d}$

$n = \text{instances}$ $d = \text{dimension}$

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



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

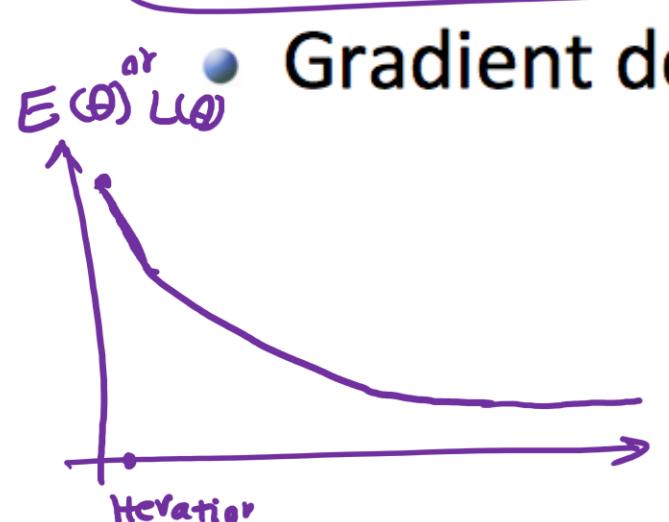
Alternative Way to Optimize

- The matrix inversion in $\hat{\theta} = (X^T X)^{-1} X^T y$ can be very expensive to compute

$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n} \sum_{i=1}^N (x^{(i)})^T (y^{(i)} - x^{(i)}\theta)$$

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

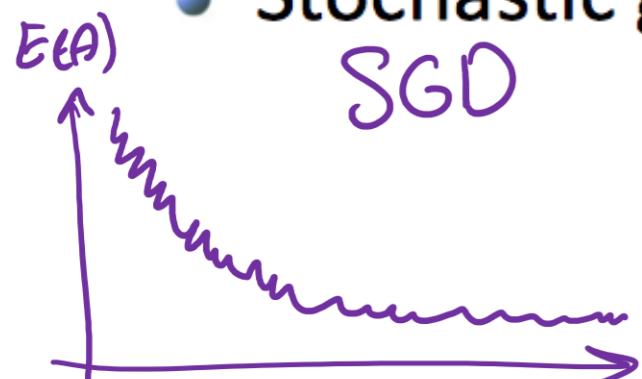
Gradient descent



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

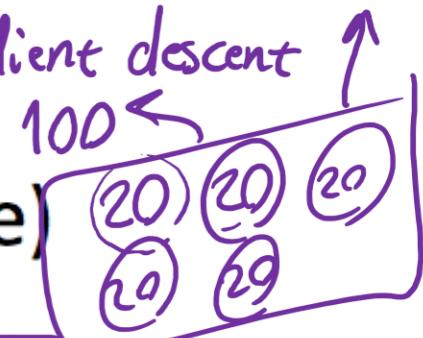
Full FGD

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



$$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \beta_t \times (x^{(i)})^T (y^{(i)} - x^{(i)}\theta)$$

BGD \rightarrow Batch Gradient descent



Recap

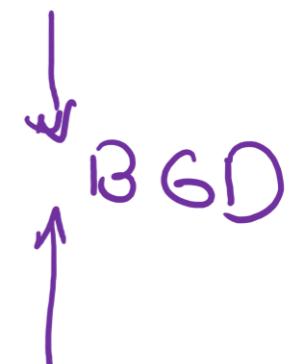
- Stochastic gradient update rule

$$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \beta_t \times (x^{\{i\}})^T (y^{\{i\}} - x^{\{i\}}\theta) \rightarrow SGD$$

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

- Gradient descent

$$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \frac{\alpha}{n} \sum_{i=1}^N (x^{\{i\}})^T (y^{\{i\}} - x^{\{i\}}\theta) \rightsquigarrow FGD$$



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

- Solve normal equations

$$\theta = (X^T X)^{-1} X^T y \rightsquigarrow \text{closed form solution}$$

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

$$\hat{y}_p = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d = \mathbf{x} \theta$$

column vector

$$LCF$$

$$\begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

$$x^{\{i\}} \theta_{d+1 \times 1} \in \mathbb{R}$$

Convex optimization

$$L(\theta) = E(\theta) = \frac{1}{N} \sum_{i=1}^N (y_a - \hat{y}_p)^2$$

avg

$$\theta_{d+1 \times 1} = \left(\begin{matrix} X^T & X \end{matrix} \right)^{-1} \begin{matrix} X^T \\ d+1 \times n \end{matrix} \begin{matrix} Y \\ n \times 1 \end{matrix} \rightarrow \text{Closed Form Solution}$$

FGD, SGD, BGD

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

Linear regression for classification

Raw Input $x = (1, x_1, \dots, x_{256})$

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

Extract useful information

$$\theta_0 \quad \theta_1 \quad \theta_2$$

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



Sum up all the pixels = intensity

Symmetry = -(difference between flip version)

“Feature Engineering”

$$x = (1, x_1, x_2)$$

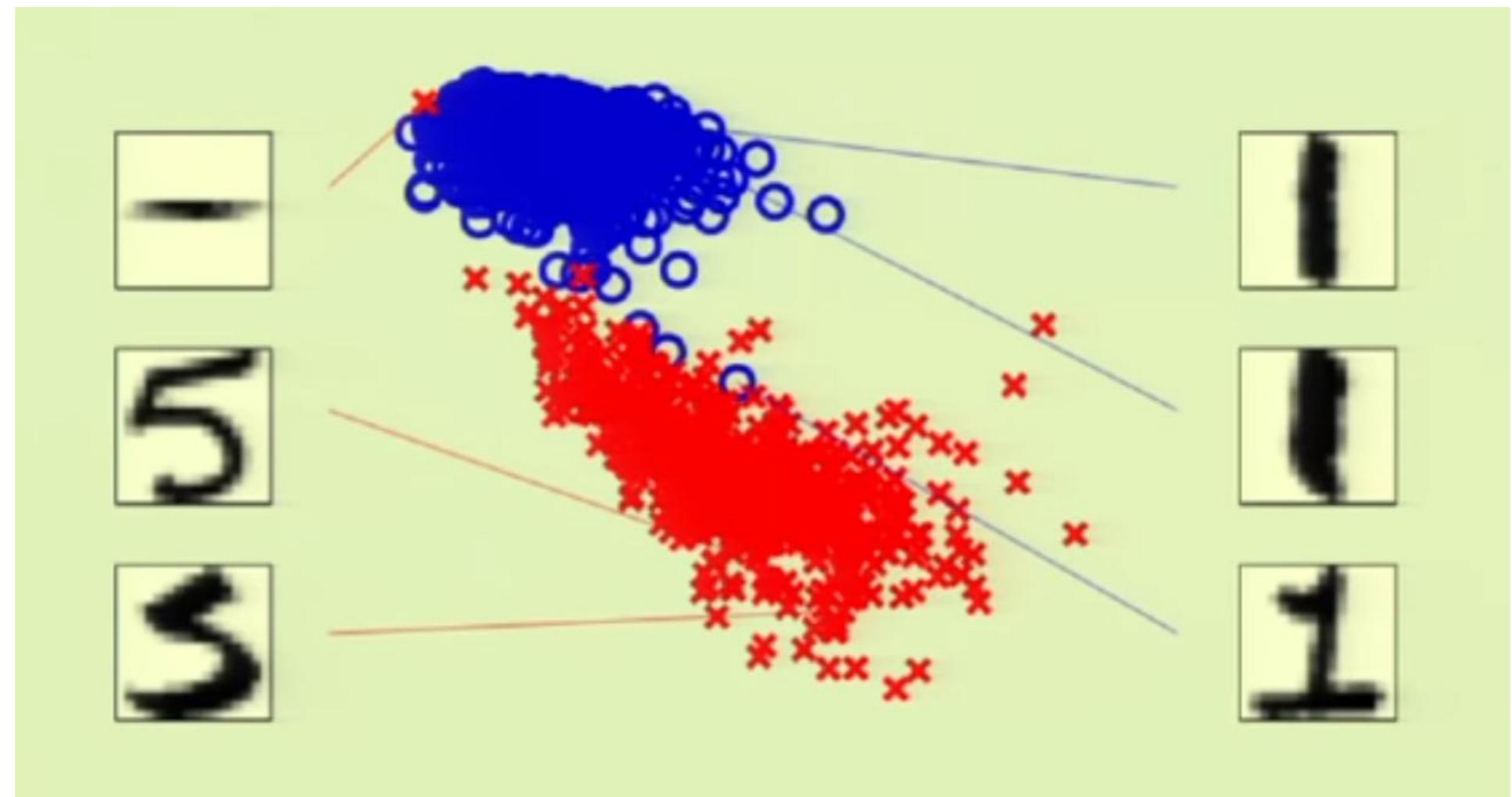
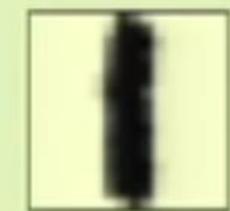
$x_1 = \text{intensity}$ $x_2 = \text{symmetry}$

It is almost linearly separable

symmetry

5

3



intensity

Linear regression for classification

$$x^{\{i\}} \theta \in \mathbb{R}$$

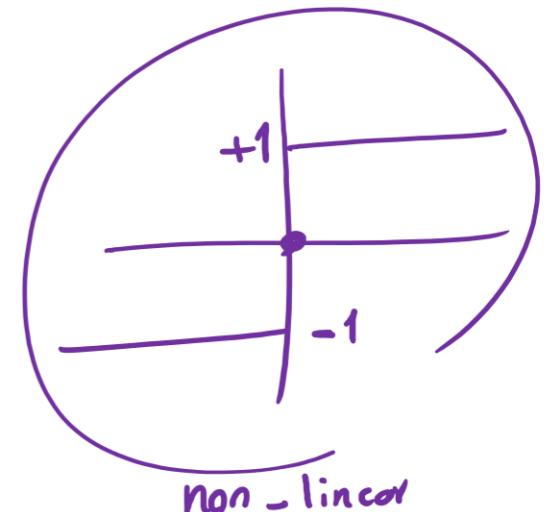
$$Y_{\text{reg}} = \begin{bmatrix} \uparrow \\ \text{continuous} \end{bmatrix} \rightarrow Y_{\text{Class}} = \begin{bmatrix} \uparrow \\ \text{discrete} \end{bmatrix}$$

Binary-valued functions are also real-valued $\pm 1 \in R$

Use linear regression $x^{\{i\}} \theta \approx y^{\{i\}} = \pm 1$ i = index of a data-point

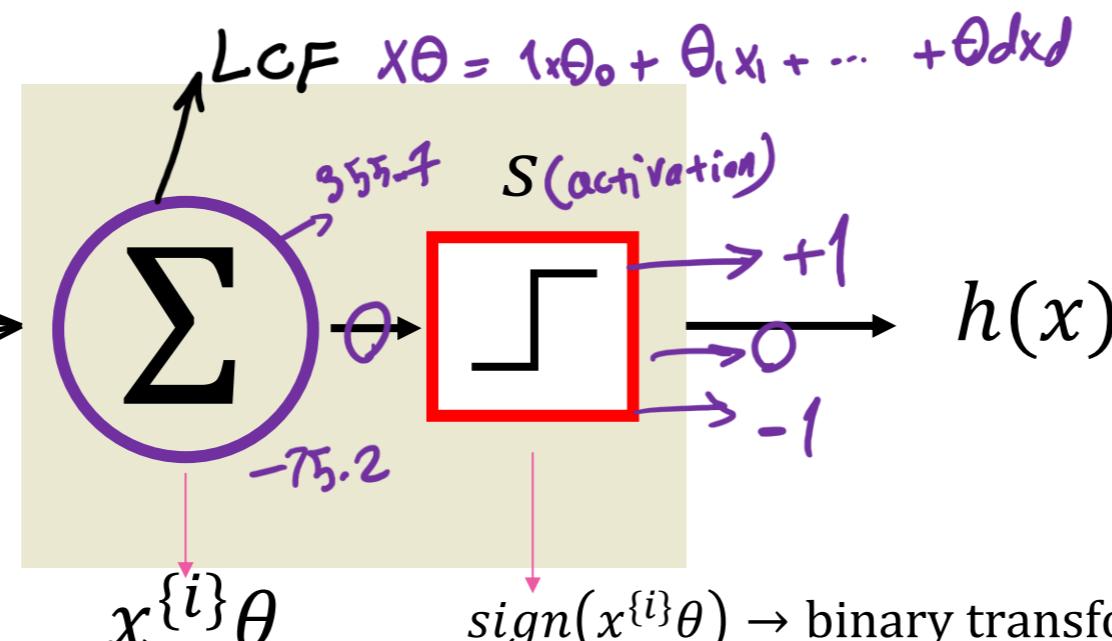
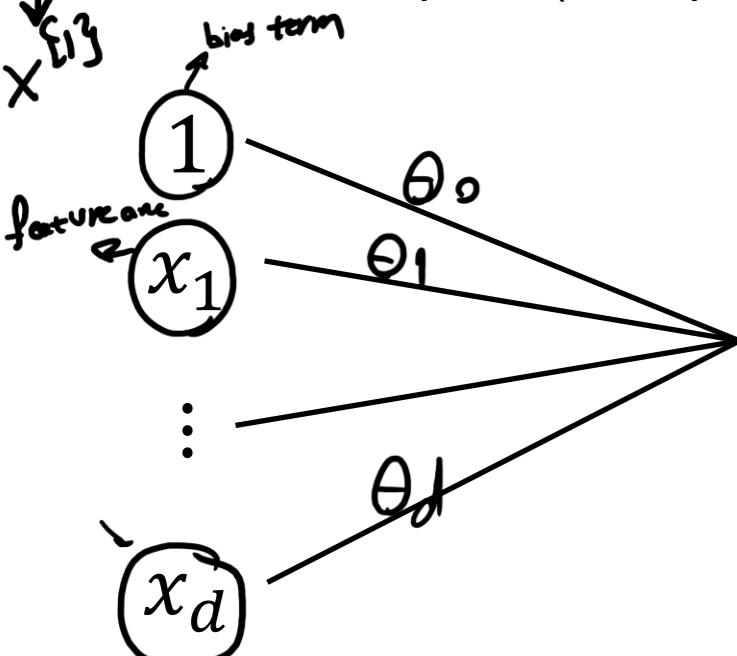
Let's calculate, $\text{sign}(x^{\{i\}} \theta) = \begin{cases} -1 & x^{\{i\}} \theta < 0 \\ 0 & x^{\{i\}} \theta = 0 \\ 1 & x^{\{i\}} \theta > 0 \end{cases}$

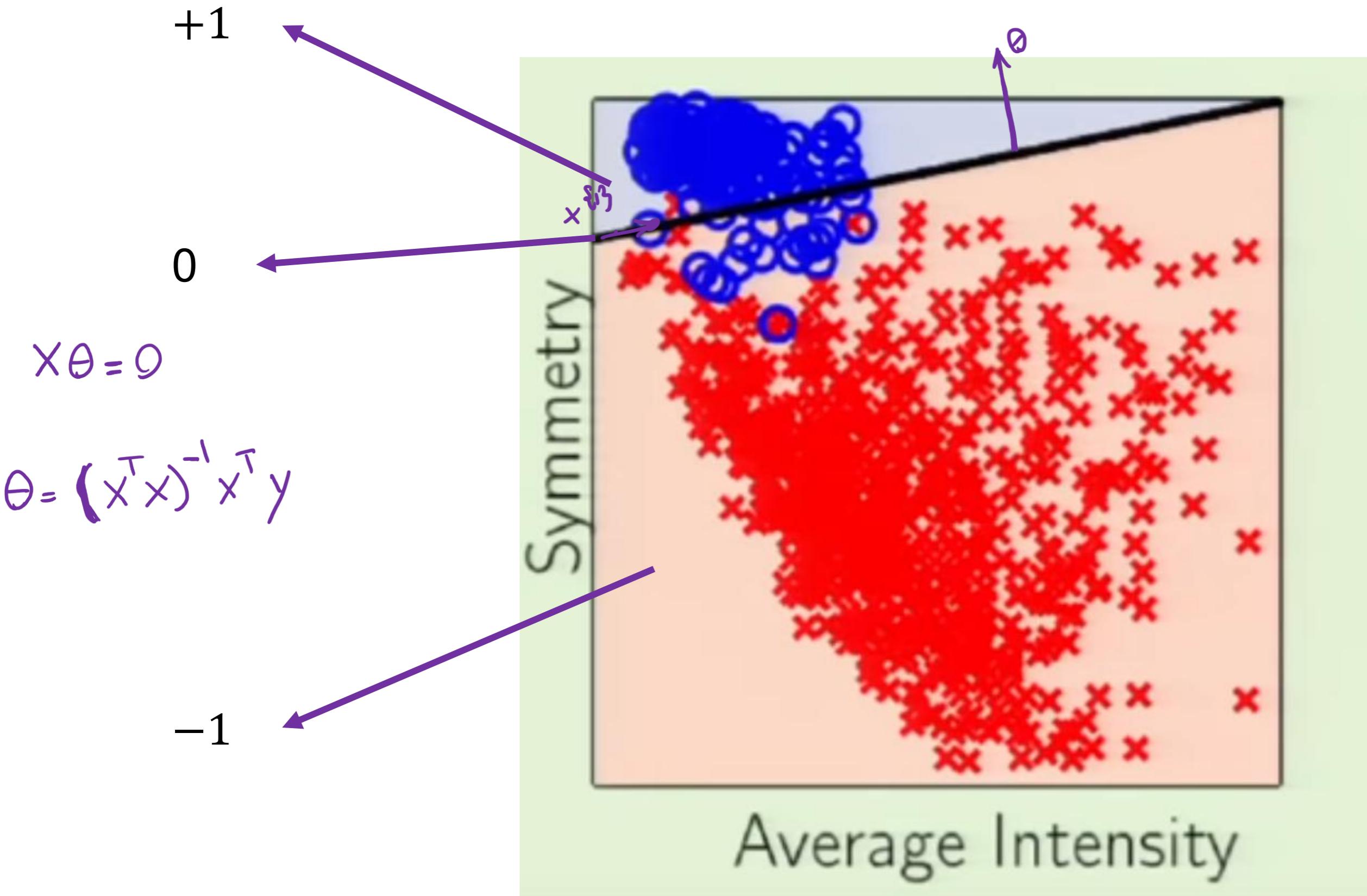
$$\begin{aligned} x^{\{i\}} \theta &< 0 \\ x^{\{i\}} \theta &= 0 \\ x^{\{i\}} \theta &> 0 \end{aligned}$$



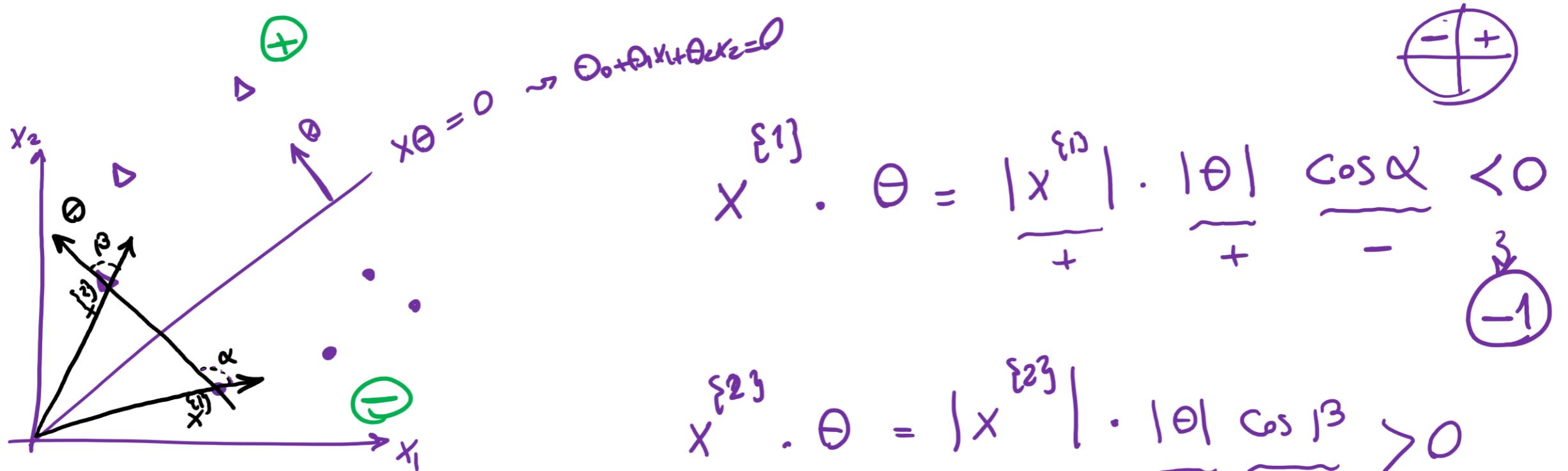
SGD

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





Not really the best for classification, but it's a good start

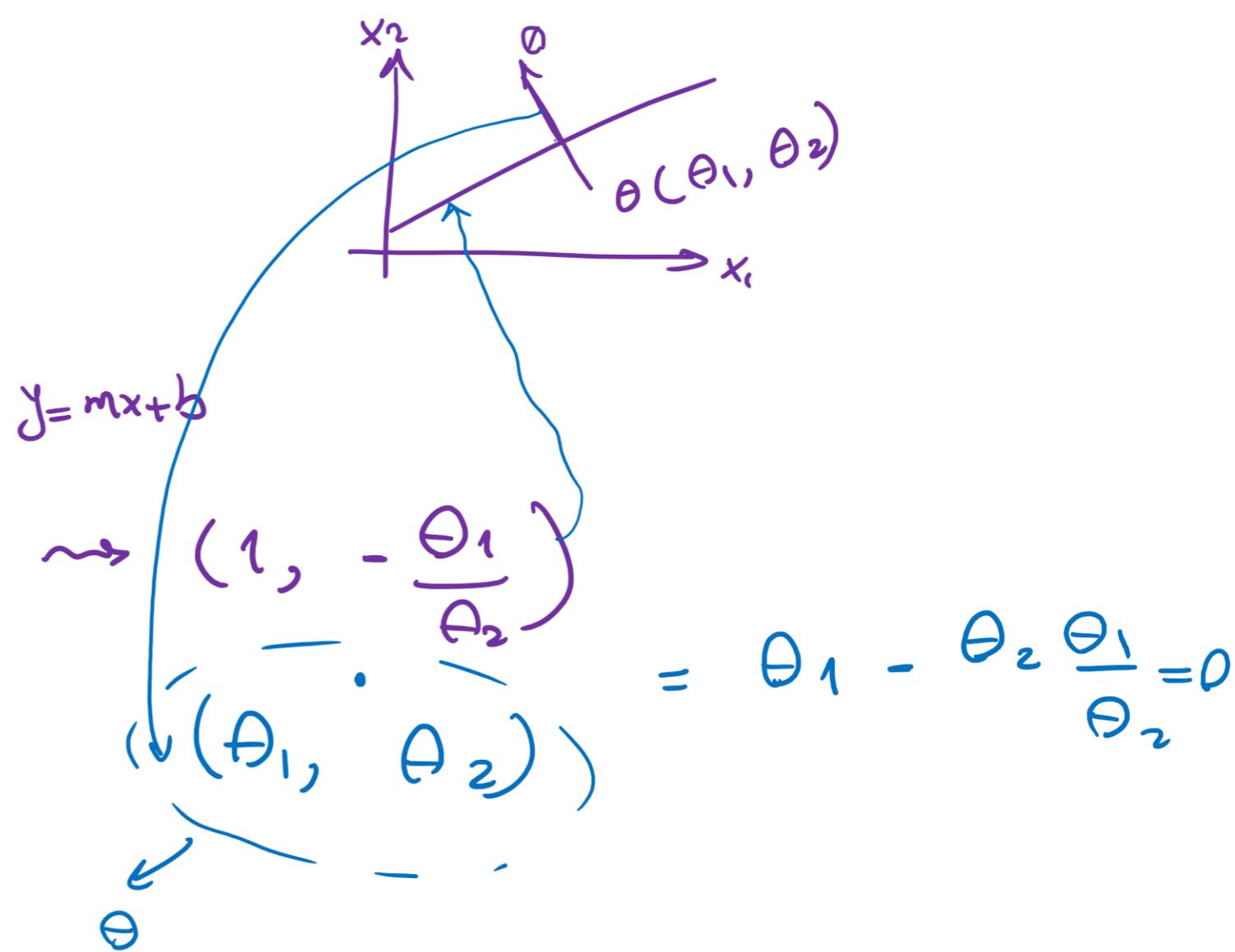


$\hat{y}_p = \theta_0 + \theta_1 x_1 + \theta_2 x_2$ plane

$$\hat{y}_P = \theta_0 + \theta_1 x_1 + \theta_2 x_2 = 0$$

$$\theta_2 x_2 = -(\theta_0 + \theta_1 x_1)$$

$$x_2 = -\frac{\theta_0}{\theta_2} - \frac{\theta_1}{\theta_2} x_1$$

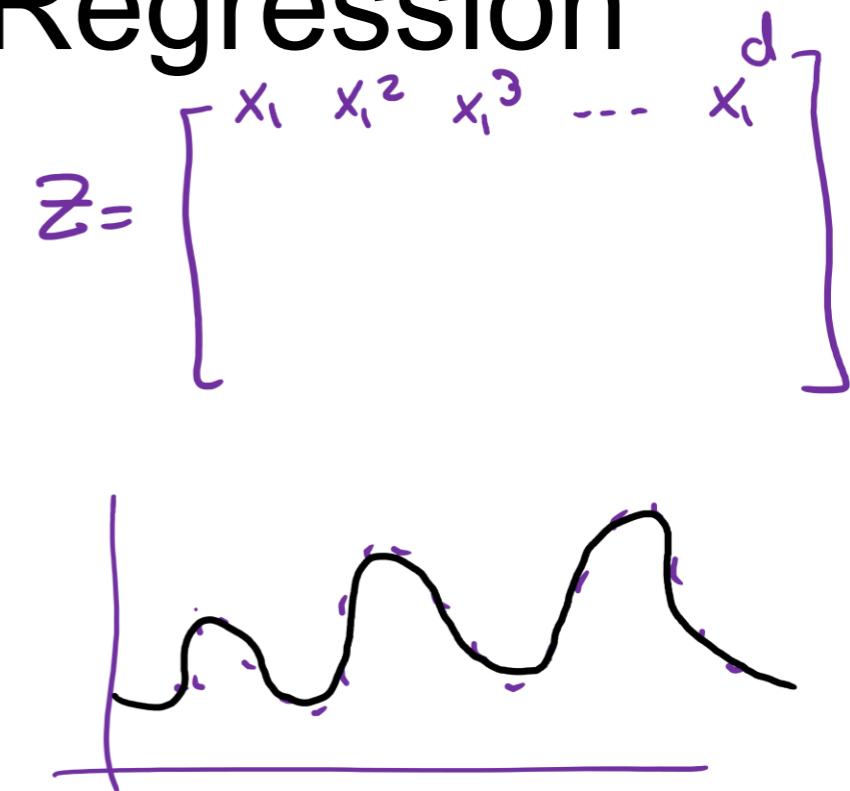
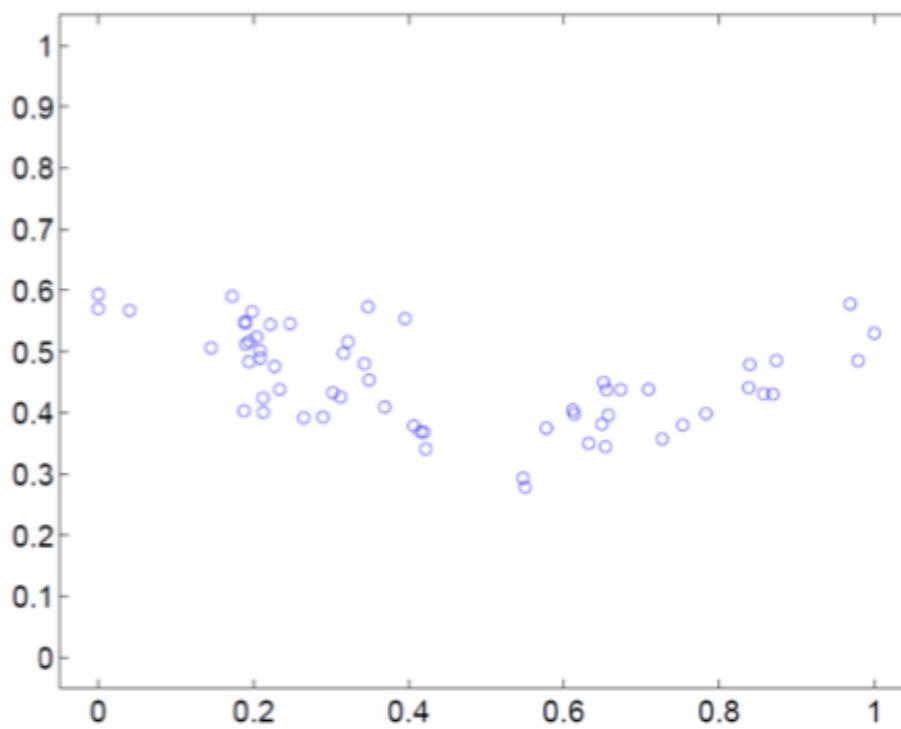
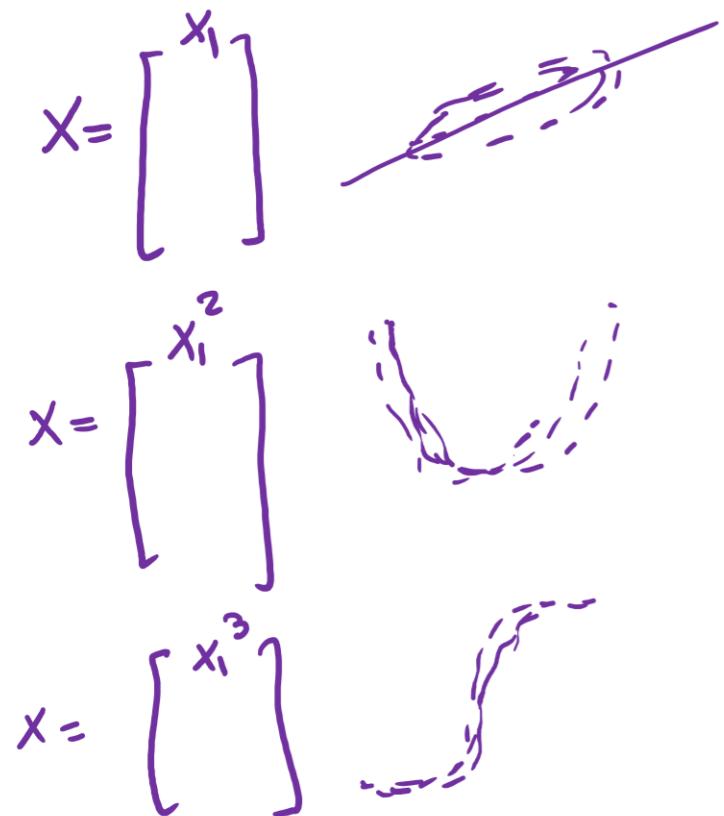


$$X\theta = 0$$

Outline

- Supervised Learning
- Linear Regression
- Extension 

Extension to Higher-Order Regression



- Want to fit a polynomial regression model

LCF

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_d x^d + \epsilon$$

$$\theta_0 z_0 + \theta_1 z_1 + \theta_2 z_2 + \dots + \theta_d z_d$$

- $z = \{1, x, x^2, \dots, x^d\} \in R^d$ and $\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_d)^T$

$$y = z\theta$$

Least Mean Square Still Works the Same

- Given n data points, find θ that minimizes the mean square error

$$\hat{\theta} = \operatorname{argmin}_{\theta} L(\theta) = \frac{1}{n} \sum_{i=1}^N (y^{(i)} - z^{(i)}\theta)^2$$

- 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

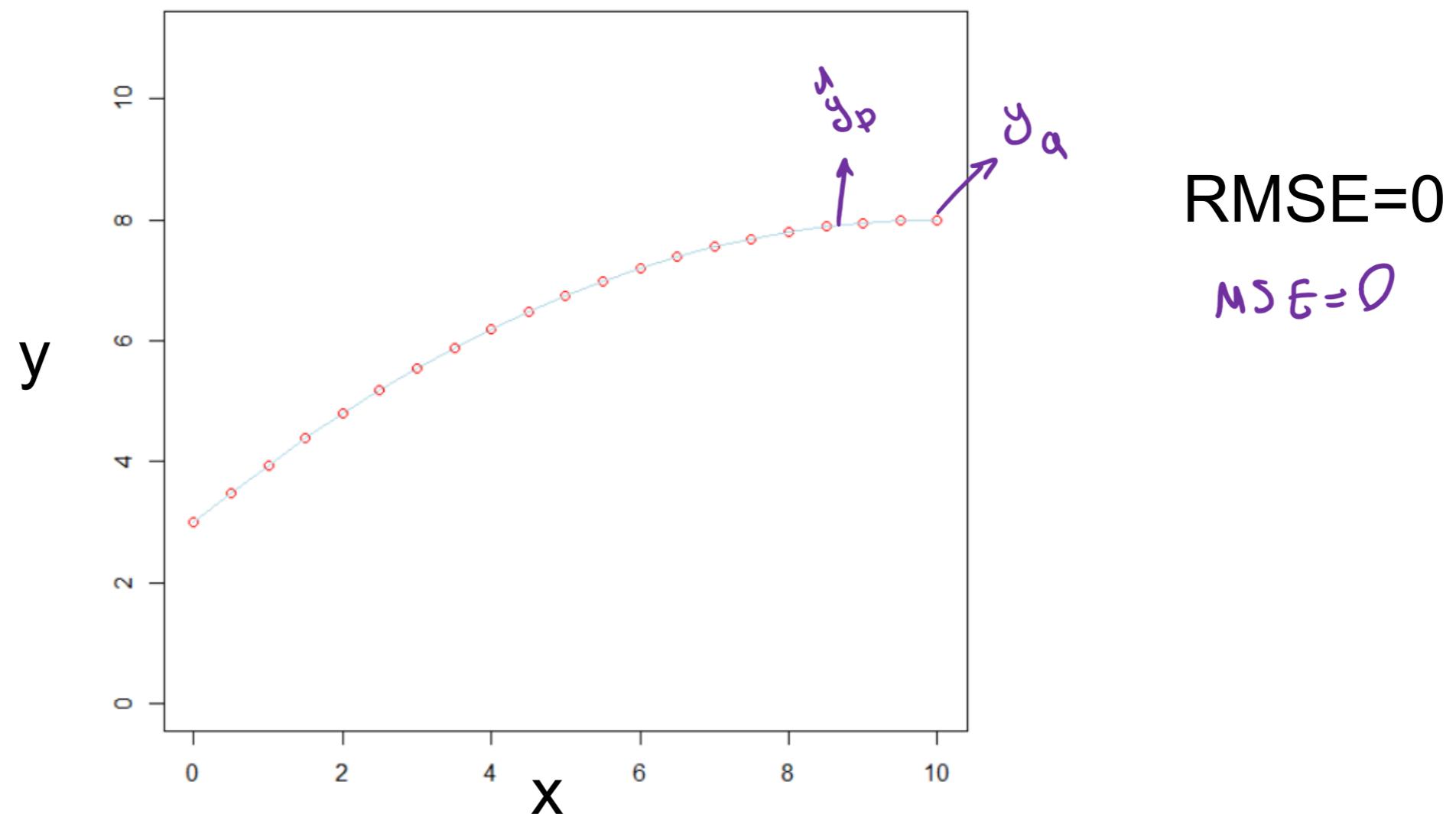
$$z = \{1, x, x^2, \dots, x^d\} \in R^d \quad y = \{y^{(1)}, y^{(2)}, \dots, y^{(n)}\}$$

$$\begin{aligned}\frac{\partial L(\theta)}{\partial \theta} &= -\frac{2}{n} z^T y + \frac{2}{n} z^T z \theta = 0 \\ \Rightarrow \theta &= (z^T z)^{-1} z^T y = z^+ y\end{aligned}$$

- If we choose a different maximal degree \mathbf{d} for the polynomial, the solution will be different.

What is happening in polynomial regression?

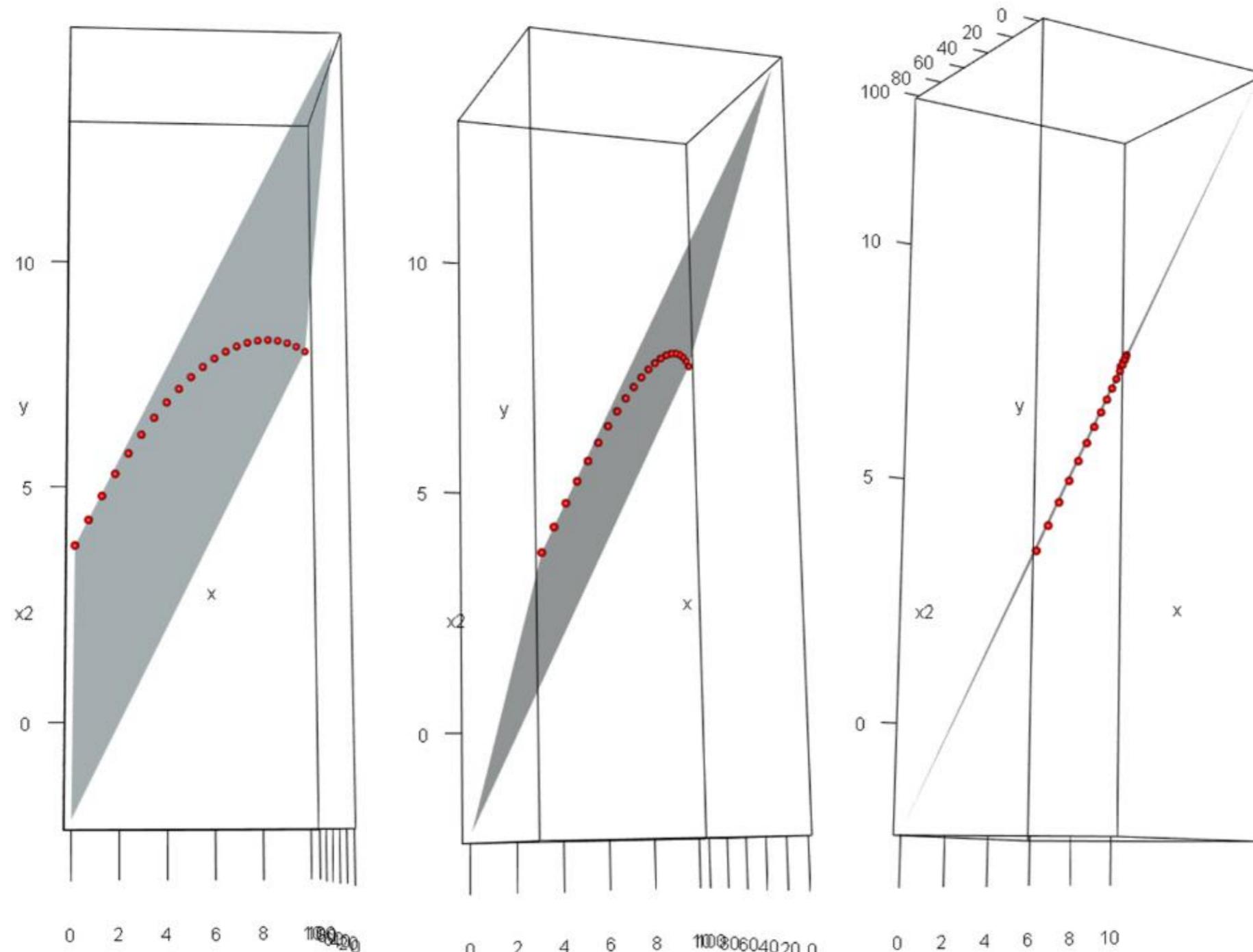
$$\left. \begin{array}{l} x = [0, 0.5, 1, \dots, 9.5, 10] \\ y = [3, 3.4875, 3.95, \dots, 7.98, 8] \end{array} \right\} \quad \begin{array}{l} f = \theta_0 + \theta_1 x + \theta_2 x^2 \\ \theta_0 = 3; \theta_1 = 1; \theta_2 = -0.5 \end{array}$$



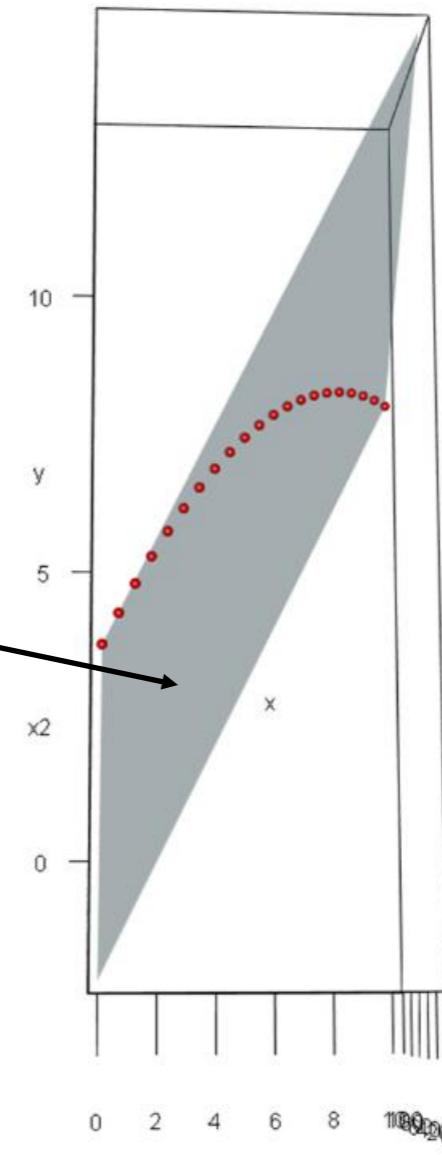
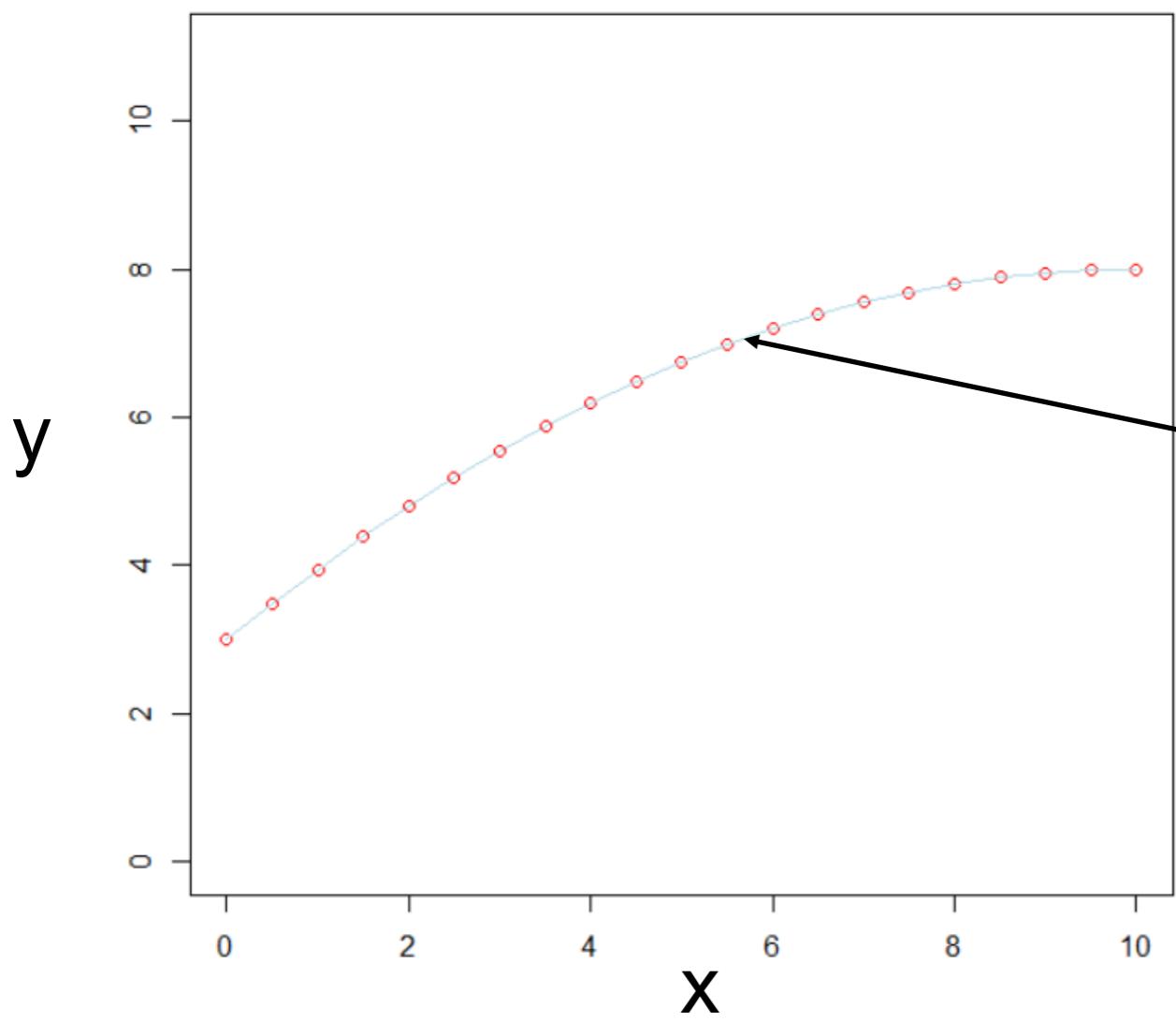
Let's add to the feature space

$$x_1 = [0, 0.5, 1, \dots, 9.5, 10] \quad x_2^2 = [0, 0.25, 1, \dots, 90.25, 100]$$

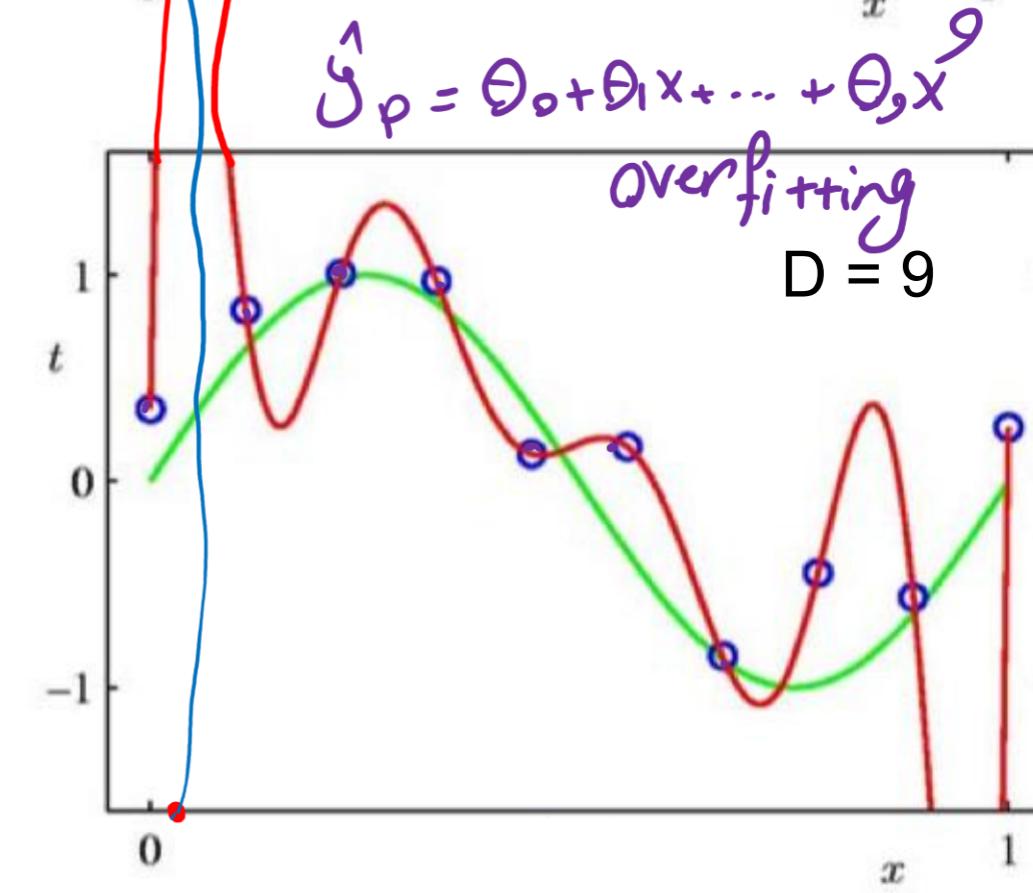
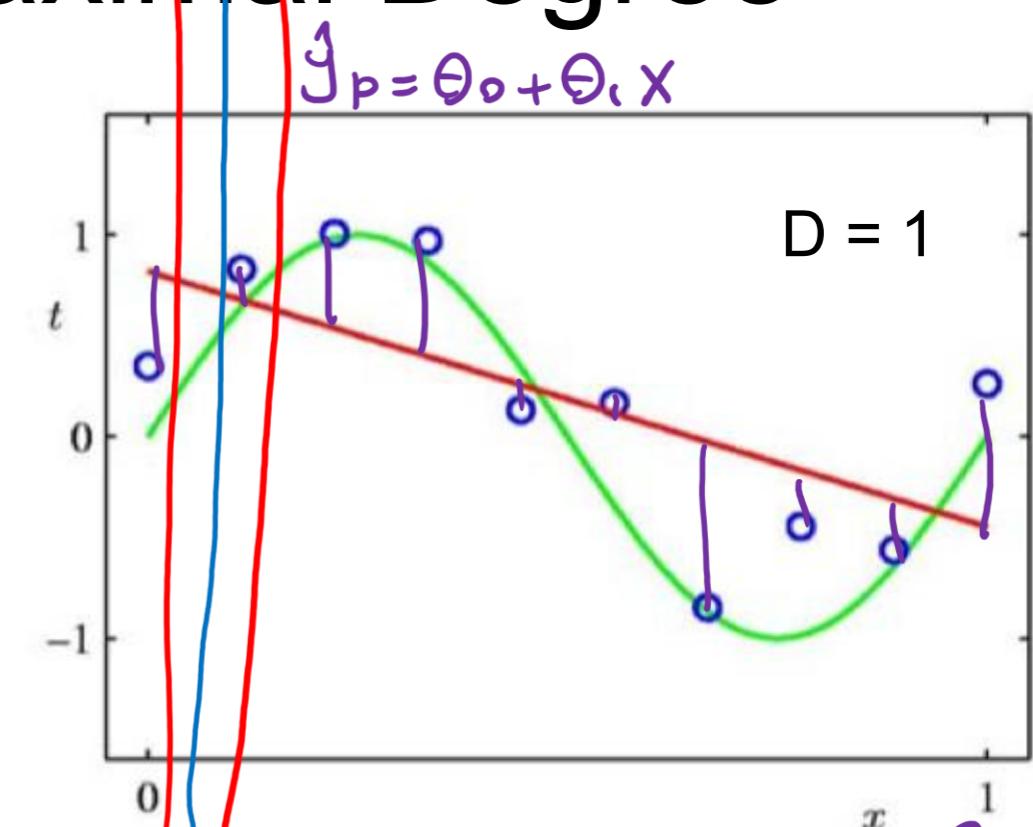
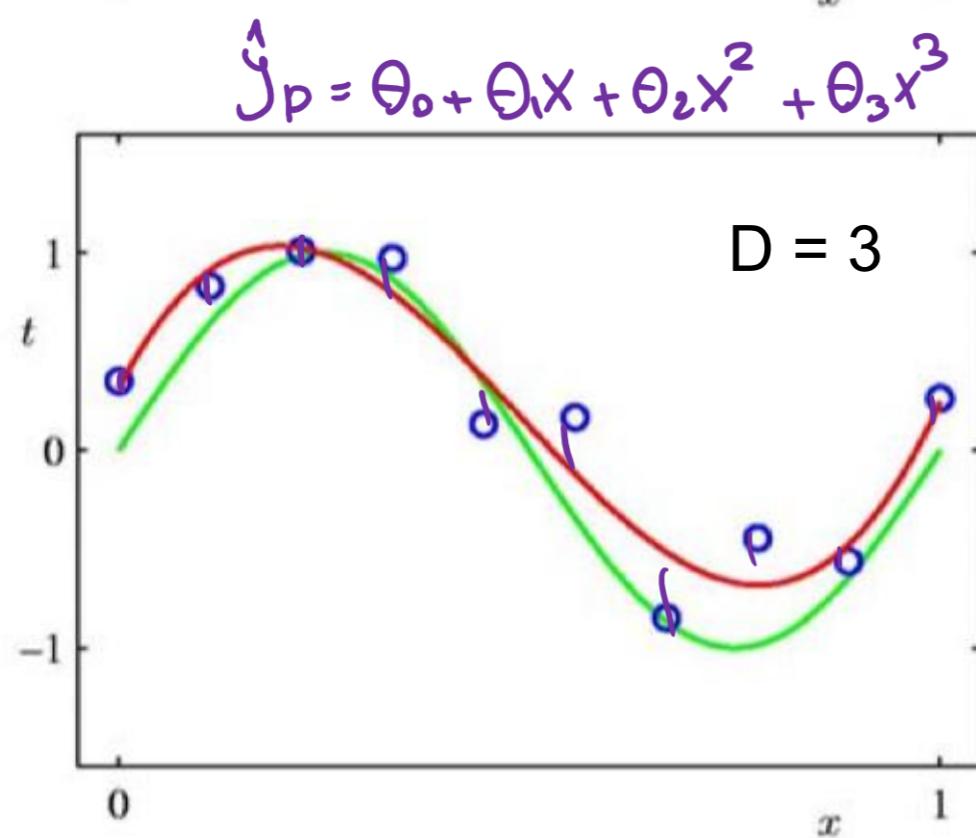
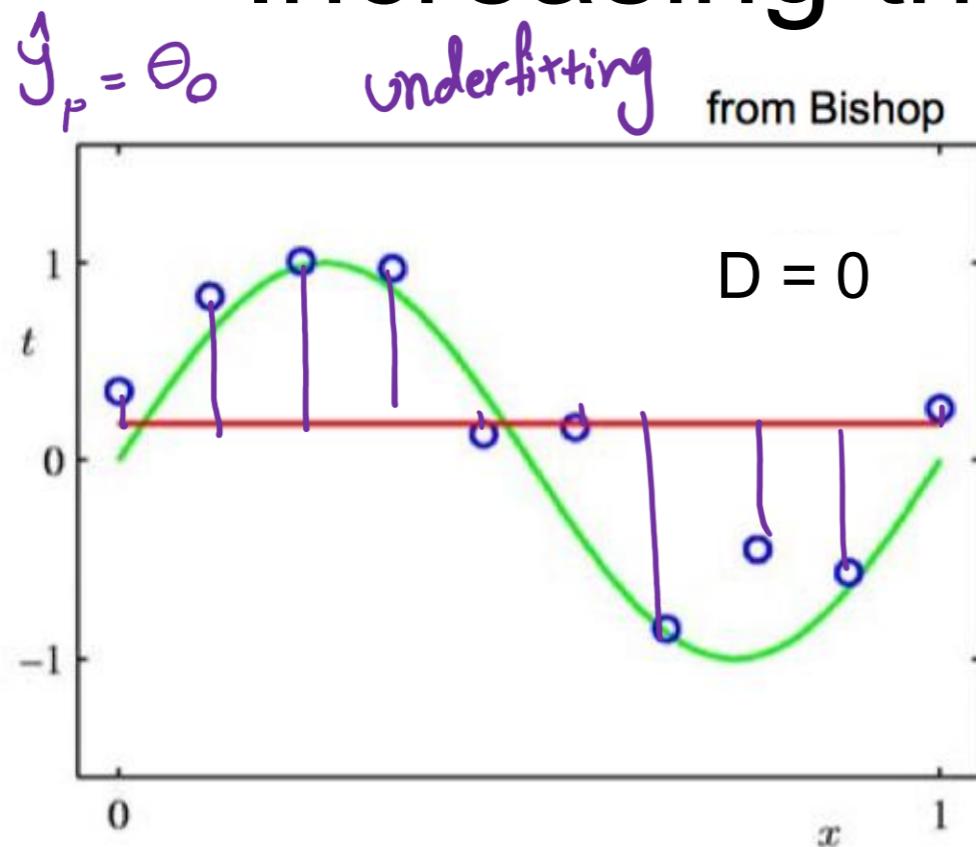
$$y = [3, 3.4875, 3.95, \dots, 7.98, 8]$$



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.

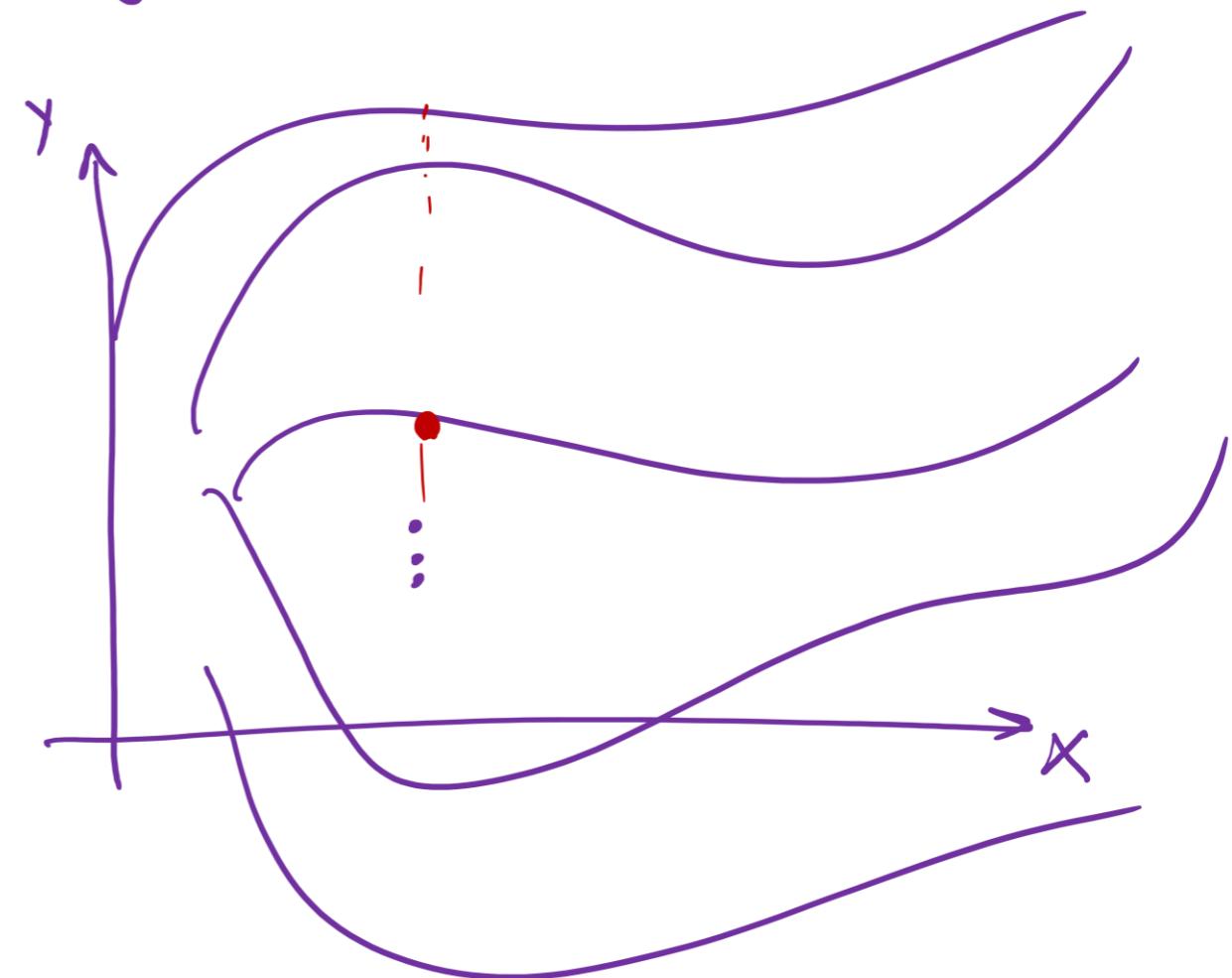
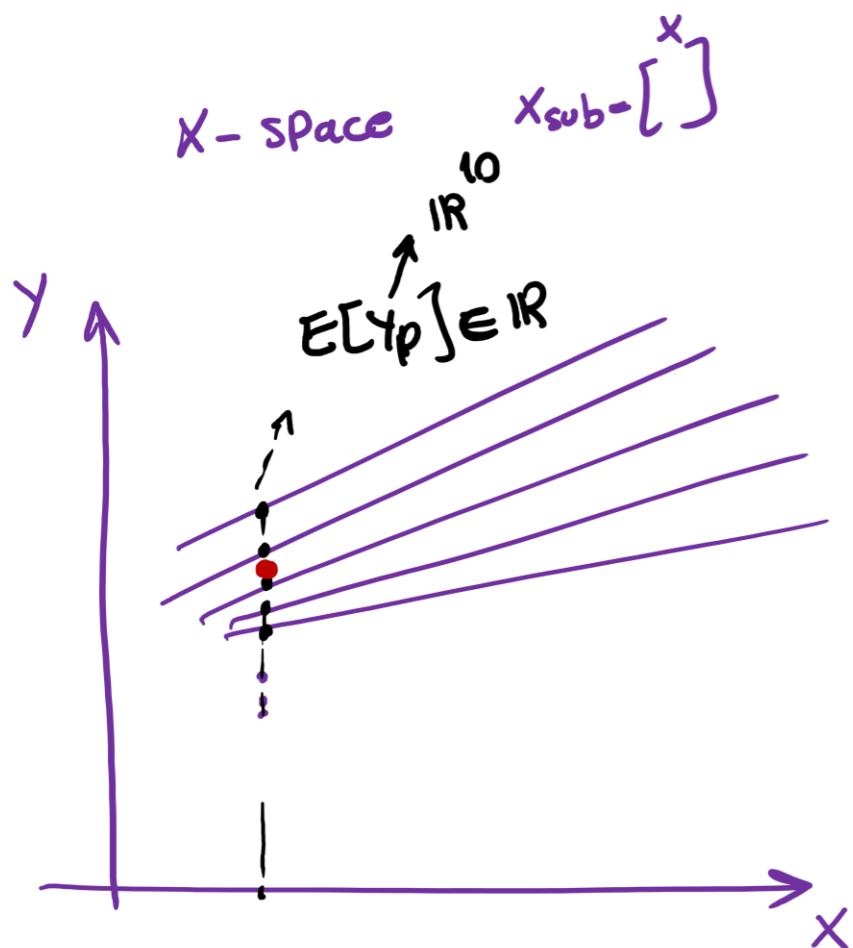


Increasing the Maximal Degree



$$X = \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}_{1000 \times 1} \rightsquigarrow X_{\text{sub}} = \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}_{100 \times 1} \quad \textcircled{1}$$

$$X_{\text{sub}} = \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}_{10 \times 1} \quad \dots \quad X_{\text{sub}} = \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix}_{1 \times 1} \quad \textcircled{10}$$



Bias-Variance Trade off

[Animation](#)

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

$$L(\theta) = \frac{1}{n} \sum_{i=1}^N (y^{(i)} - x^{(i)}\theta)^2 = E[(y_a - y_p)^2] = \text{bias}^2 + \text{Variance}$$

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

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

$$\cancel{E[(y_a - E[y_p])^2]} + \cancel{E[(E[y_p] - y_p)^2]} + 2 E[y_a - E[y_p]] E[E[y_p] - y_p]$$
$$\underline{(y_a - E[y_p])^2} + \frac{1}{N} \sum_{i=1}^N (E[y_p] - y_p)^2 + \dots \dots (E[\cancel{E[y_p]}] - E[y_p])$$

bias² + Variance

Bias-Variance Trade off

[Animation](#)

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

$$L(\theta) = \frac{1}{n} \sum_{i=1}^N (y^{\{i\}} - x^{\{i\}}\theta)^2 = E[(y_a - y_p)^2]$$

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

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

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

$$= [Bias]^2 + Variance$$

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

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

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

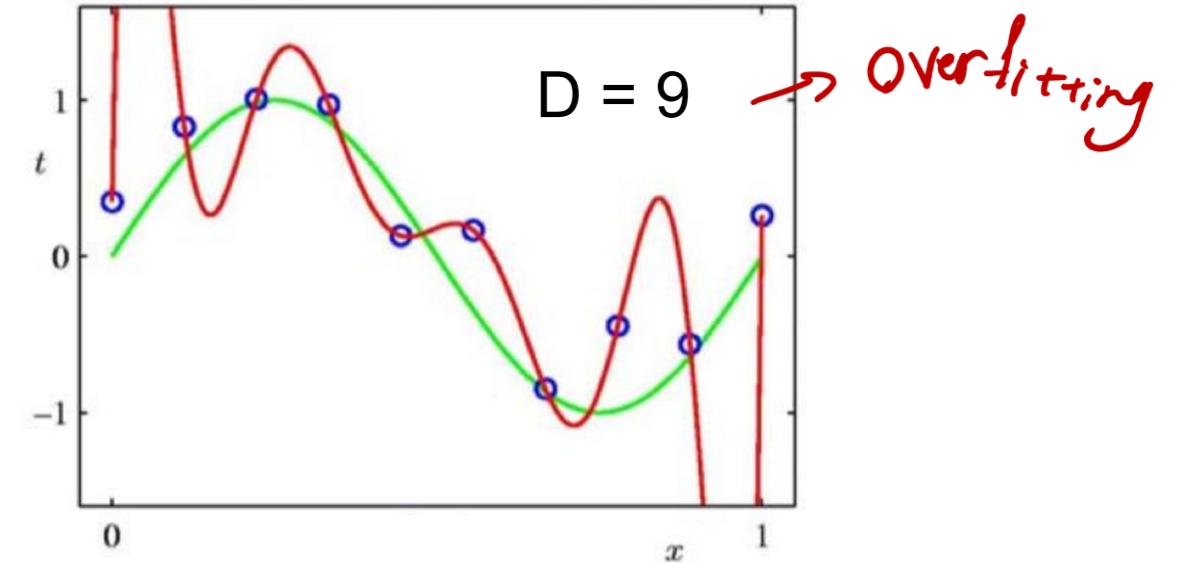
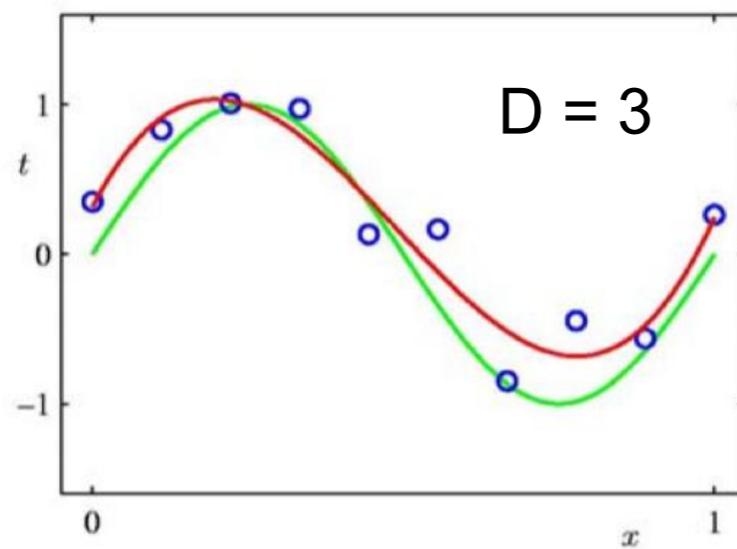
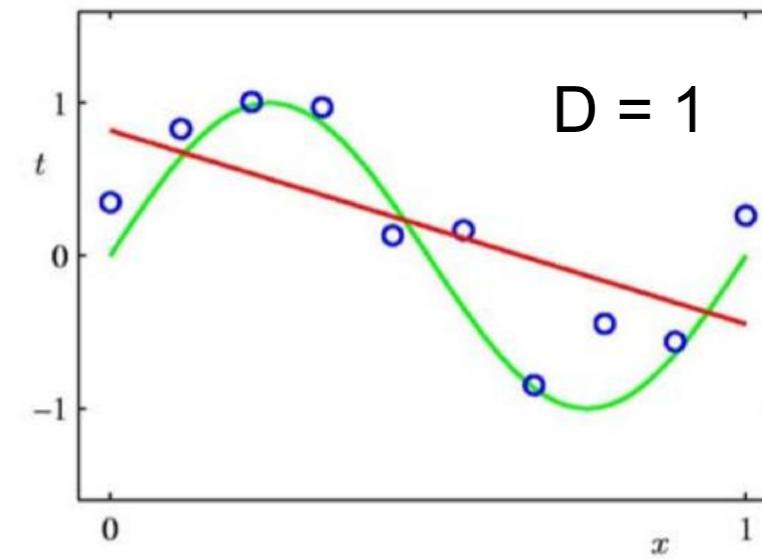
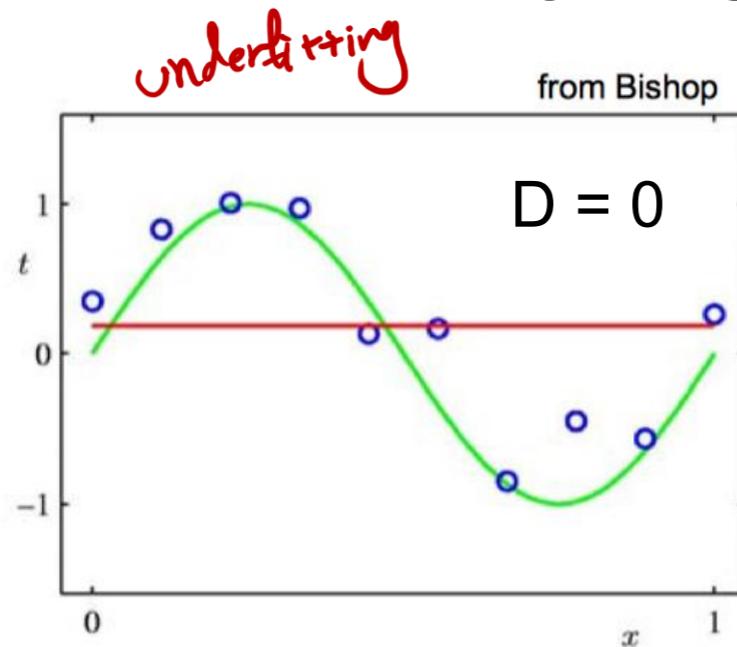
$$E[2(y_a - E[y_p])(E[y_p] - y_p)]$$

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

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

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

Which One is Better?



- 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