

Linear Regression

Mahdi Roozbahani
Georgia Tech

The Moment


© Yongqing Bao



Incoming ML HW3

You; after HW2

Outline

- Supervised Learning ← 
- Linear Regression
- Extension

Training data X & $Y \Rightarrow 80\% \rightsquigarrow$ Training data $20\% \rightsquigarrow$ Testing data

Supervised Learning: Overview

70% training data 10% validation 20% testing

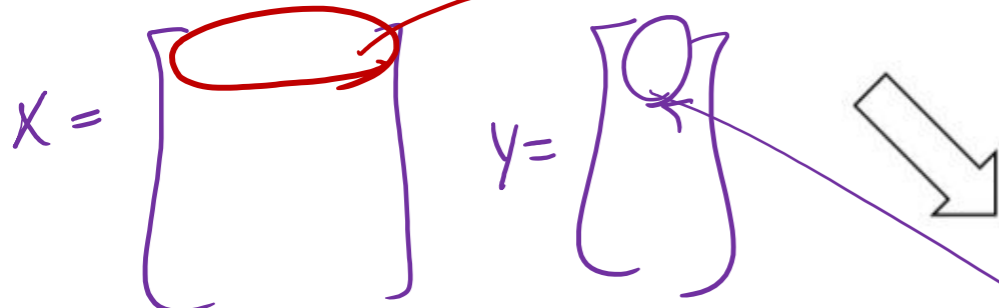
Functions \mathcal{F}

Training data

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

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

$\underbrace{\hspace{2cm}}_{n \times d} \quad \underbrace{\hspace{2cm}}_{n \times 1}$



LEARNING

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

actual \hat{y}_p



Learning machine

PREDICTION

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

New data

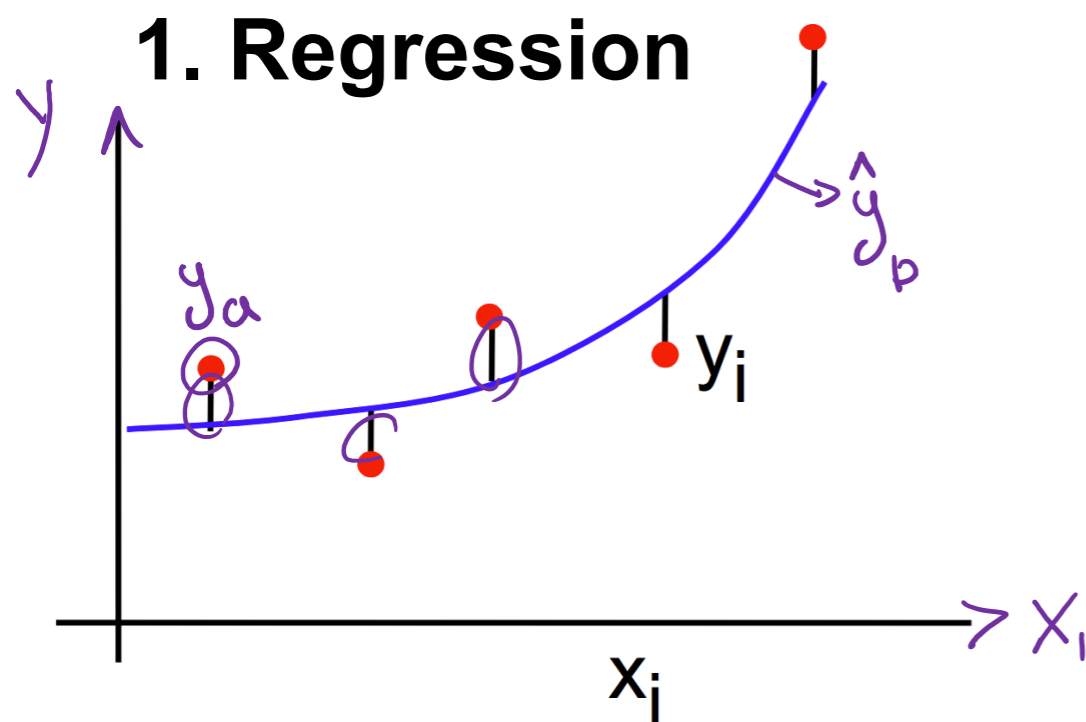
$$x$$

Supervised Learning: Two Types of Tasks

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

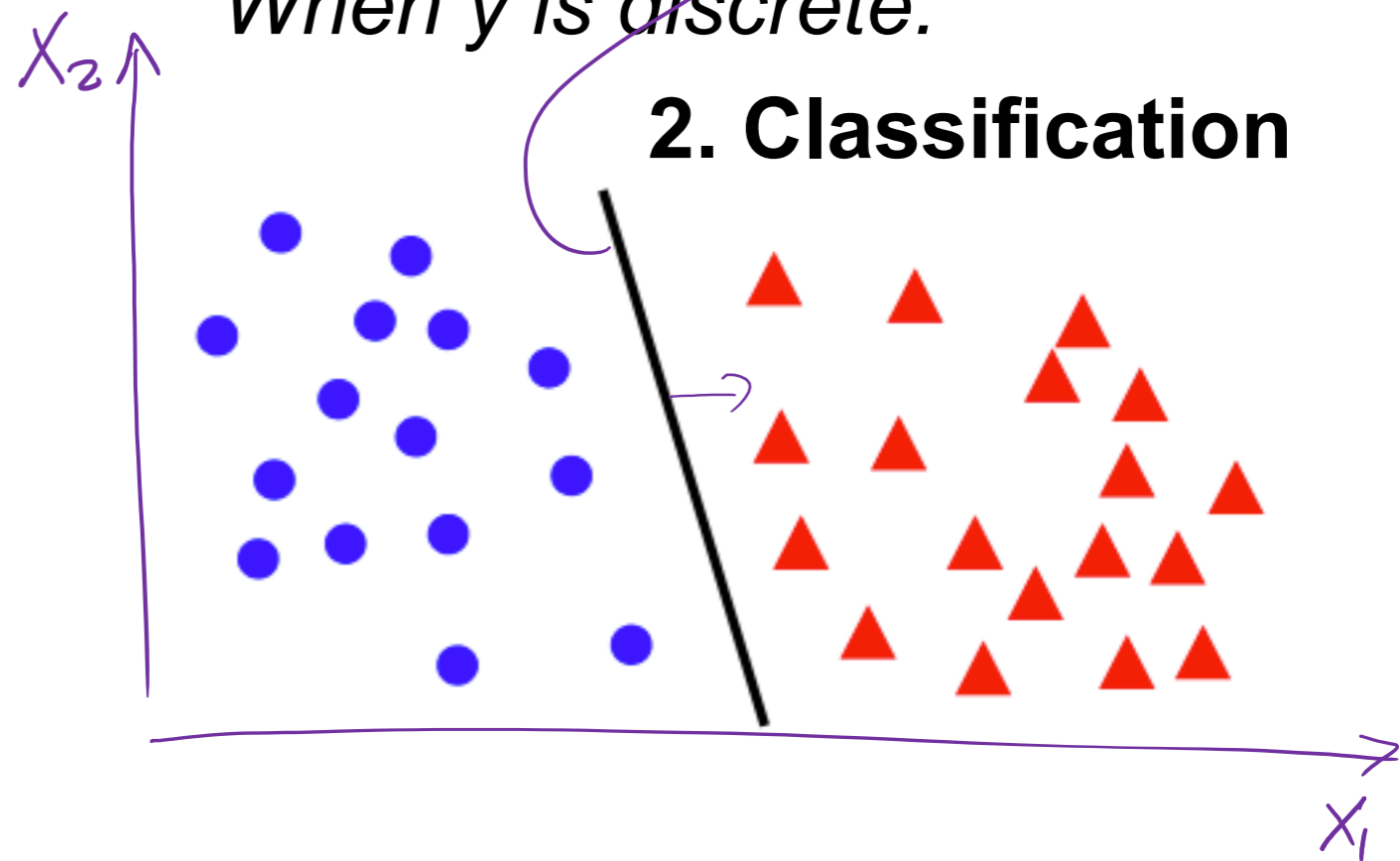
Learn: a function $f(\mathbf{x}) : y = f(\mathbf{x})$

When y is continuous:



Minimize the difference between y_a & \hat{y}_p

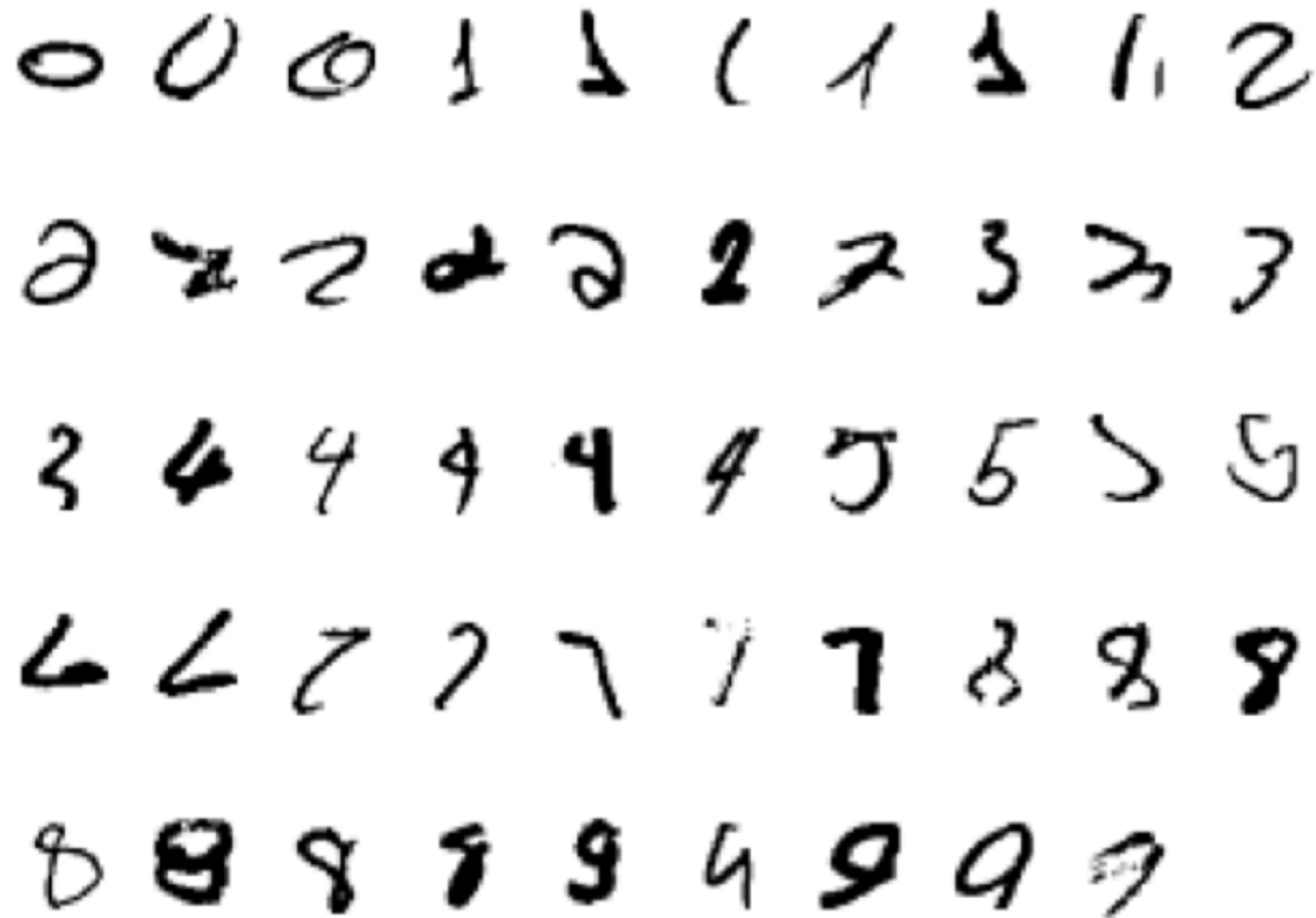
When y is discrete:



Classification Example 1: Handwritten digit recognition

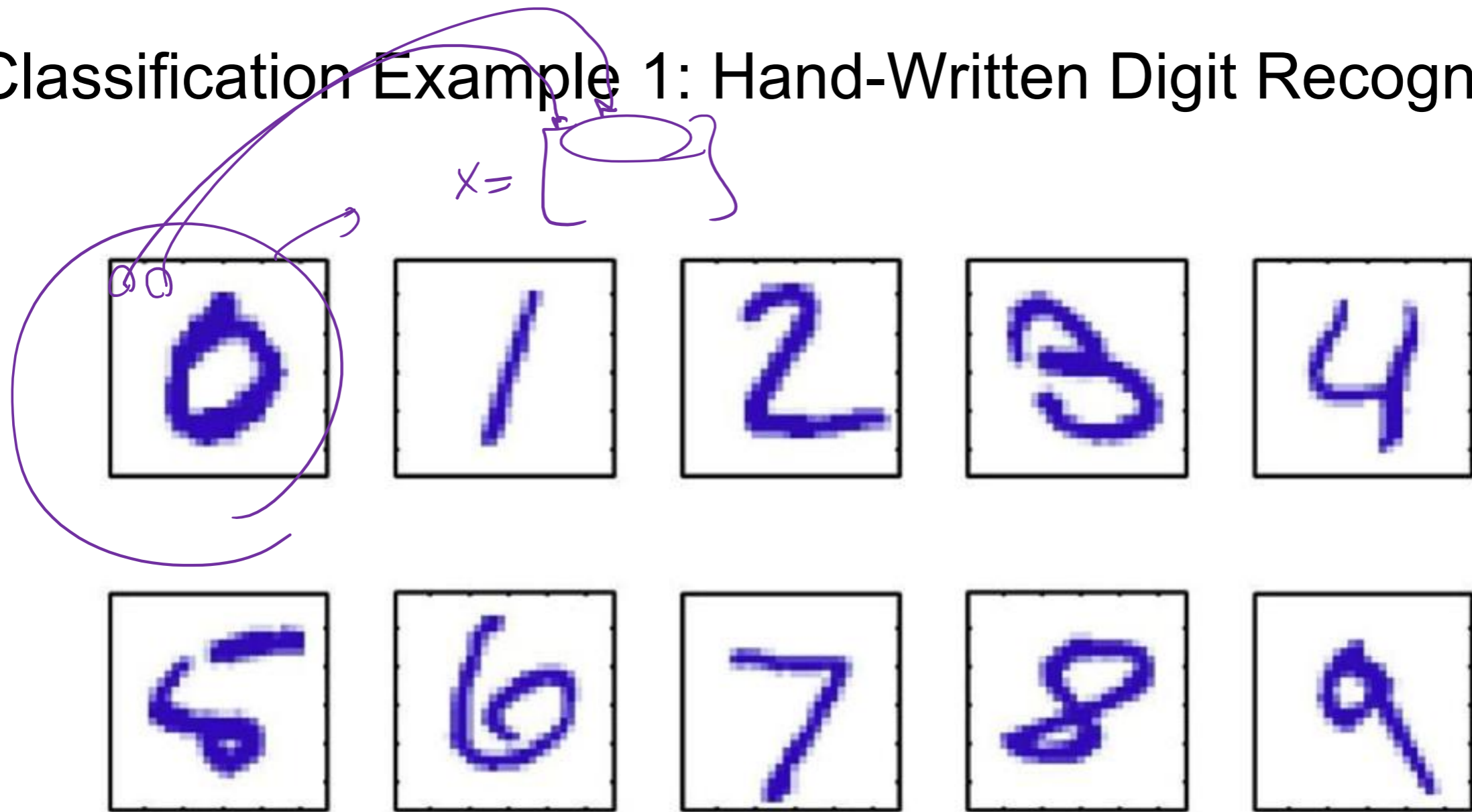
As a supervised classification problem

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

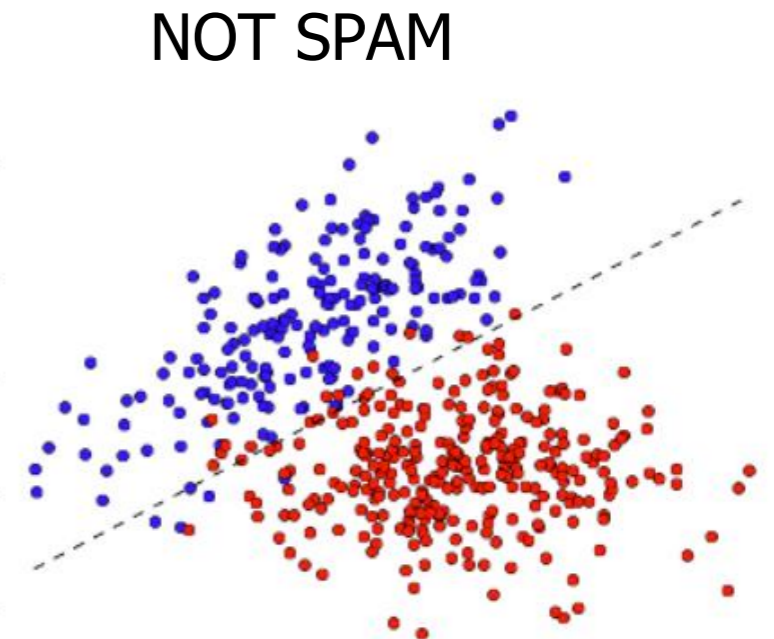
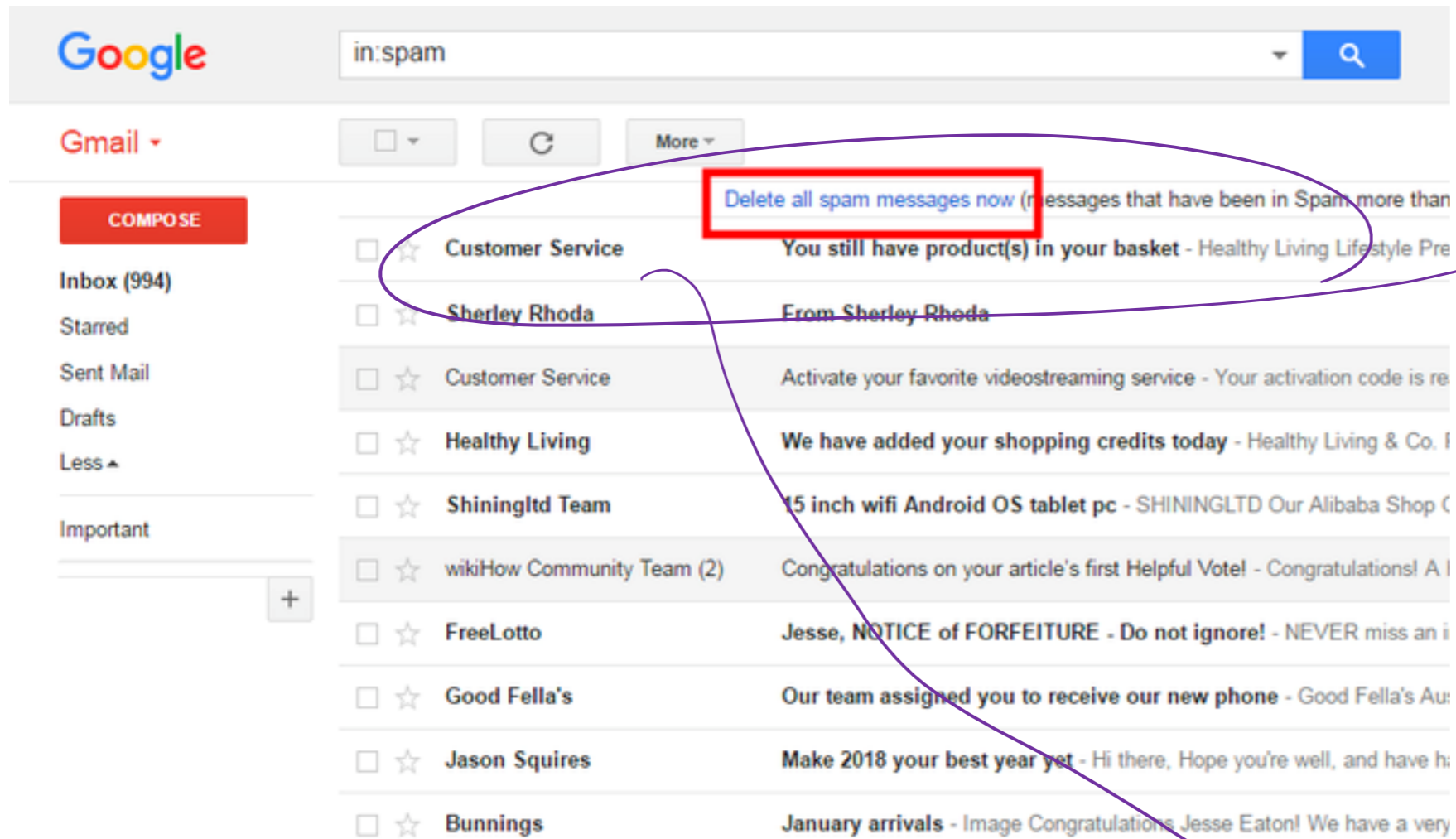
A classification problem

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

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

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

Classification Example 2: Spam Detection



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

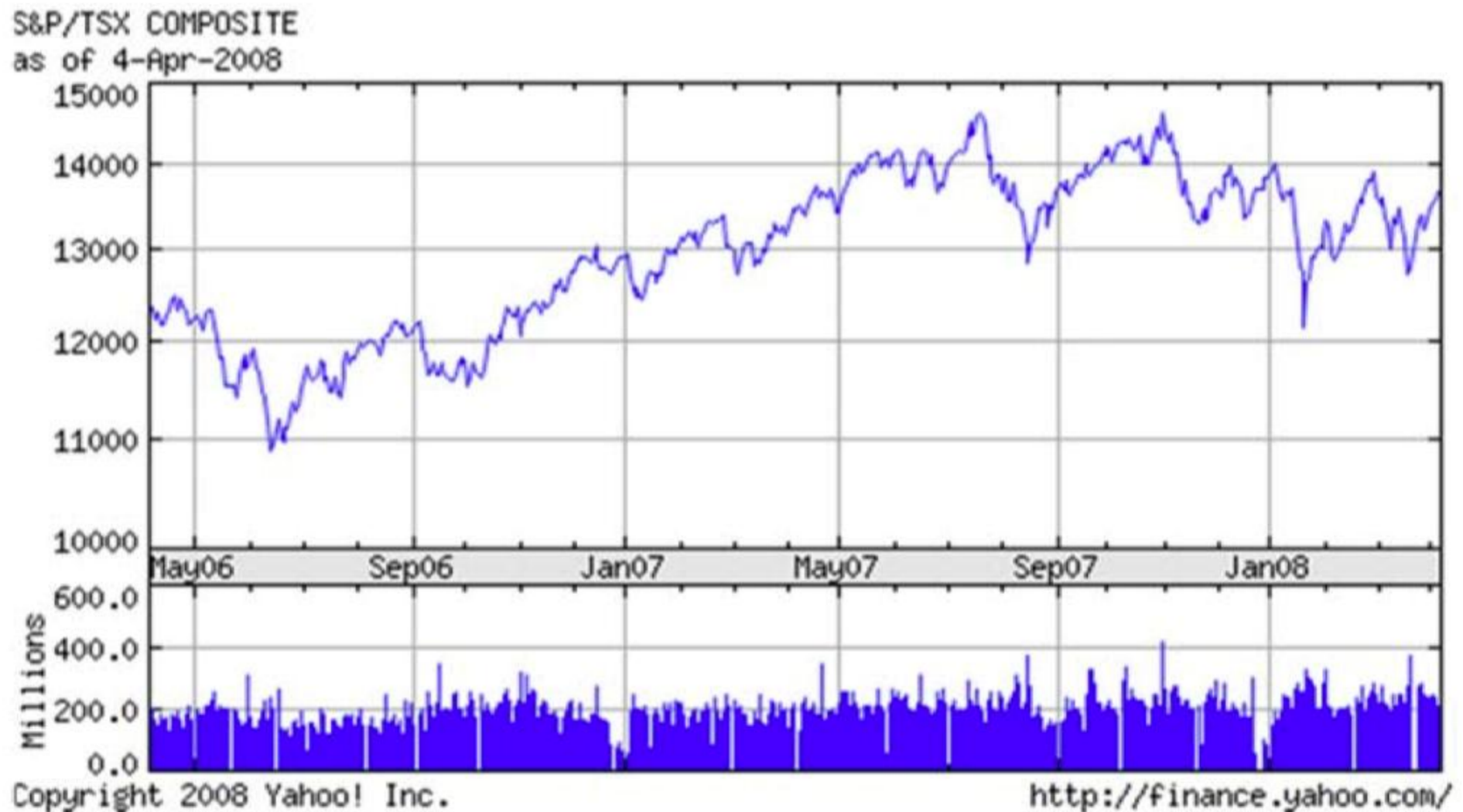
- Suppose you are to move to Atlanta
- And you want to find the **most reasonably priced** apartment satisfying your **needs**:

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



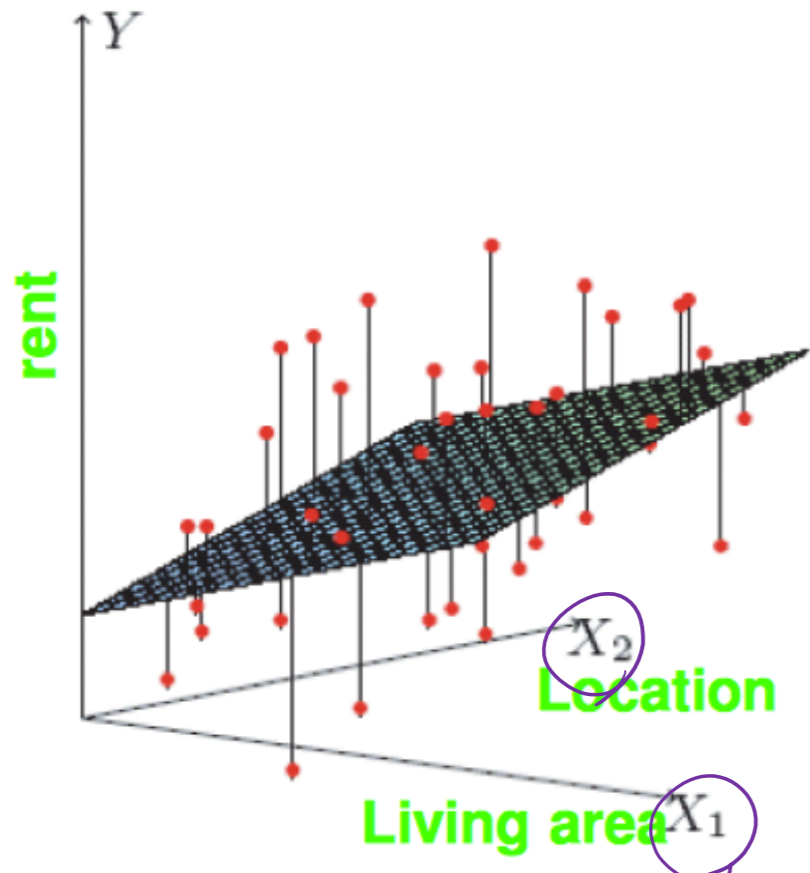
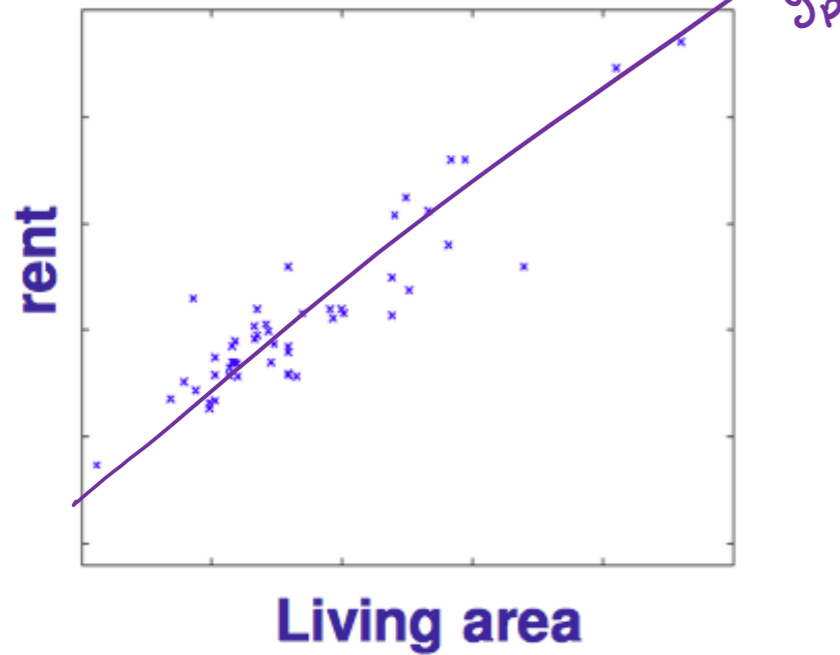
- Task is to predict stock price at future date

A regression problem

$$f(x^{(i)}) = m x^{(i)} + \underline{b}$$

$$f(x^{(i)}) = \underbrace{\Theta_0}_{\text{bias term}} + \underbrace{\Theta_1}_{\text{parameter}} x_1$$

the parameter that defines the behavior of feature x_1



• Features: Parameter

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

$$x^{(i)} = [1, x_1^{(i)}, x_2^{(i)}, \dots, x_d^{(i)}]$$

$$\Theta = \begin{bmatrix} \Theta_0 \\ \Theta_1 \\ \vdots \\ \Theta_d \end{bmatrix}_{d+1 \times 1}$$

• Target:

- Rent $f(x^{(i)}) = x^{(i)} \cdot \Theta$
- Denoted as y

$$f(x) = \hat{y}_p = x \cdot \Theta$$

• Training set:

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

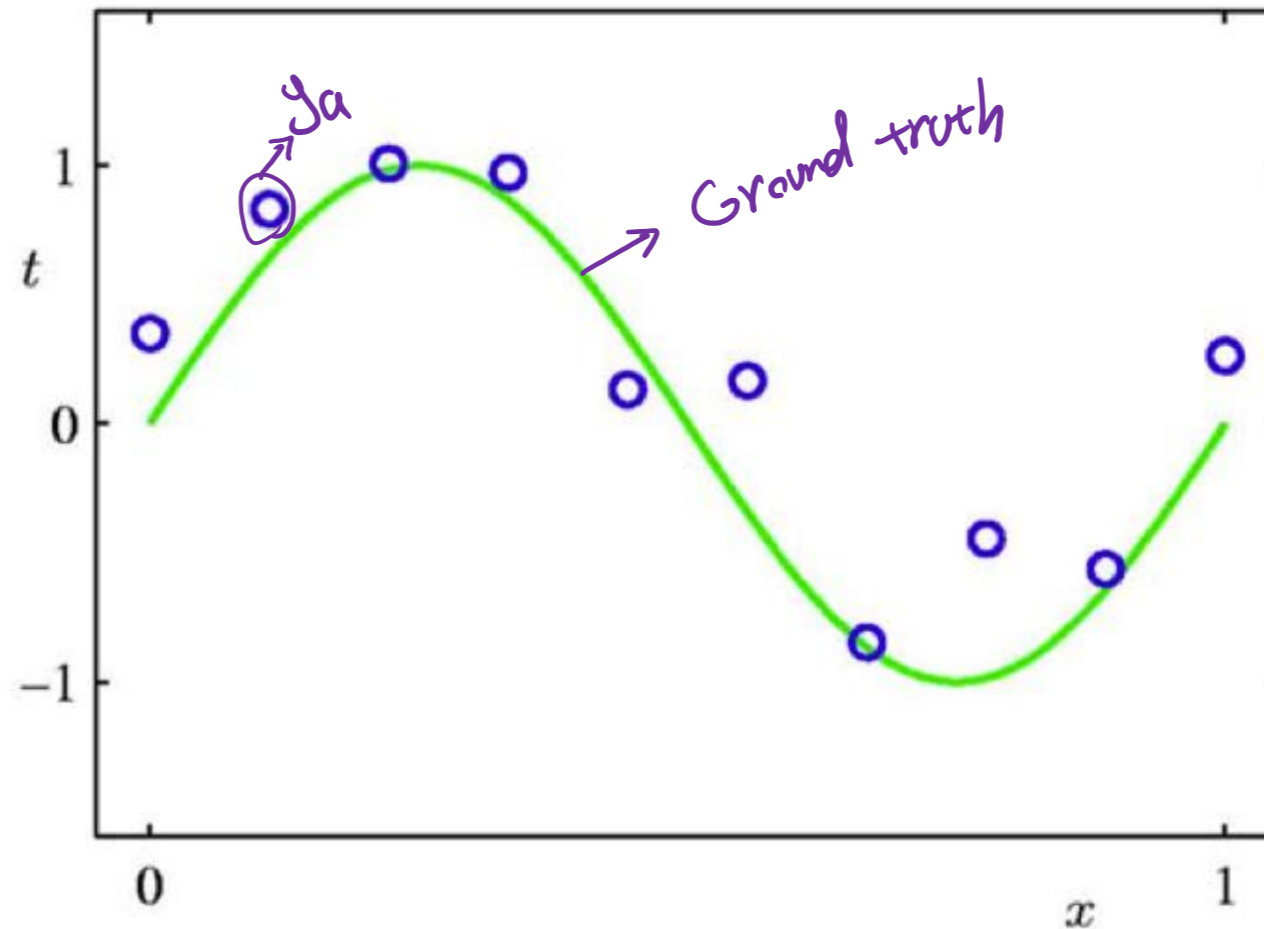
$$y = \{y^{(1)}, y^{(2)}, \dots, y^{(n)}\}$$

LCF

$$f(x^{(i)}) = \Theta_0 + \Theta_1 x_1 + \Theta_2 x_2$$


LCF \rightarrow linear combination of features

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

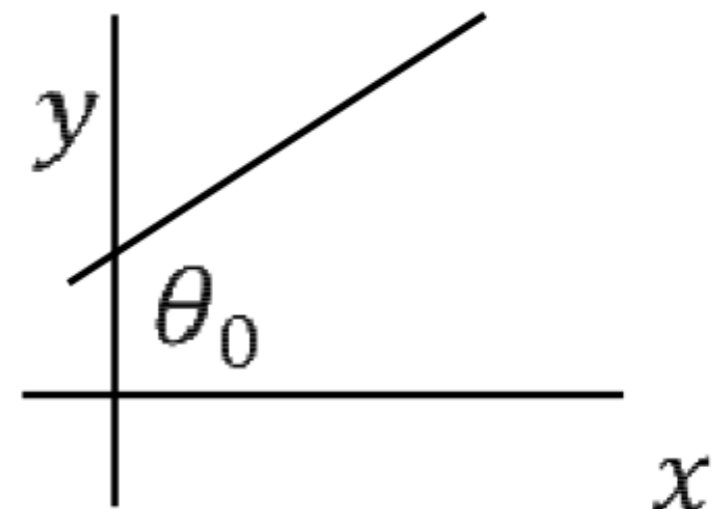
- Assume y is a linear function of x (features) plus noise ϵ

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

- 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

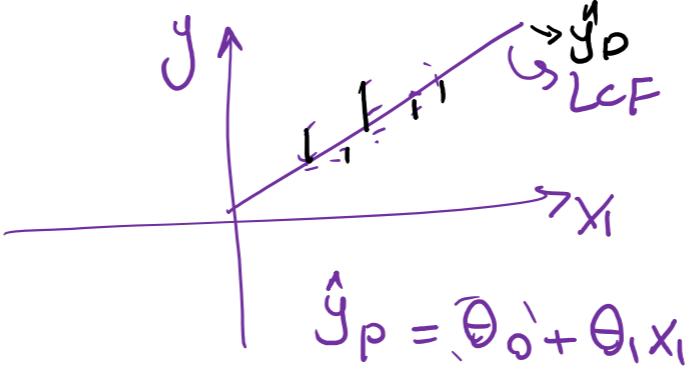
$$f(x) = \hat{y}_p = x\theta \rightarrow LCF$$

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



LCF $f(x^{[i]}) = \hat{y}_p = \theta_0 + \theta_1 x_1^{[i]} + \dots + \theta_d x_d^{[i]} \Rightarrow \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}_{d+1 \times 1}$

$f(x) = \hat{y}_p = X \cdot \theta \rightsquigarrow$ LCF



Objective function $L(\theta) = E(\theta) = \frac{1}{N} \left(\sum_{i=1}^N \left(\underline{y_a^{[i]}} - \underline{\hat{y}_p^{[i]}} \right)^2 \right)$

$\underline{\hat{y}_p^{[i]}} = X^{[i]} \cdot \theta$

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

~~$E[\theta]$~~

$f(x) = -x^2$ $f(x) = -2x$

$|y^{[i]} - x^{[i]}\theta|$ $y = |x| \checkmark$ $y = x^2 \cup$

Least Mean Square Method

$x^{[t+1]} \leftarrow x^{[t]} + \alpha \frac{\partial f(x)}{\partial x} = -2x$

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

Training

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

Annotations: $L(\theta)$ is circled and labeled "Loss" and "E(L(theta)) error". $y^{[i]}$ is labeled "actual y_a 1x1". $x^{[i]}\theta$ is labeled " \hat{y}_p 1x1". A note above says " $1 \times d+1$ " with an arrow pointing to " $d+1 \times 1$ ". A note to the right says "Convex".

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

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

GD

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

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

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

Matrix form

$$x = \begin{bmatrix} 1 & x_1^{\{1\}} & \dots & x_d^{\{1\}} \\ 1 & x_1^{\{2\}} & \ddots & x_d^{\{2\}} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_1^{\{n\}} & \dots & x_d^{\{n\}} \end{bmatrix} \quad y = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix} \quad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$$

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

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

$$x\theta = \begin{bmatrix} \theta_0 + \theta_1 x_1^{\{1\}} + \theta_2 x_2^{\{1\}} + \dots + \theta_d x_d^{\{1\}} \\ \theta_0 + \theta_1 x_1^{\{2\}} + \theta_2 x_2^{\{2\}} + \dots + \theta_d x_d^{\{2\}} \\ \vdots \\ \theta_0 + \theta_1 x_1^{\{n\}} + \theta_2 x_2^{\{n\}} + \dots + \theta_d x_d^{\{n\}} \end{bmatrix} \begin{matrix} \rightarrow \hat{y}_p^{(1)} \\ \\ \\ \hat{y}_p^{(n)} \end{matrix}$$

$n \times 1$

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:

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

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

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

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

Handwritten annotations:
 X is circled in purple.
 $(X^T X)^{-1}$ is circled in purple.
 X^T is circled in purple.
 y is circled in purple.
 X^+ is circled in purple.
Dimensions are written below: $d+1 \times 1$ for θ , $d+1 \times n$ for $X^T X$, $n \times d+1$ for X^T , $n \times 1$ for y , and $d+1 \times n$ for X^+ .

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

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

$X_{n \times d}$

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

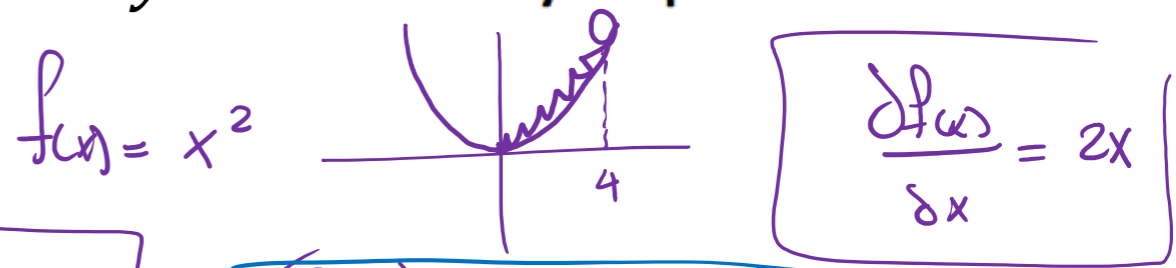
$$X^T X = \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \quad d \times n \quad \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \quad n \times d = \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \quad d \times d$$

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.

$f(x) = x^2$ $f(x) = 10x^2$ $f(x) = -x^2$

Alternative Way to Optimize

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



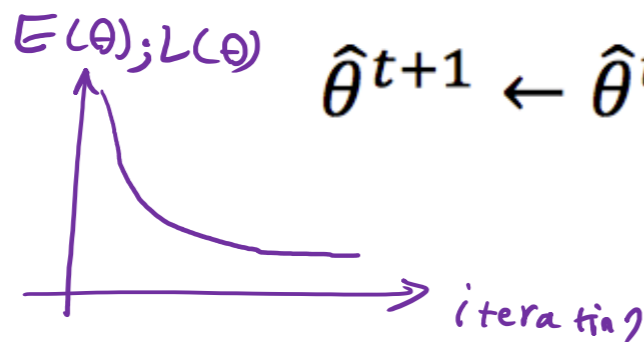
$$\frac{\partial L(\theta)}{\partial \theta} = -\frac{2}{n} \sum_{i=1}^N (x^{i})^T (y^{i} - x^{i} \theta)$$

$$x^{t+1} \leftarrow x^t - \alpha \frac{\partial f(x)}{\partial x} = (x^t) - \alpha (2x)$$

$$\theta \leftarrow \theta - \alpha \frac{\partial L(\theta)}{\partial \theta} = \theta - \alpha$$

- Gradient descent

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

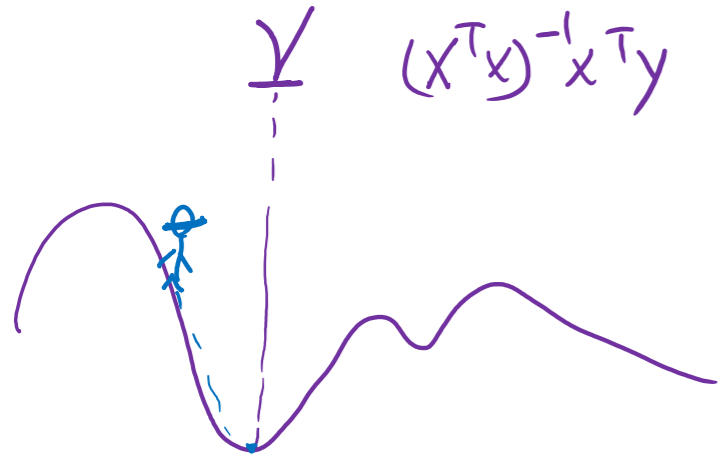


$$\hat{\theta}^{t+1} \leftarrow \hat{\theta}^t + \frac{\alpha}{n} \sum_{i=1}^N (x^{i})^T (y^{i} - x^{i} \theta) \quad \text{"FGD"}$$

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

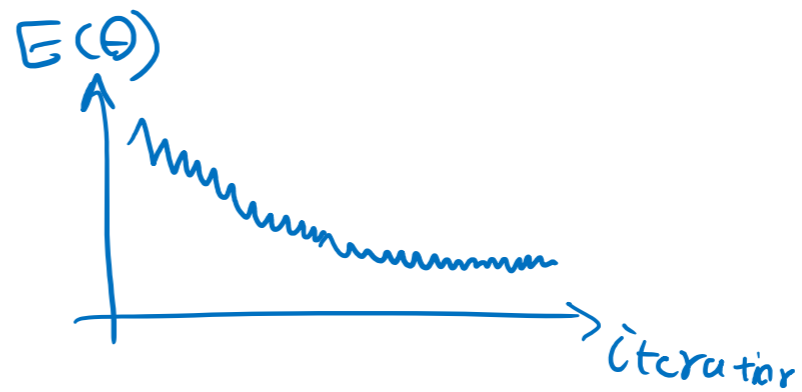
SGD

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



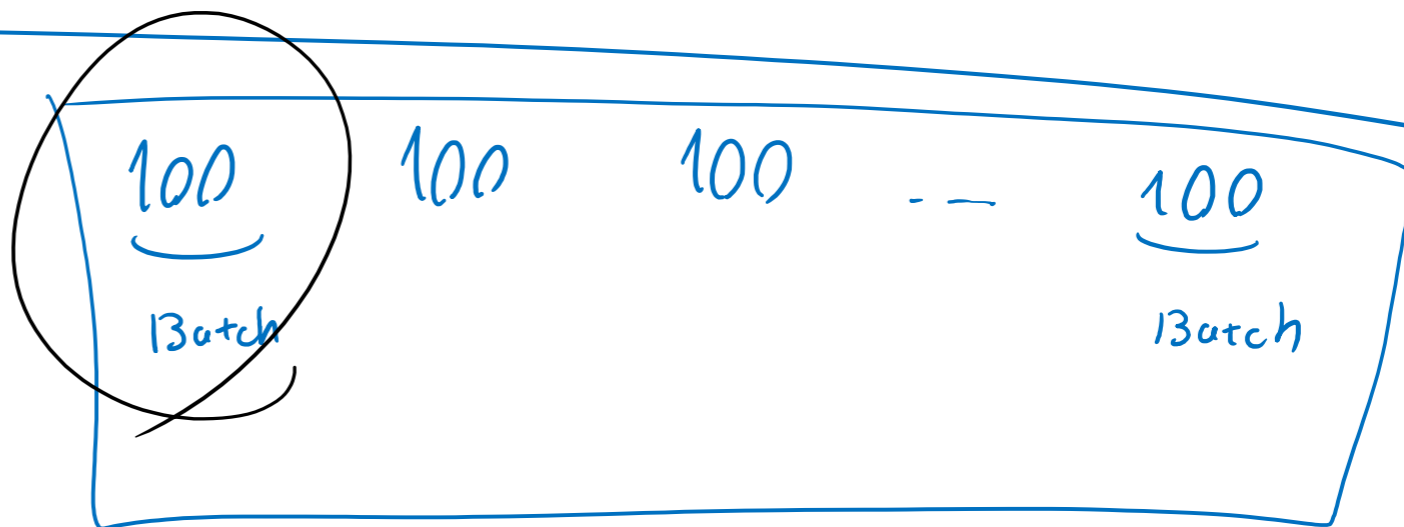
F60

SGD

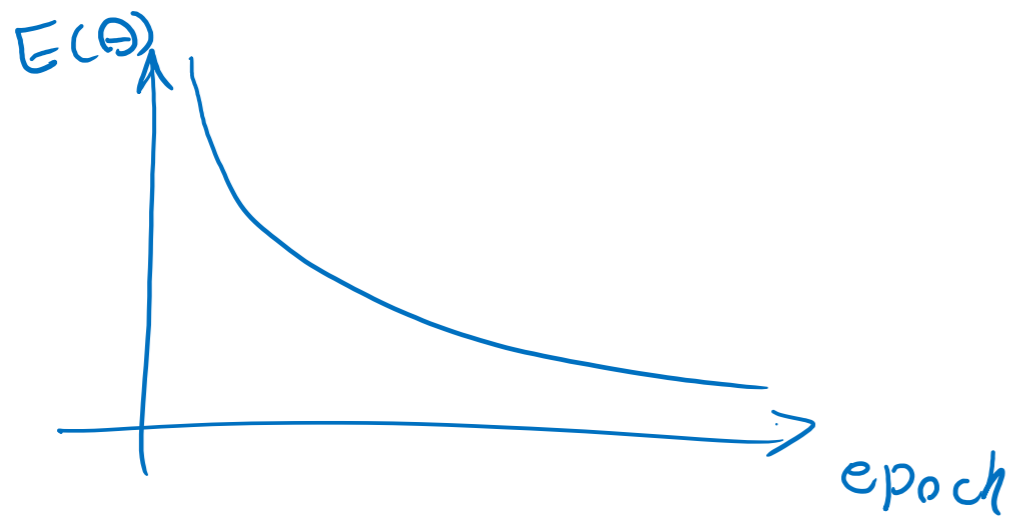


BGD

1000 datapoint



Batch



Recap

- Stochastic gradient update rule

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

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

↓
BGD

- 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

↓
FGD

- 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 ..)

Linear regression for classification

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

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

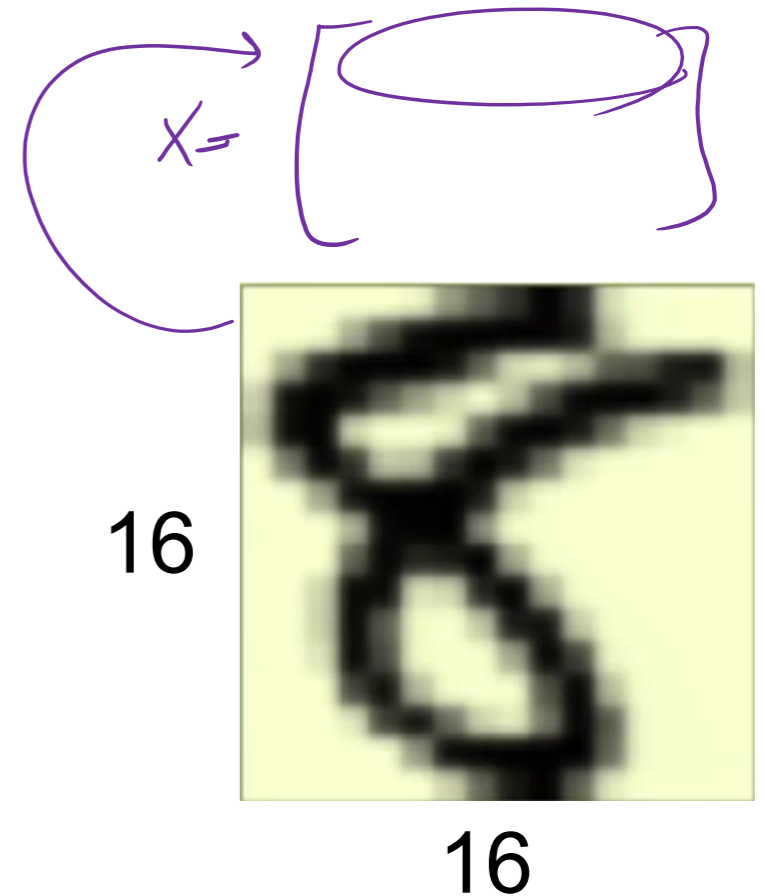
feature engineering

Extract useful information

intensity and symmetry $x = (1, x_1, x_2)$
 $\theta_0 \theta_1 \theta_2$

Sum up all the pixels = intensity

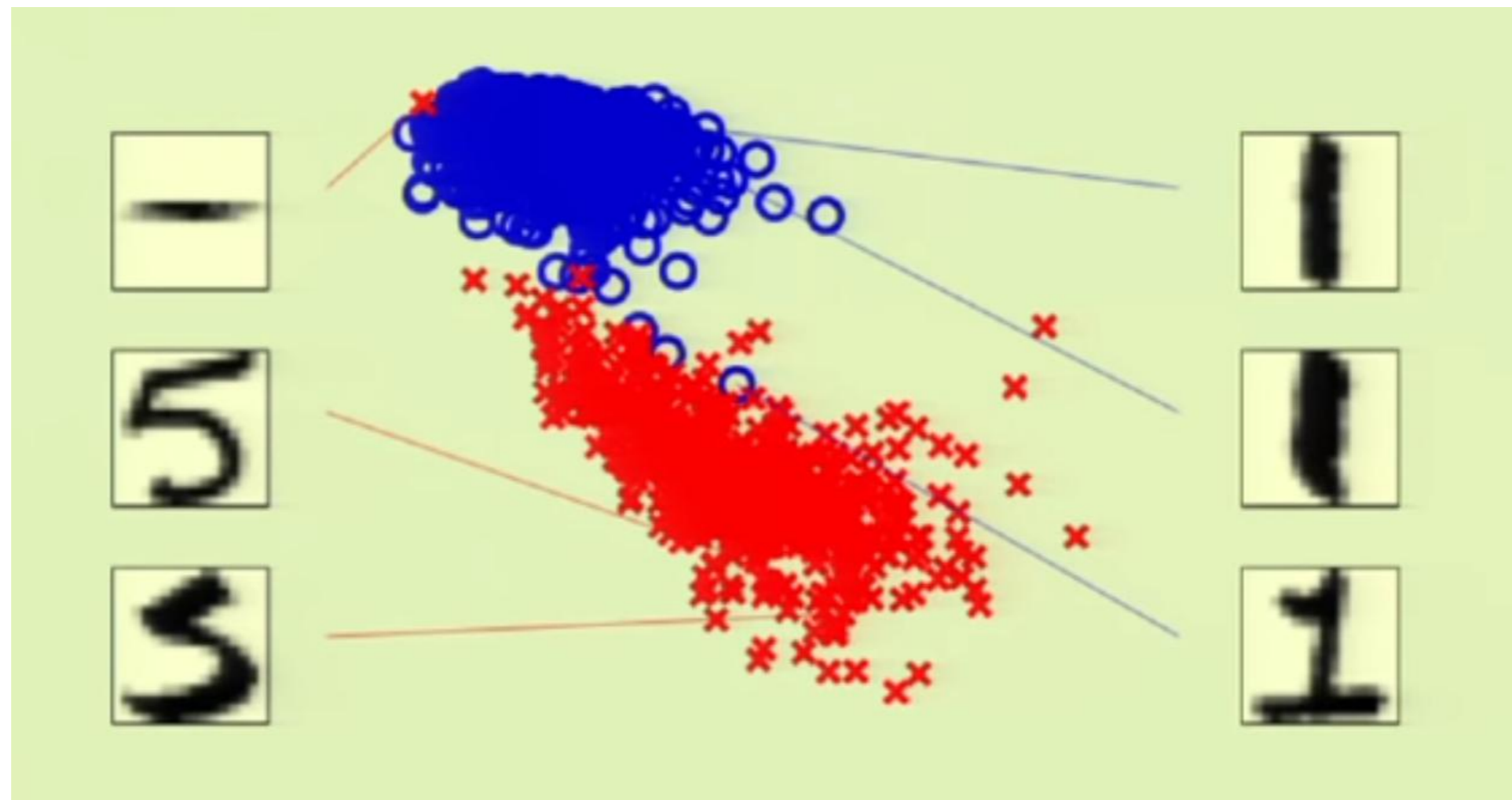
Symmetry = -(difference between flip version)



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

$$x_1 = \textit{intensity} \quad x_2 = \textit{symmetry}$$

It is almost linearly separable



symmetry

intensity

Linear regression for classification

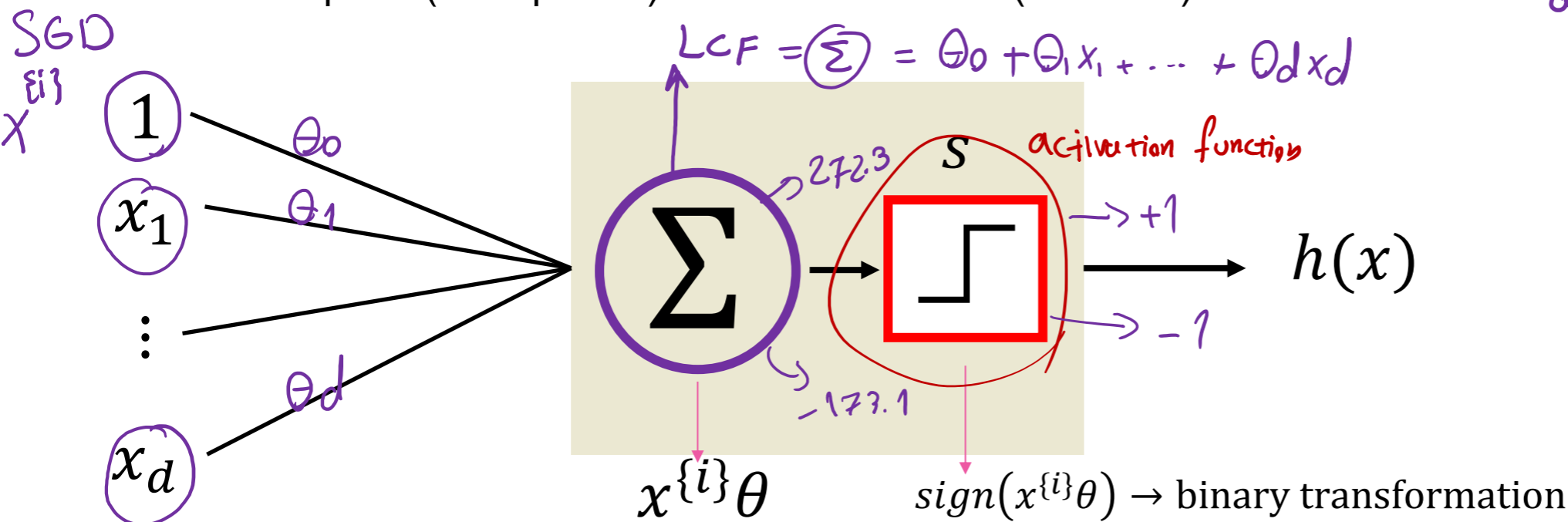
$y_p = X \cdot \theta \in \mathbb{R}$
 $-\infty, +\infty$
 LCF y_p

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

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

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

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



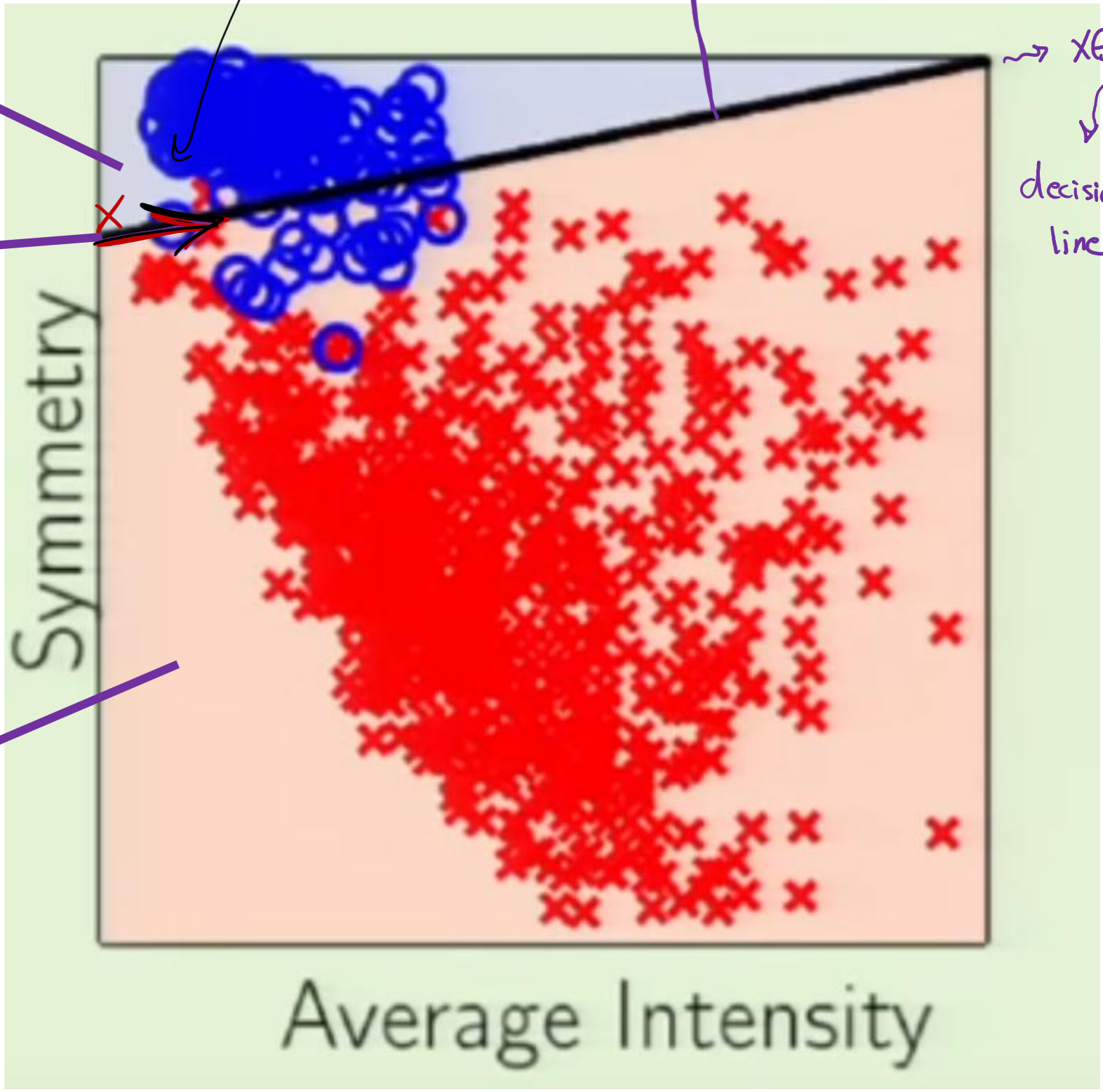
$$\Theta = (X^T X)^{-1} X^T y$$

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

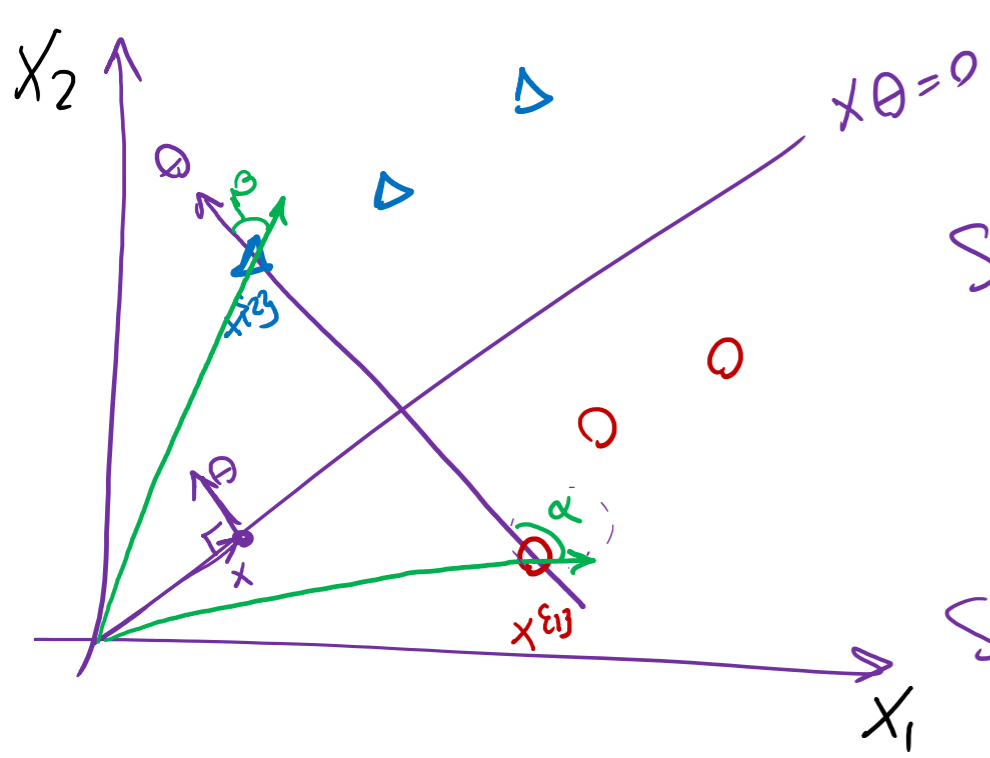
$$(X) \Theta = 0$$

0

-1



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



$$LCF = x\theta = 0$$

$$\text{Sign}(LCF(x^{\xi 1})) = x^{\xi 1} \cdot \theta = |x^{\xi 1}| |\theta| \cos \alpha = -1$$

$$\text{Sign}(LCF(x^{\xi 2})) = x^{\xi 2} \cdot \theta = |x^{\xi 2}| |\theta| \cos \beta = +1$$

5

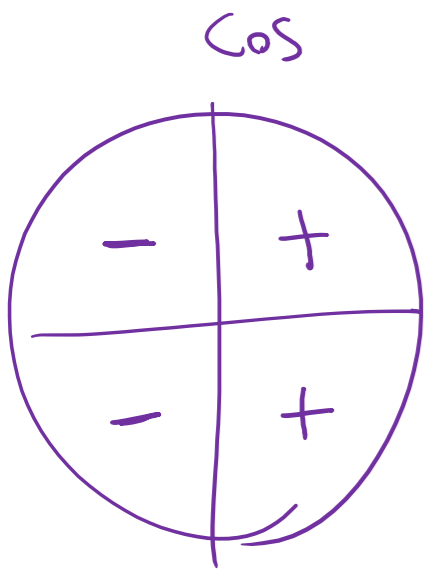
-1

+1


1

$$LCF = f(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

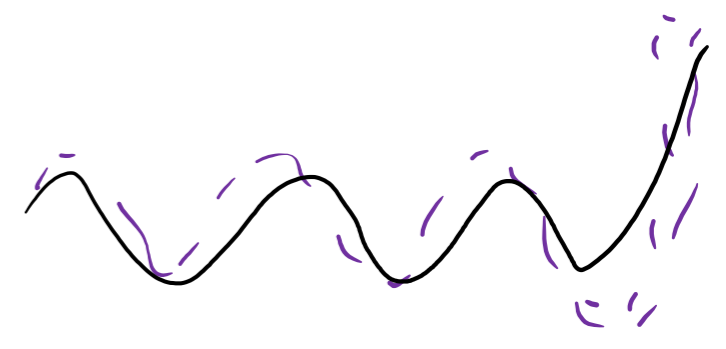
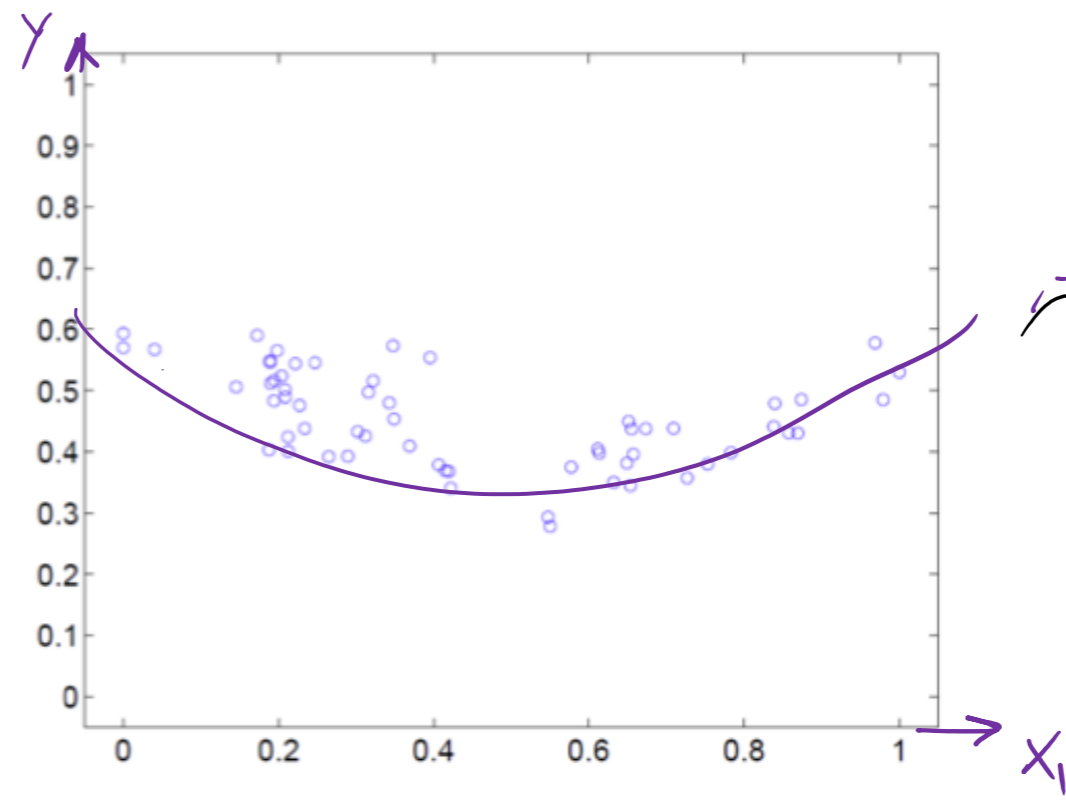
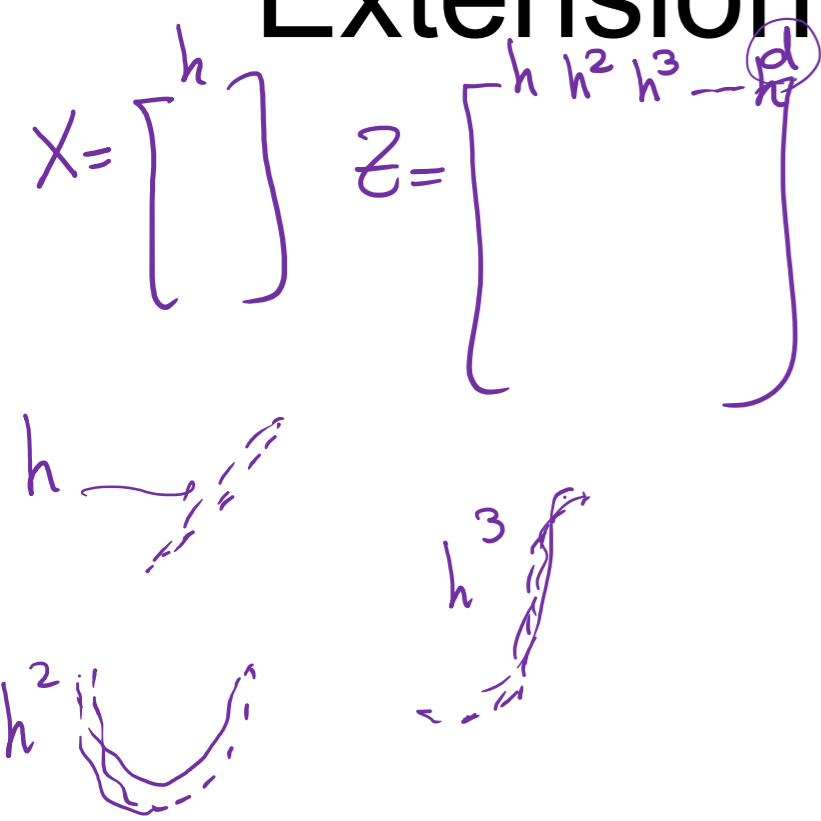
Plane



Outline

- Supervised Learning
- Linear Regression
- Extension ← 

Extension to Higher-Order Regression



- Want to fit a polynomial regression model

$$y = \theta_0 + \theta_1 x + \theta_2 x^2 + \dots + \theta_d x^d + \epsilon \quad \text{LCF}$$

$\hat{y}_p = \theta_0 + \theta_1 x$

- $z = \{1, x, x^2, \dots, x^d\} \in R^d$ 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

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

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

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

- If we choose a different maximal degree **d** for the polynomial, the solution will be different.

What is happening in polynomial regression?

$$x = [0, 0.5, 1, \dots, 9.5, 10]$$

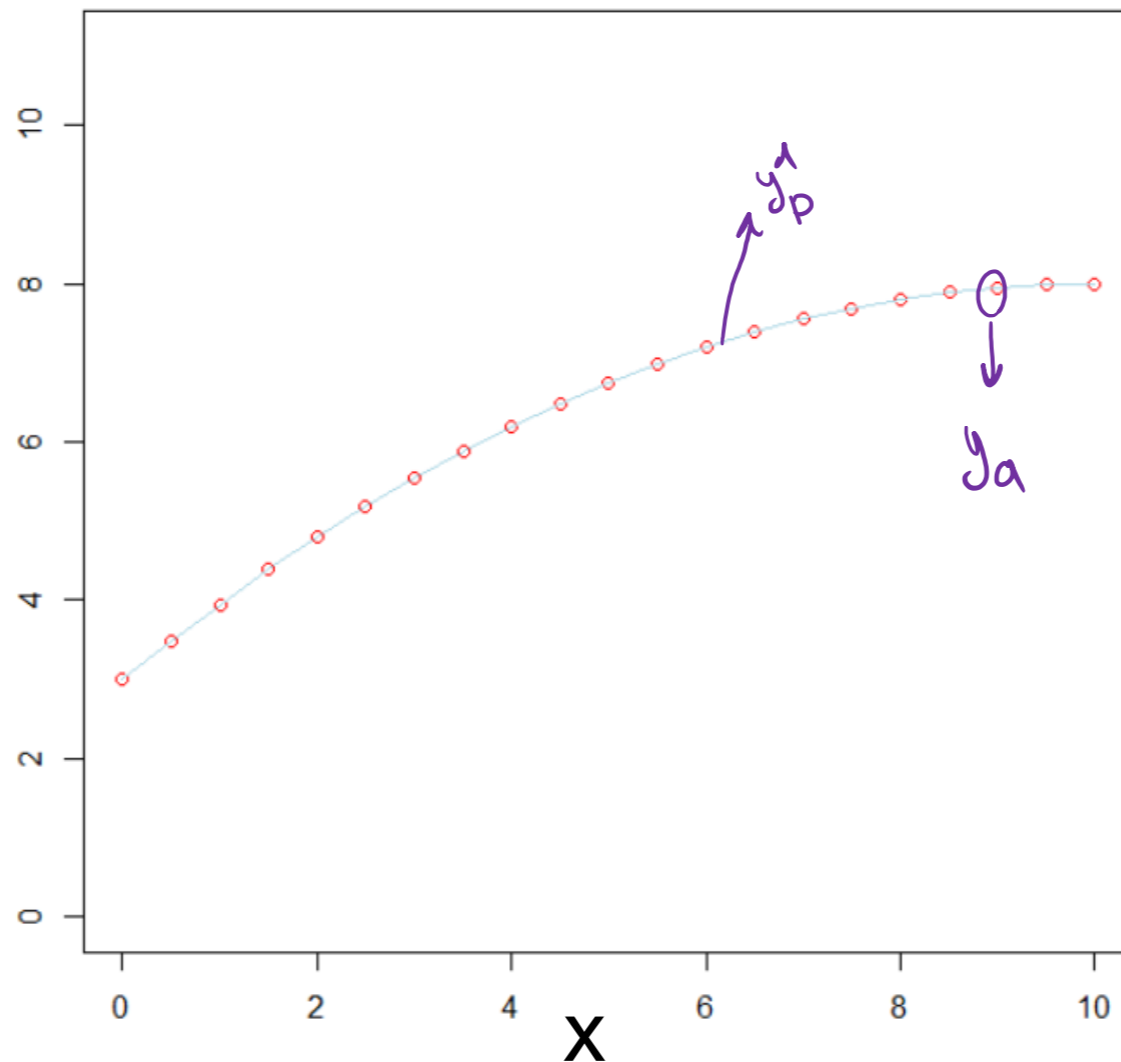
$$y = [3, 3.4875, 3.95, \dots, 7.98, 8]$$

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

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

$$Z = \begin{bmatrix} x & x^2 \end{bmatrix} \quad y$$

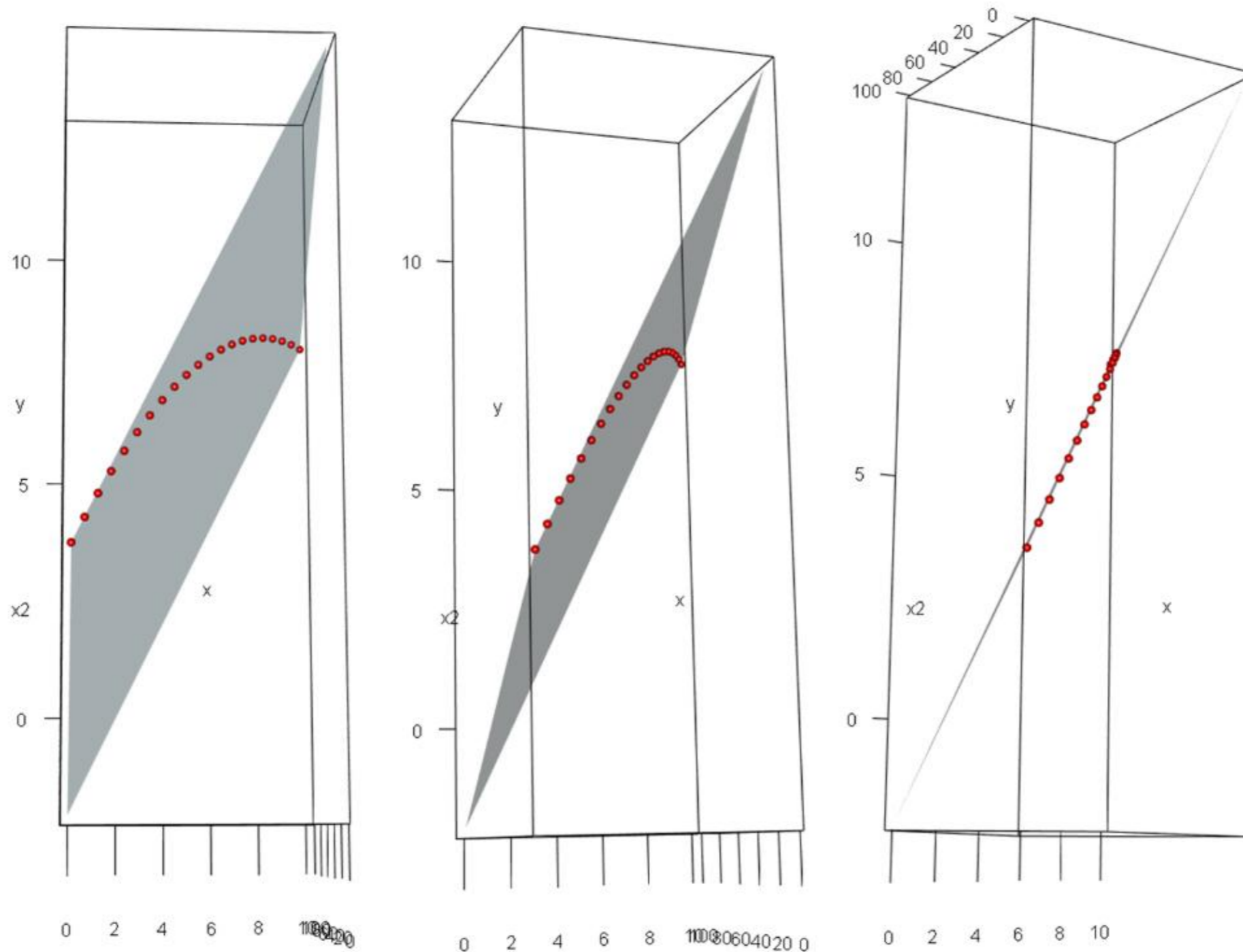
$$X = \begin{bmatrix} x \end{bmatrix}$$



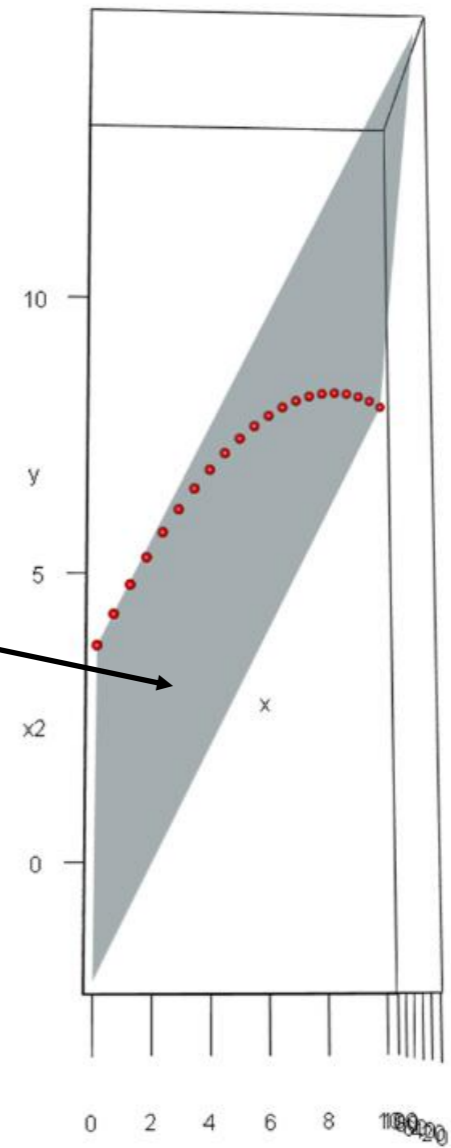
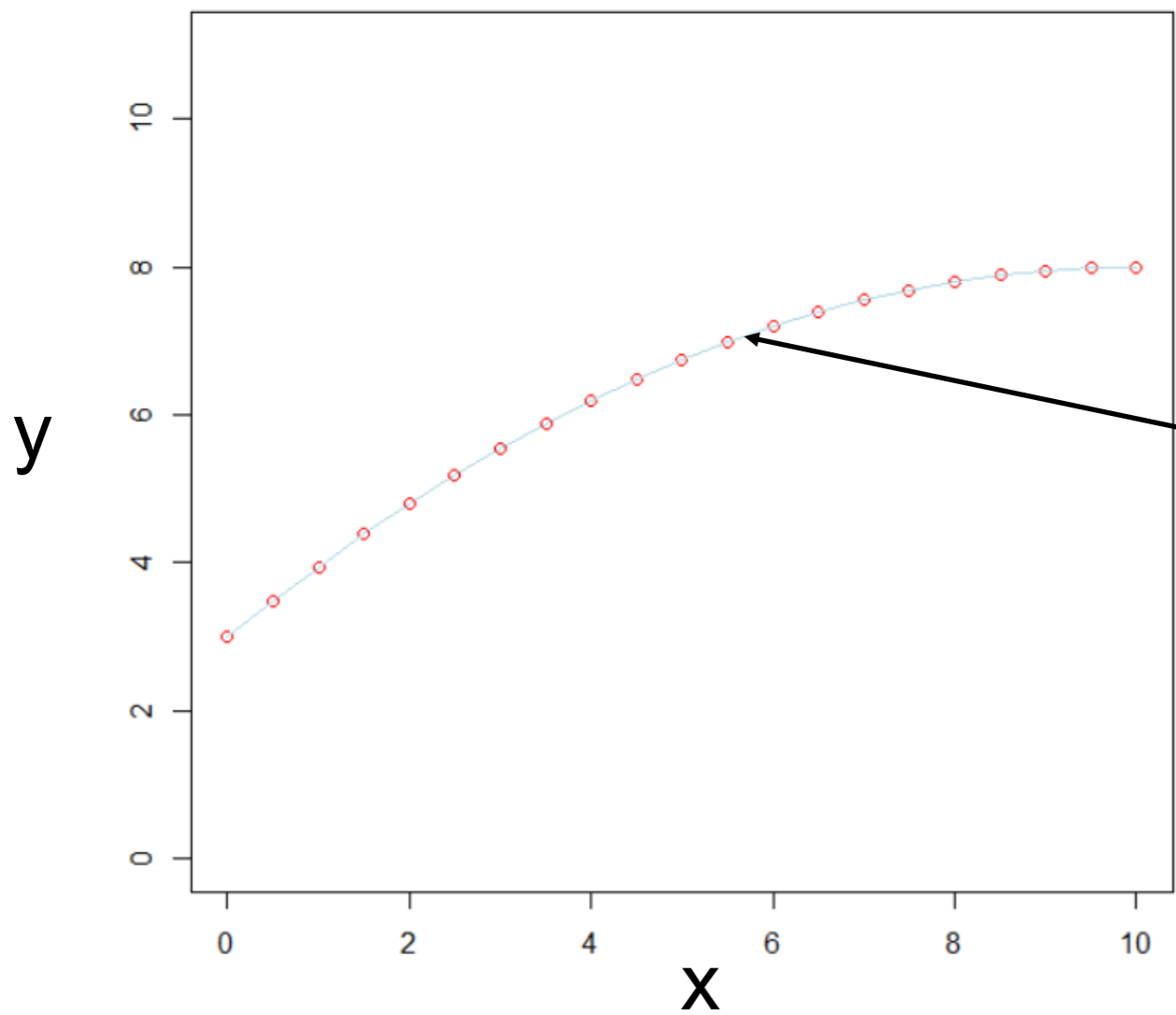
RMSE=0

Let's add to the feature space

$$x_1 = [0, 0.5, 1, \dots, 9.5, 10] \quad x_2 = [0, 0.25, 1, \dots, 90.25, 100]$$
$$y = [3, 3.4875, 3.95, \dots, 7.98, 8]$$



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.

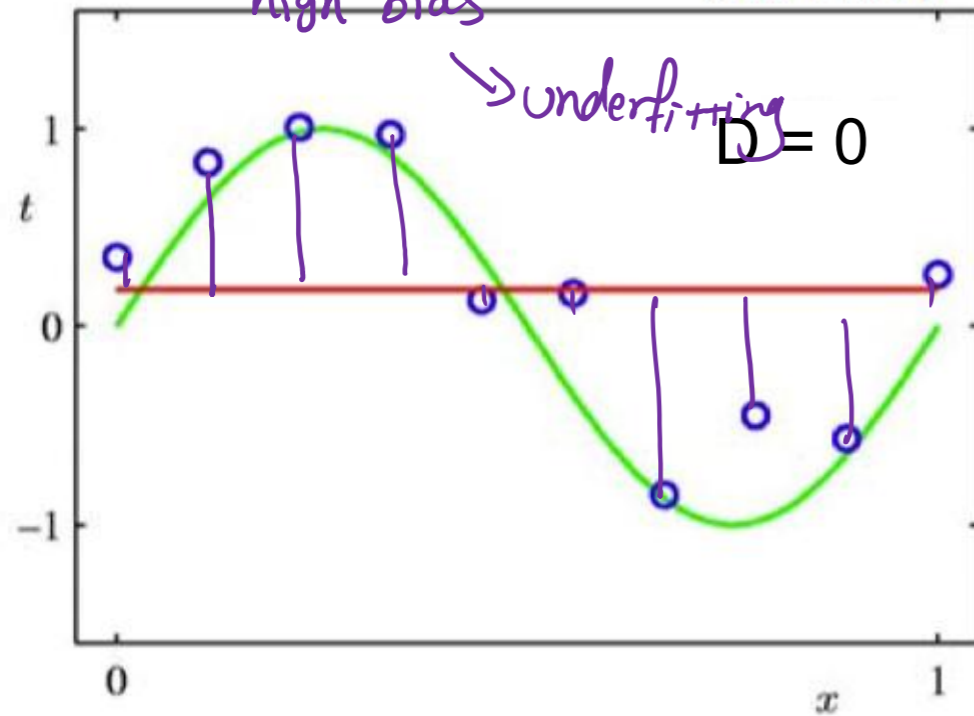


Increasing the Maximal Degree

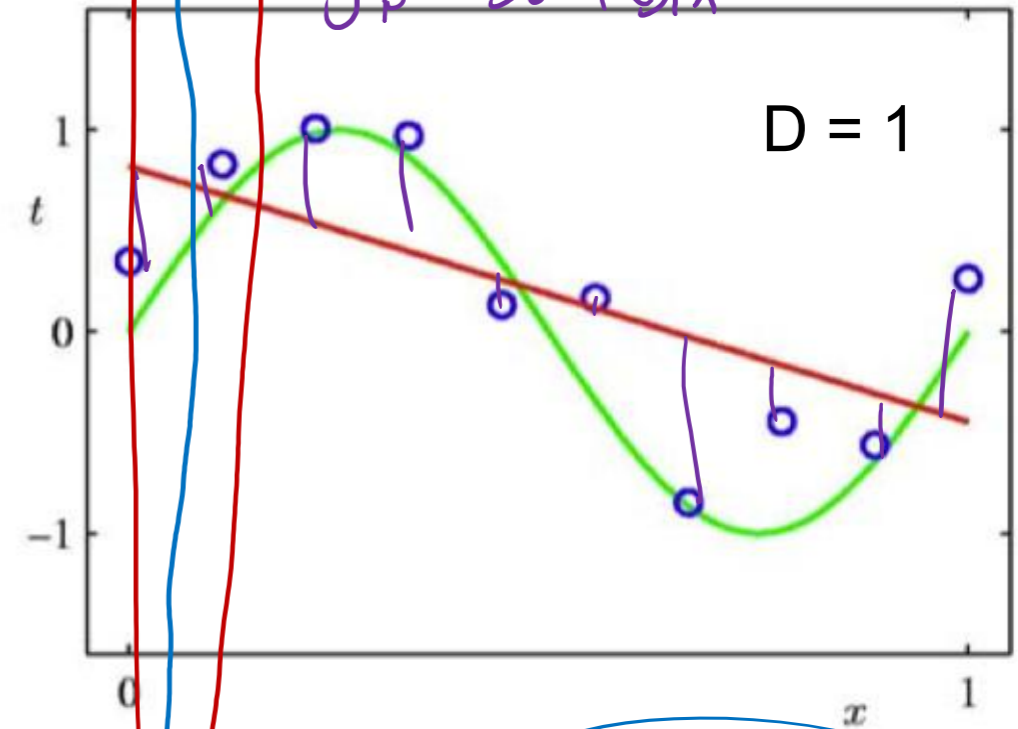
$$\hat{y}_p = \Theta_0$$

high bias

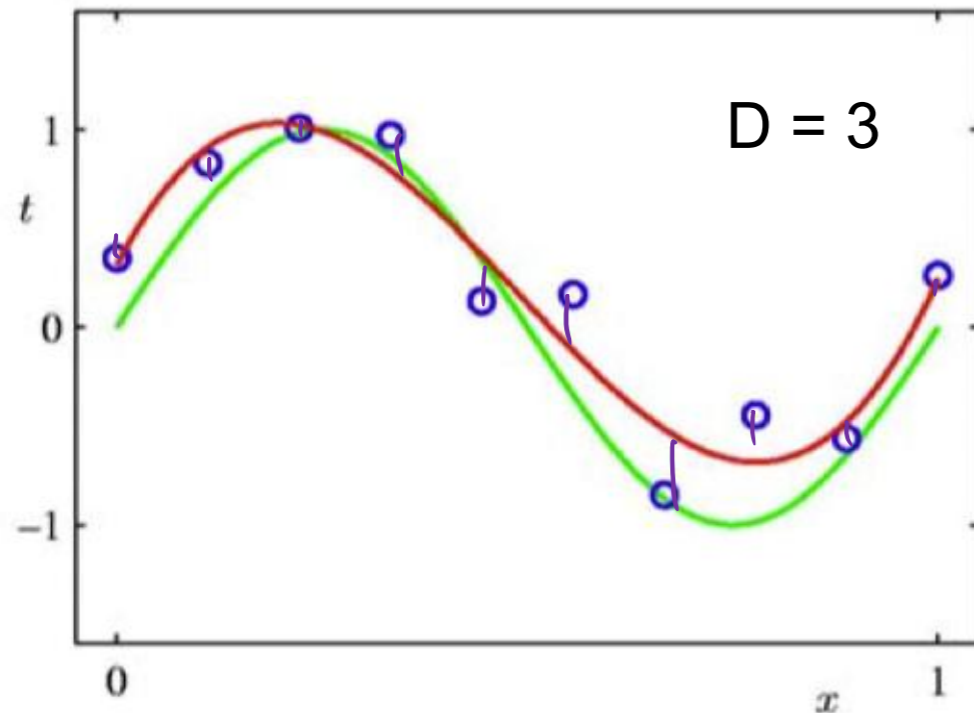
from Bishop



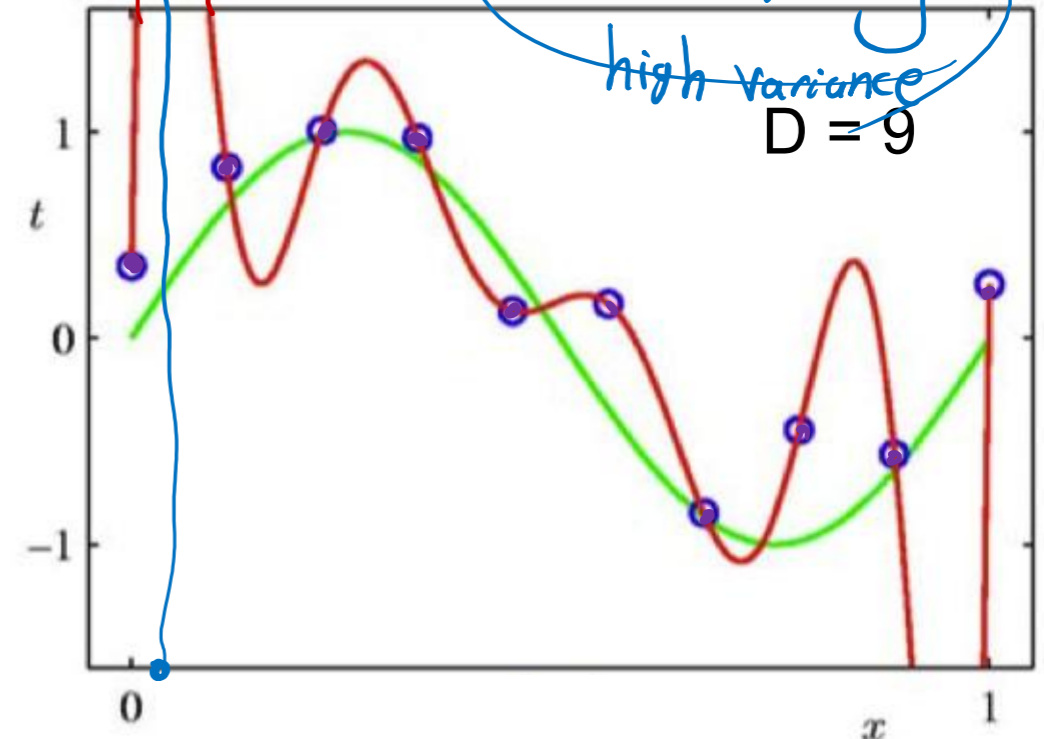
$$\hat{y}_p = \Theta_0 + \Theta_1 x$$

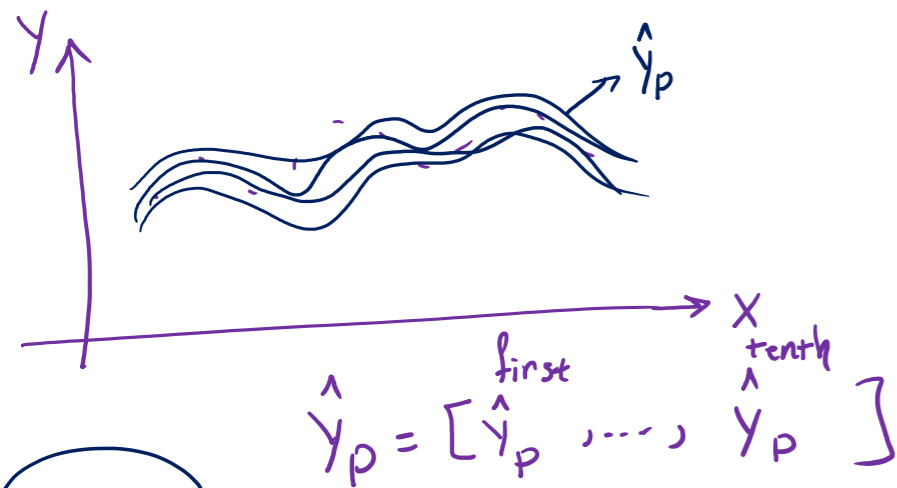
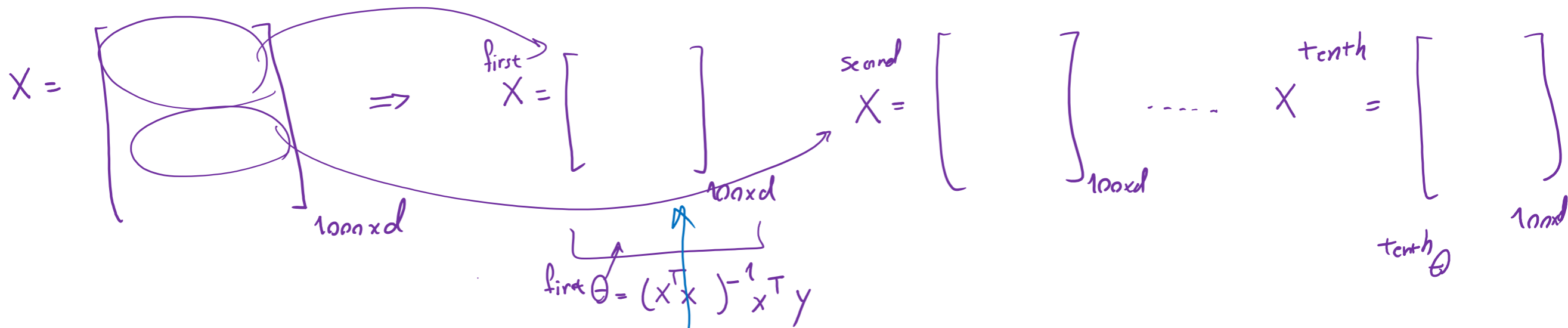


$$\hat{y}_p = \Theta_0 + \Theta_1 x + \Theta_2 x^2 + \Theta_3 x^3$$

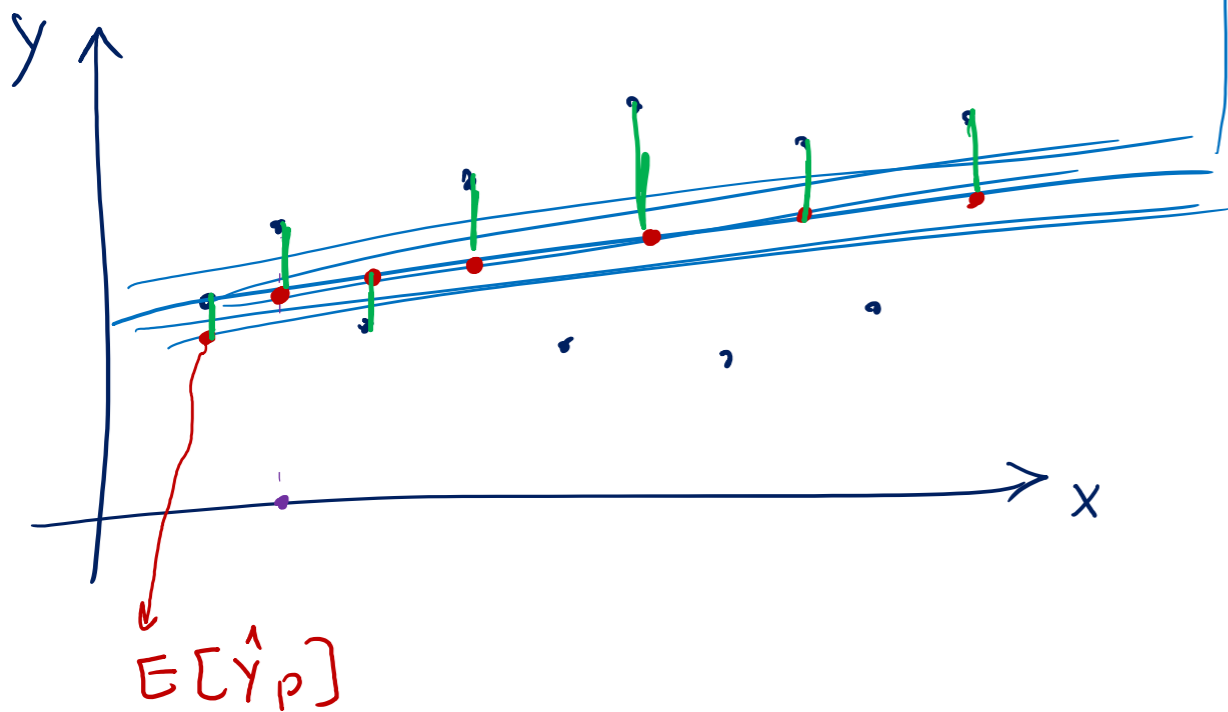


Overfitting
high variance





Bias



$$\frac{1}{N} \sum (\cdot) = E[\cdot]$$

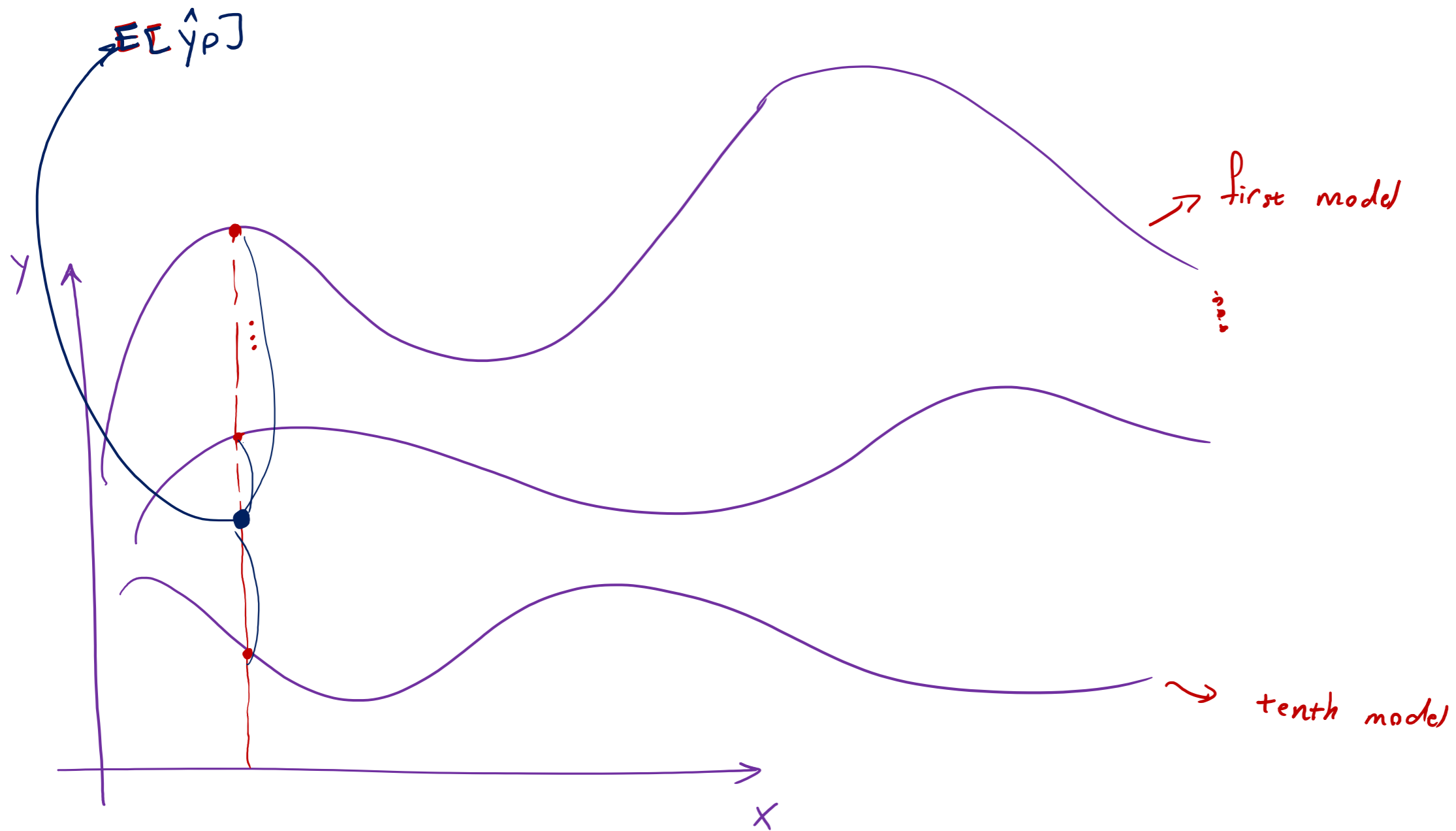
$$E[\hat{y}_p] = \frac{1}{10} \sum_{j=1}^{10} \hat{y}_p^{(j)} \in \mathbb{R}$$

$\{\text{j}^{\text{th}} \text{ dataset}\}$

$$\text{Bias} = \hat{y}_a - E[\hat{y}_p]$$

$$E[E[\hat{y}_p]] = E[\hat{y}_p]$$

$\in \mathbb{R}$



$$\text{Var}(x) = \frac{1}{N} \sum_{i=1}^N (x^{i3} - \mu)^2 = E[(x^{i3} - \mu)^2] = E[(x^{i3} - E[x])^2]$$

$$\text{Variance} = \frac{1}{10} \sum_{j=1}^N \left(\overset{\text{jth model}}{\hat{y}_p} - E[\hat{y}_p] \right)^2 = E \left[\left(\overset{\text{jth model}}{\hat{y}_p} - E[\hat{y}_p] \right)^2 \right]$$

Bias-Variance Trade off

[Animation](#)

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

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

$$(y_a - y_p)^2 = \overset{a}{(y_a - E[y_p])} + \overset{b}{E[y_p] - y_p} \Big)^2$$

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

$$\cancel{E \left[\underbrace{(y_a - E[y_p])^2}_{\text{Scalar}} \right]} + E \left[\underbrace{(E[y_p] - y_p)^2}_{\text{Scalar}} \right] + 2 E \left[\underbrace{(y_a - E[y_p])}_{\text{Scalar}} \right] \underbrace{E \left[\underbrace{(E[y_p] - y_p)}_{\text{Scalar}} \right]}_{\text{Scalar}}$$

$$(y_a - E[y_p])^2 + E \left[(E[y_p] - y_p)^2 \right] + 0 = \text{bias}^2 + \text{Variance}$$

$$\frac{\cancel{E[E[y_p]]} - E[y_p]}{E[y_p] - E[y_p]} = 0$$

Bias-Variance Trade off

[Animation](#)

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

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

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

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

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

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

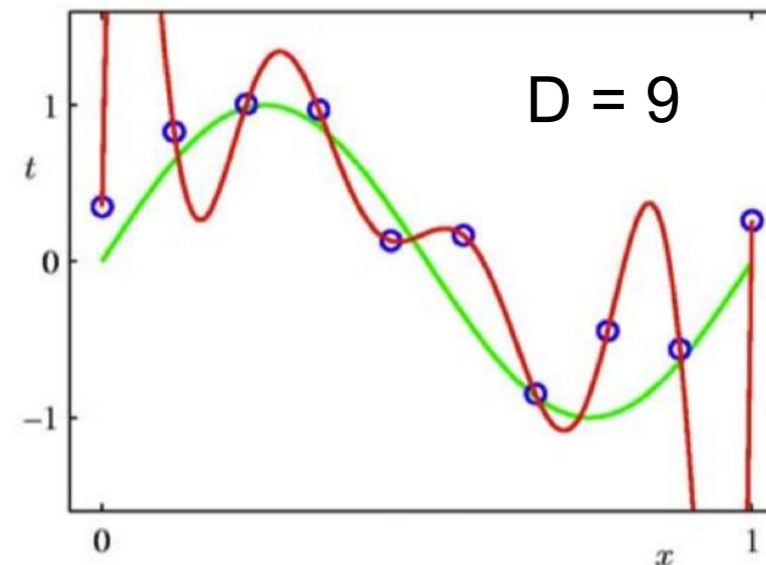
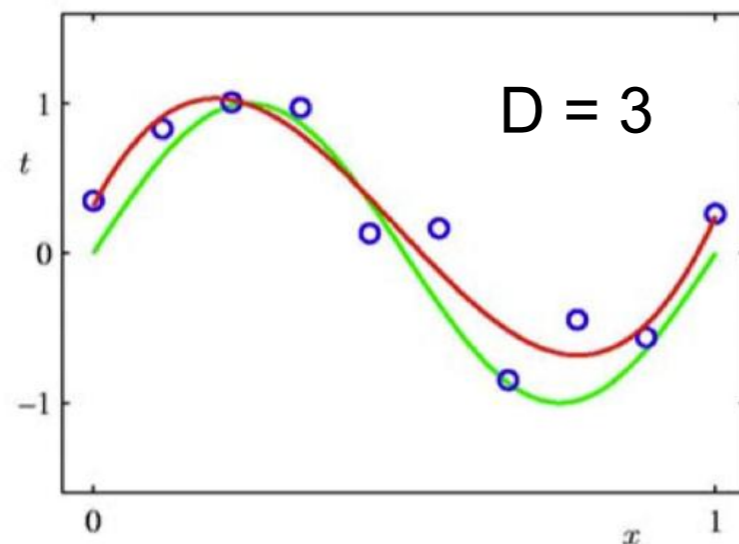
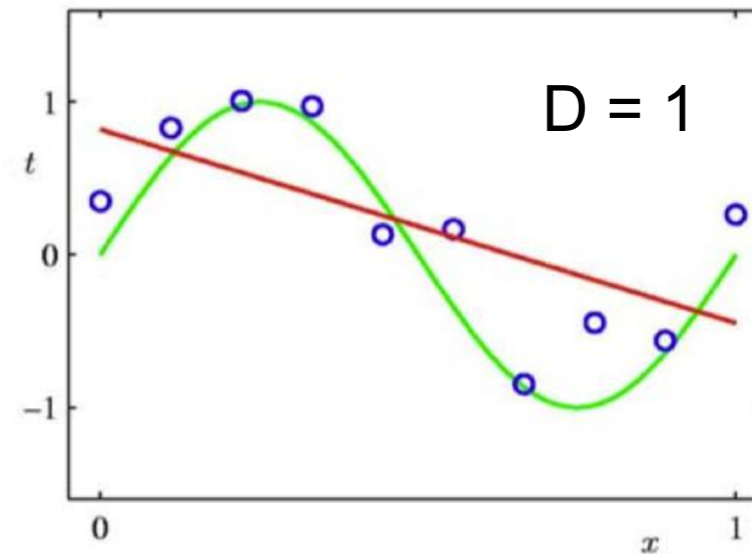
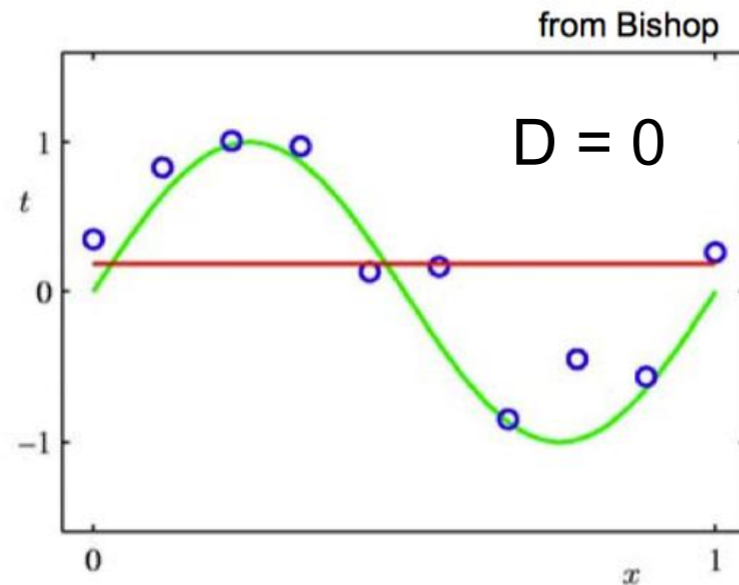
$$= [true\ value - mean(predictions)]^2 - mean[(mean(prediction) - prediction)^2]$$

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

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

$$\begin{aligned} & E[2(y_a - E[y_p])(E[y_p] - y_p)] \\ &= 2(y_a - E[y_p])E[E[y_p] - y_p] \\ &= 2(y_a - E[y_p])\left(E[E[y_p]] - E[y_p]\right) \\ &= 2(y_a - E[y_p])(E[y_p] - E[y_p]) = 0 \end{aligned}$$

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