



香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

COMP 5212  
Machine Learning  
Lecture 2

# Math Basics and Supervised Learning: Regression

Junxian He  
Feb 10, 2026

# Announcement

Lecture videos till next week will be released considering the Lunar New Year

# Example Distributions

Distribution	PDF or PMF	Mean	Variance
<i>Bernoulli</i> ( $p$ )	$\begin{cases} p, & \text{if } x = 1 \\ 1 - p, & \text{if } x = 0. \end{cases}$	$p$	$p(1 - p)$
<i>Binomial</i> ( $n, p$ )	$\binom{n}{k} p^k (1 - p)^{n-k}$ for $k = 0, 1, \dots, n$	$np$	$np(1 - p)$
<i>Geometric</i> ( $p$ )	$p(1 - p)^{k-1}$ for $k = 1, 2, \dots$	$\frac{1}{p}$	$\frac{1-p}{p^2}$
<i>Poisson</i> ( $\lambda$ )	$\frac{e^{-\lambda} \lambda^k}{k!}$ for $k = 0, 1, \dots$	$\lambda$	$\lambda$
<i>Uniform</i> ( $a, b$ )	$\frac{1}{b-a}$ for all $x \in (a, b)$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
<i>Gaussian</i> ( $\mu, \sigma^2$ )	$\frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$ for all $x \in (-\infty, \infty)$	$\mu$	$\sigma^2$
<i>Exponential</i> ( $\lambda$ )	$\lambda e^{-\lambda x}$ for all $x \geq 0, \lambda \geq 0$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$

# Joint and Marginal Distributions

# Joint and Marginal Distributions

- Joint PMF for discrete RV's  $X, Y$ :

$$p_{XY}(x, y) = P(X = x, Y = y)$$

Note that  $\sum_{x \in \text{Val}(X)} \sum_{y \in \text{Val}(Y)} p_{XY}(x, y) = 1$

# Joint and Marginal Distributions

- Joint PMF for discrete RV's  $X, Y$ :

$$p_{XY}(x, y) = P(X = x, Y = y)$$

Note that  $\sum_{x \in \text{Val}(X)} \sum_{y \in \text{Val}(Y)} p_{XY}(x, y) = 1$

- Marginal PMF of  $X$ , given joint PMF of  $X, Y$ :

$$p_X(x) = \sum_y p_{XY}(x, y)$$

$p(x, y)$        $p(x)$        $p(y)$

*marginalizing  
out  $Y$*

# Joint and Marginal Distributions

- **Joint PDF** for continuous RV's  $X_1, \dots, X_n$ :

$$f(x_1, \dots, x_n) = \frac{\delta^n F(x_1, \dots, x_n)}{\delta x_1 \delta x_2 \dots \delta x_n}$$

Note that  $\int_{x_1} \int_{x_2} \dots \int_{x_n} f(x_1, \dots, x_n) dx_1 \dots dx_n = 1$

- **Marginal PDF** of  $X_1$ , given joint PDF of  $X_1, \dots, X_n$ :

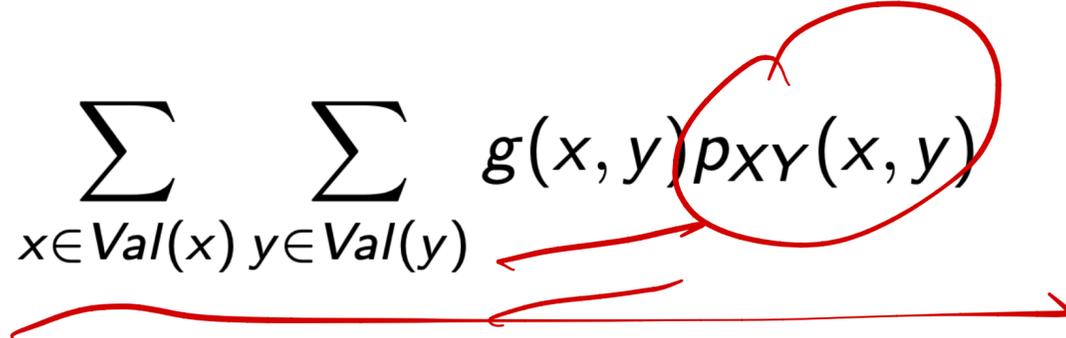
$$f_{X_1}(x_1) = \int_{x_2} \dots \int_{x_n} f(x_1, \dots, x_n) dx_2 \dots dx_n$$

# Expectation for multiple random variables

# Expectation for multiple random variables

Given two RV's  $X, Y$  and a function  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$  of  $X, Y$ ,

- for discrete  $X, Y$ :

$$\mathbb{E}[g(X, Y)] := \sum_{x \in \text{Val}(x)} \sum_{y \in \text{Val}(y)} g(x, y) p_{XY}(x, y)$$


# Expectation for multiple random variables

Given two RV's  $X, Y$  and a function  $g : \mathbb{R}^2 \rightarrow \mathbb{R}$  of  $X, Y$ ,

- for discrete  $X, Y$ :

$$\mathbb{E}[g(X, Y)] := \sum_{x \in \text{Val}(x)} \sum_{y \in \text{Val}(y)} g(x, y) p_{XY}(x, y)$$

- for continuous  $X, Y$ :

$$\mathbb{E}[g(X, Y)] := \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} g(x, y) f_{XY}(x, y) dx dy$$

# Covariance

$$\mathbb{E}[(X - \mathbb{E}[X])^2]$$

**Intuitively:** measures how much one RV's value tends to move with another RV's value. For RV's  $X, Y$ :

$$\begin{aligned}\text{Cov}[X, Y] &:= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$

# Covariance

**Intuitively:** measures how much one RV's value tends to move with another RV's value. For RV's  $X, Y$ :

$$\begin{aligned}\text{Cov}[X, Y] &:= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]\end{aligned}$$



- If  $\text{Cov}[X, Y] < 0$ , then  $X$  and  $Y$  are negatively correlated
- If  $\text{Cov}[X, Y] > 0$ , then  $X$  and  $Y$  are positively correlated
- If  $\text{Cov}[X, Y] = 0$ , then  $X$  and  $Y$  are uncorrelated

# Variance of two variables

$$\text{Var}[X + Y] = \text{Var}[X] + \text{Var}[Y] + 2\text{Cov}[X, Y]$$

# Conditional distributions for RVs

Works the same way with *RV*'s as with events:

- For discrete  $X, Y$ :

$$p_{Y|X}(y|x) = \frac{p_{XY}(x, y)}{p_X(x)}$$

*joint*  
*marginal*

- For continuous  $X, Y$ :

$$f_{Y|X}(y|x) = \frac{f_{XY}(x, y)}{f_X(x)}$$

- In general, for continuous  $X_1, \dots, X_n$ :

$$f_{X_1|X_2, \dots, X_n}(x_1|x_2, \dots, x_n) = \frac{f_{X_1, X_2, \dots, X_n}(x_1, x_2, \dots, x_n)}{f_{X_2, \dots, X_n}(x_2, \dots, x_n)}$$

# Bayes' Rule for RVs

Also works the same way for *RV's* as with events:

- For discrete  $X, Y$ :

$$p_{Y|X}(y|x) = \frac{p_{X|Y}(x|y)p_Y(y)}{\sum_{y' \in \text{Val}(Y)} p_{X|Y}(x|y')p_Y(y')}$$

- For continuous  $X, Y$ :

$$f_{Y|X}(y|x) = \frac{f_{X|Y}(x|y)f_Y(y)}{\int_{-\infty}^{\infty} f_{X|Y}(x|y')f_Y(y')dy'}$$

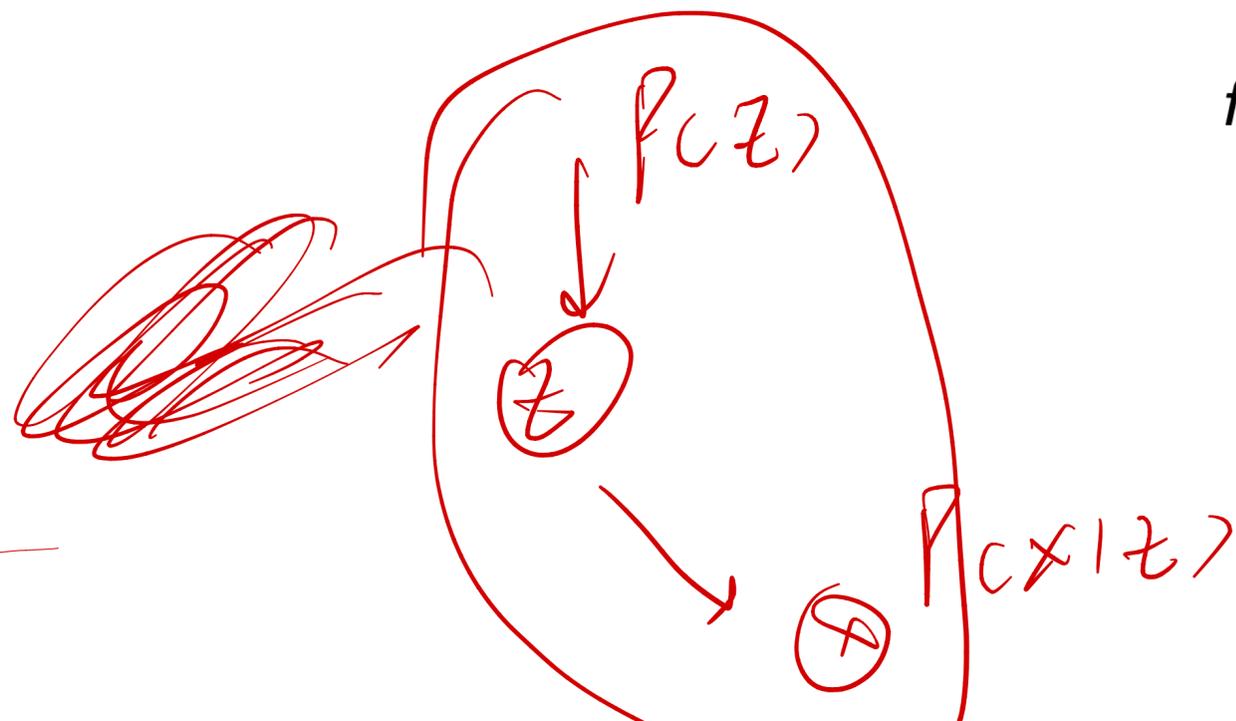
$$\frac{p_Y(y) p_{X|Y}(x|y)}{p_X(x)} = \frac{p_{X|Y}(x|y) p_Y(y)}{p_X(x)}$$

$$p_Y(y|x)$$

$$p_X(x) = \sum_y p_{X,Y}(x,y)$$

$$= \sum_y p_Y(y) p_{X|Y}(x|y)$$

$$p_{Y|X}(y|x)$$



# Random Vectors

# Random Vectors

Given  $n$  RV's  $X_1, \dots, X_n$ , we can define a random vector  $X$  s.t.

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$$

# Random Vectors

Given  $n$  RV's  $X_1, \dots, X_n$ , we can define a random vector  $X$  s.t.

$$X = \begin{bmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{bmatrix}$$

Note: all the notions of joint PDF/CDF will apply to  $X$ .

Given  $g : \mathbb{R}^n \rightarrow \mathbb{R}^m$ , we have:

$$g(x) = \begin{bmatrix} g_1(x) \\ g_2(x) \\ \vdots \\ g_m(x) \end{bmatrix}, \mathbb{E}[g(X)] = \begin{bmatrix} \mathbb{E}[g_1(X)] \\ \mathbb{E}[g_2(X)] \\ \vdots \\ \mathbb{E}[g_m(X)] \end{bmatrix}$$

# Covariance Matrices

# Covariance Matrices

For a random vector  $X \in \mathbb{R}^n$ , we define its **covariance matrix**  $\Sigma$  as the  $n \times n$  matrix whose  $ij$ -th entry contains the covariance between  $X_i$  and  $X_j$ .

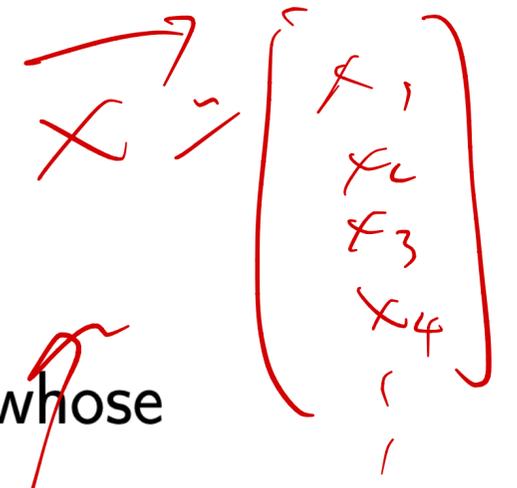
# Covariance Matrices

For a random vector  $X \in \mathbb{R}^n$ , we define its **covariance matrix**  $\Sigma$  as the  $n \times n$  matrix whose  $ij$ -th entry contains the covariance between  $X_i$  and  $X_j$ .

$$\Sigma = \begin{bmatrix} \text{Cov}[X_1, X_1] & \dots & \text{Cov}[X_1, X_n] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_n, X_1] & \dots & \text{Cov}[X_n, X_n] \end{bmatrix} \rightarrow n \times n$$

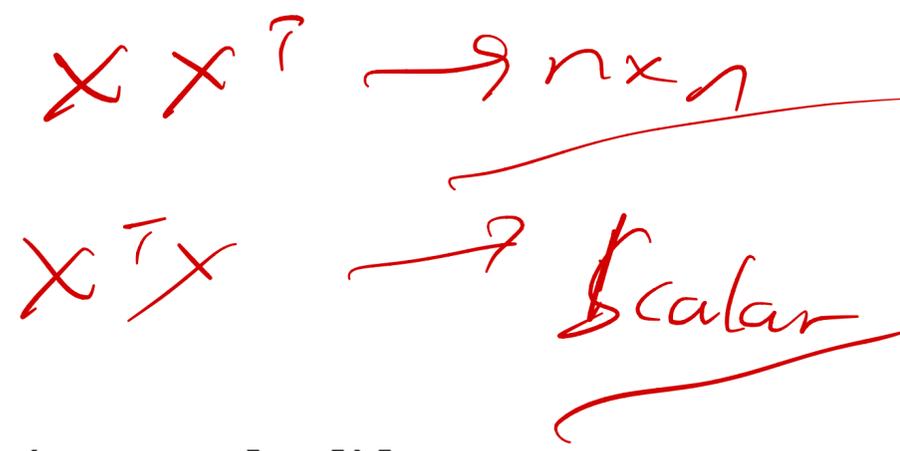
# Covariance Matrices

For a random vector  $X \in \mathbb{R}^n$ , we define its **covariance matrix**  $\Sigma$  as the  $n \times n$  matrix whose  $ij$ -th entry contains the covariance between  $X_i$  and  $X_j$ .



*high-dim  
Gaussian*

$$\Sigma = \begin{bmatrix} \text{Cov}[X_1, X_1] & \dots & \text{Cov}[X_1, X_n] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_n, X_1] & \dots & \text{Cov}[X_n, X_n] \end{bmatrix}$$



applying linearity of expectation and the fact that  $\text{Cov}[X_i, X_j] = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]$ , we obtain

$$\Sigma = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T]$$

*outer product*

$X - \mathbb{E}[X]$  (with  $n$  below it)

$X - \mathbb{E}[X]$  (with  $n$  below it)

# Covariance Matrices

For a random vector  $X \in \mathbb{R}^n$ , we define its **covariance matrix**  $\Sigma$  as the  $n \times n$  matrix whose  $ij$ -th entry contains the covariance between  $X_i$  and  $X_j$ .

$$\Sigma = \begin{bmatrix} \text{Cov}[X_1, X_1] & \dots & \text{Cov}[X_1, X_n] \\ \vdots & \ddots & \vdots \\ \text{Cov}[X_n, X_1] & \dots & \text{Cov}[X_n, X_n] \end{bmatrix}$$

applying linearity of expectation and the fact that  $\text{Cov}[X_i, X_j] = \mathbb{E}[(X_i - \mathbb{E}[X_i])(X_j - \mathbb{E}[X_j])]$ , we obtain

$$\Sigma = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^T]$$

## Properties:

- $\Sigma$  is symmetric and PSD
- If  $X_i \perp X_j$  for all  $i, j$ , then  $\Sigma = \text{diag}(\text{Var}[X_1], \dots, \text{Var}[X_n])$

$z^T A z \geq 0 \xrightarrow{A \text{ is}} \text{PSD}$

# Multivariate Gaussian

The multivariate Gaussian  $X \sim \mathcal{N}(\mu, \Sigma)$ ,  $X \in \mathbb{R}^n$ :

$$p(x; \mu, \Sigma) = \frac{1}{\det(\Sigma)^{\frac{1}{2}} (2\pi)^{\frac{n}{2}}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

variance =  $\sigma^2$

$$\frac{1}{\sqrt{2\pi} \sigma} \exp\left[-\frac{1}{2} \frac{(x - \mu)^2}{\sigma^2}\right]$$

# Multivariate Gaussian

The multivariate Gaussian  $X \sim \mathcal{N}(\mu, \Sigma)$ ,  $X \in \mathbb{R}^n$ :

$$p(x; \mu, \Sigma) = \frac{1}{\det(\Sigma)^{\frac{1}{2}} (2\pi)^{\frac{n}{2}}} \exp\left(-\frac{1}{2}(x - \mu)^T \Sigma^{-1}(x - \mu)\right)$$

Gaussian when  $n = 1$ .

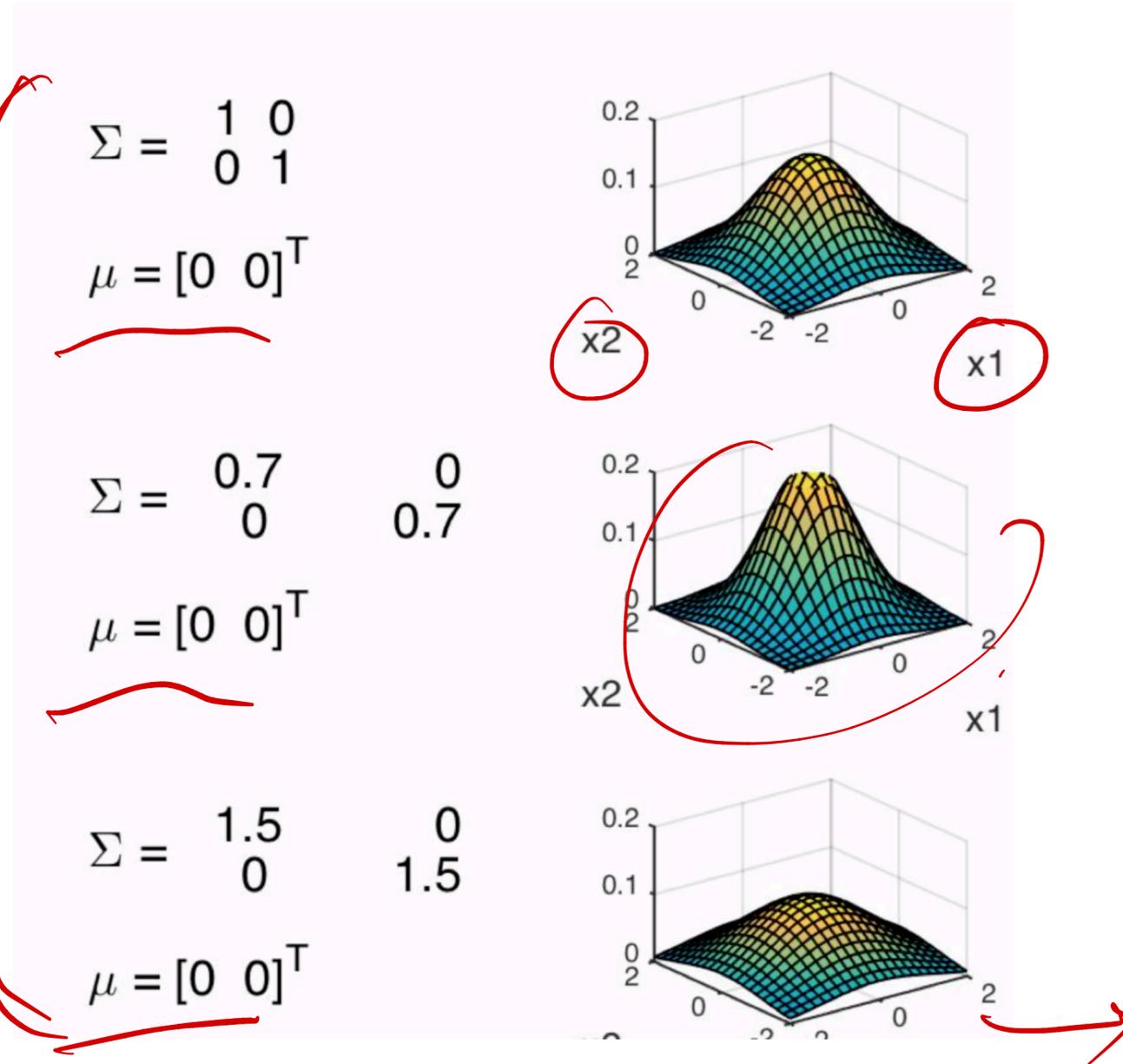
$$p(x; \mu, \sigma^2) = \frac{1}{\sigma(2\pi)^{\frac{1}{2}}} \exp\left(-\frac{1}{2\sigma^2}(x - \mu)^2\right)$$

Notice that if  $\Sigma \in \mathbb{R}^{1 \times 1}$ , then  $\Sigma = \text{Var}[X_1] = \sigma^2$ , and so  $\Sigma^{-1} = \frac{1}{\sigma^2}$  and  $\det(\Sigma)^{\frac{1}{2}} = \sigma$

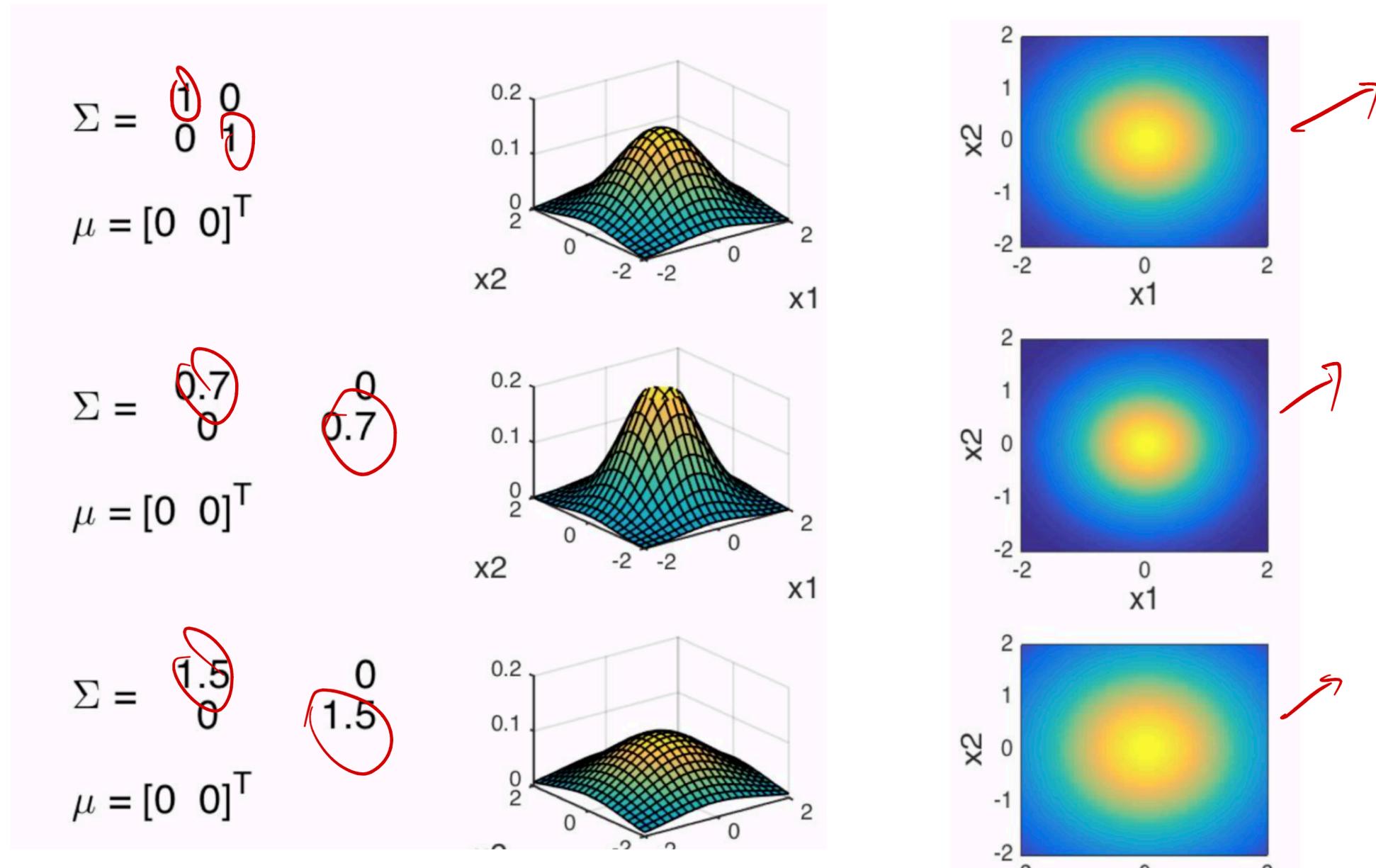
# MV Gaussian Visualization

# MV Gaussian Visualization

Effect of changing variance



# MV Gaussian Visualization

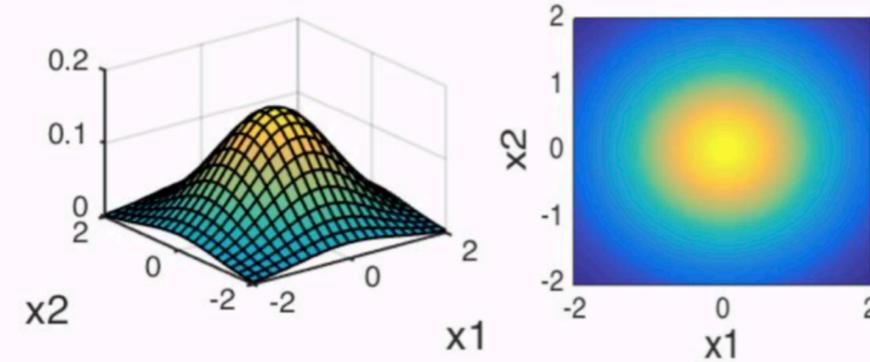


Effect of changing variance

# MV Gaussian Visualization

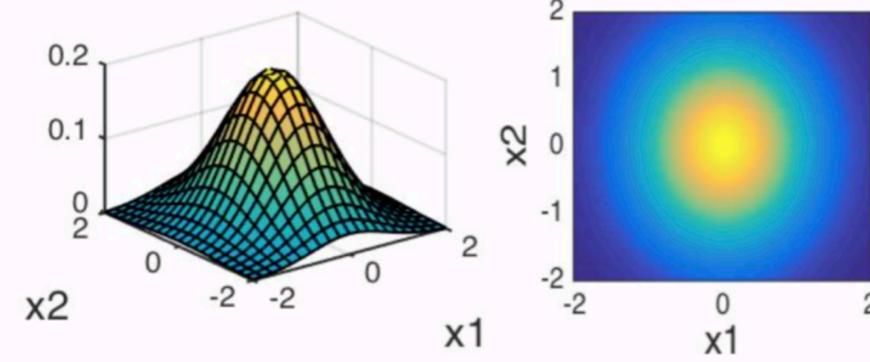
$$\Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\mu = [0 \ 0]^T$$



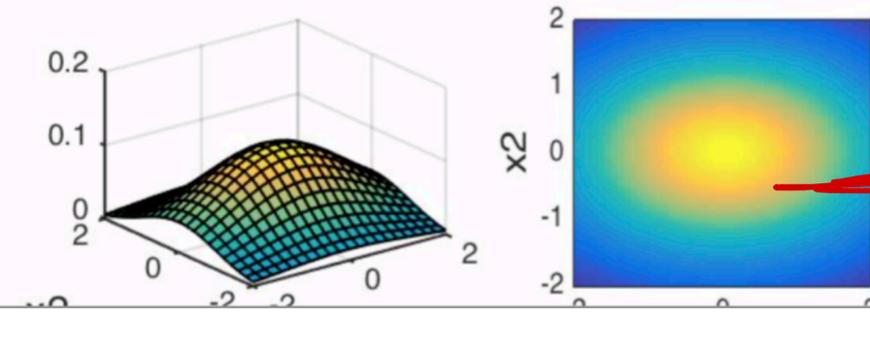
$$\Sigma = \begin{pmatrix} 0.6 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\mu = [0 \ 0]^T$$



$$\Sigma = \begin{pmatrix} 2 & 0 \\ 0 & 1 \end{pmatrix}$$

$$\mu = [0 \ 0]^T$$



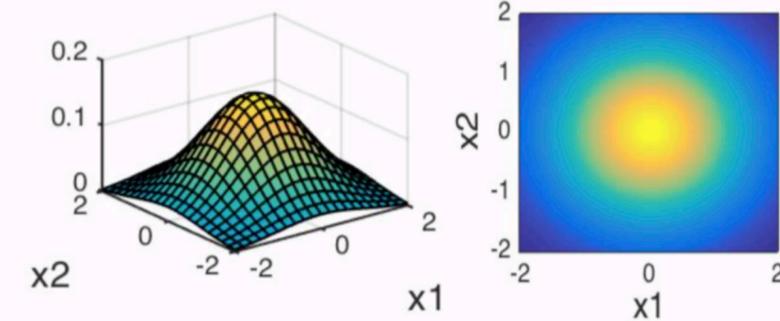
If  $Var[X_1] \neq Var[X_2]$ :

# MV Gaussian Visualization

$x_1$  ↑  $x_2$  ↓

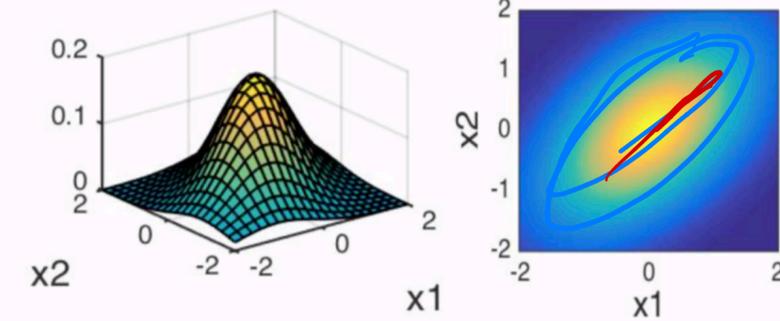
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

$$\mu = [0 \ 0]^T$$



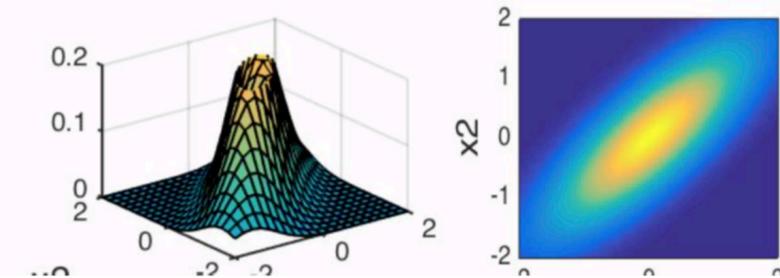
$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$

$$\mu = [0 \ 0]^T$$



$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

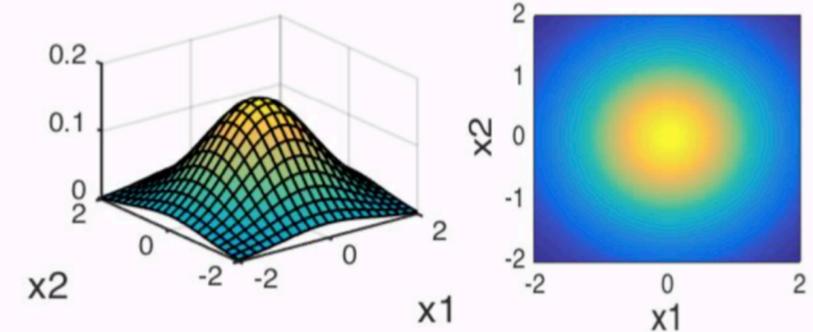
$$\mu = [0 \ 0]^T$$



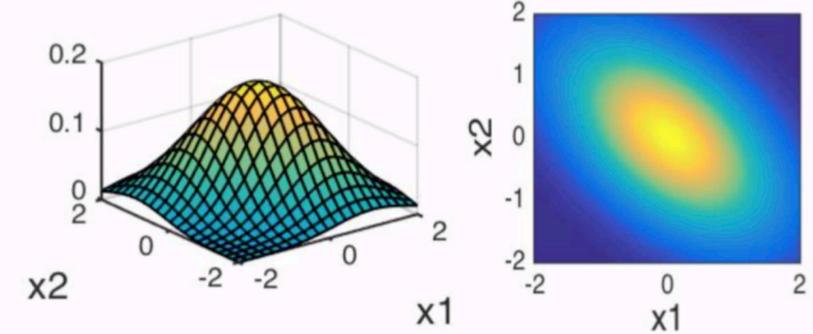
If  $X_1$  and  $X_2$  are positively correlated:

# MV Gaussian Visualization

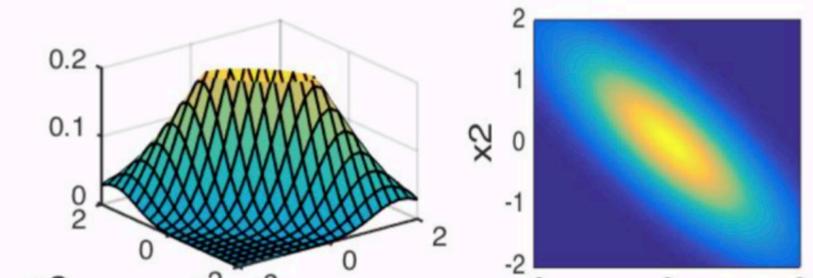
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$
$$\mu = [0 \ 0]^T$$



$$\Sigma = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$
$$\mu = [0 \ 0]^T$$



$$\Sigma = \begin{bmatrix} 1 & -0.8 \\ -0.8 & 1 \end{bmatrix}$$
$$\mu = [0 \ 0]^T$$



If  $X_1$  and  $X_2$  are negatively correlated:

**The purpose of computation is  
insight, not numbers.**

**- Richard Hamming**

# The purpose of computation is insight, not numbers.

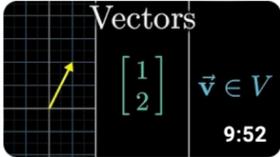
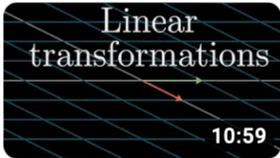
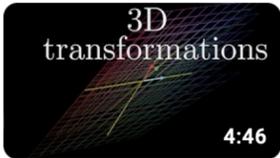
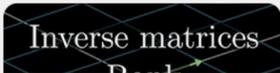
- Richard Hamming

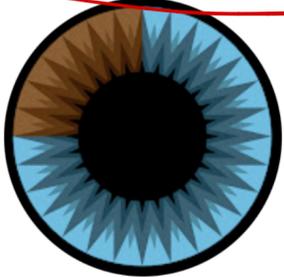
*matrix*

*voluntary embedding*

*ROPE*

<https://www.youtube.com/@3blue1brown/courses>

- 1  Vectors | Chapter 1, Essence of linear algebra 9:52
- 2  Span | Linear combinations, span, and basis vectors | Chapter 2, Essence of linear algebra 9:59
- 3  Linear transformations | Linear transformations and matrices | Chapter 3, Essence of linear algebra 10:59
- 4  Matrix multiplication | Matrix multiplication as composition | Chapter 4, Essence of linear algebra 10:04
- 5  3D transformations | Three-dimensional linear transformations | Chapter 5, Essence of linear algebra 4:46
- 6  Determinant | The determinant | Chapter 6, Essence of linear algebra 10:03
-  Inverse matrices | Inverse matrices, column space and null space | Chapter 7, Essence of linear algebra



**3Blue1Brown**

@3blue1brown · 5.88M subscribers · 172 videos

My name is Grant Sanderson. Videos here cover a variety of topics in math, or adjacent fields...

[3blue1brown.com](https://3blue1brown.com) and 7 more links

Subscribed



香港科技大學  
THE HONG KONG  
UNIVERSITY OF SCIENCE  
AND TECHNOLOGY

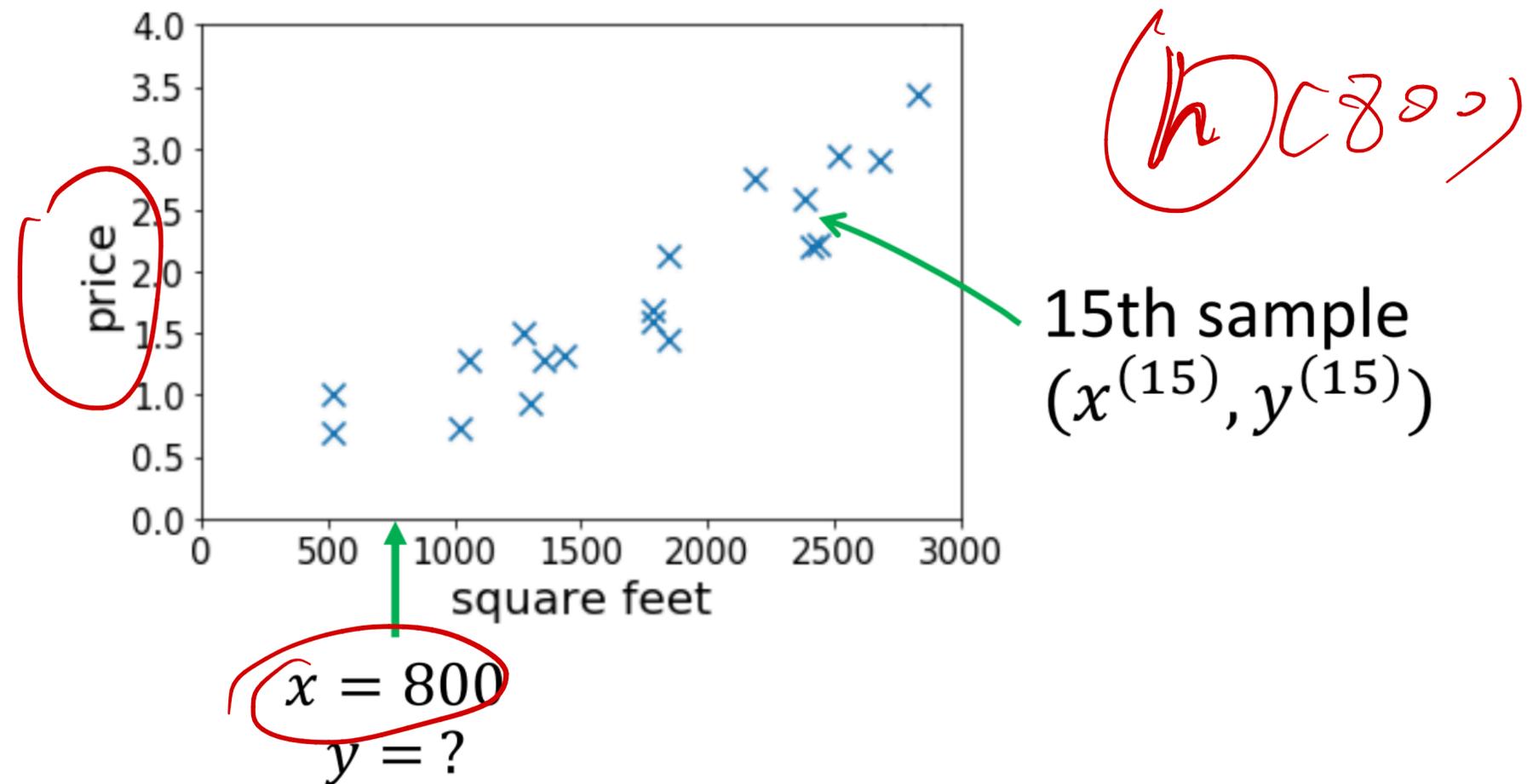
# Supervised Learning: Regression

# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$

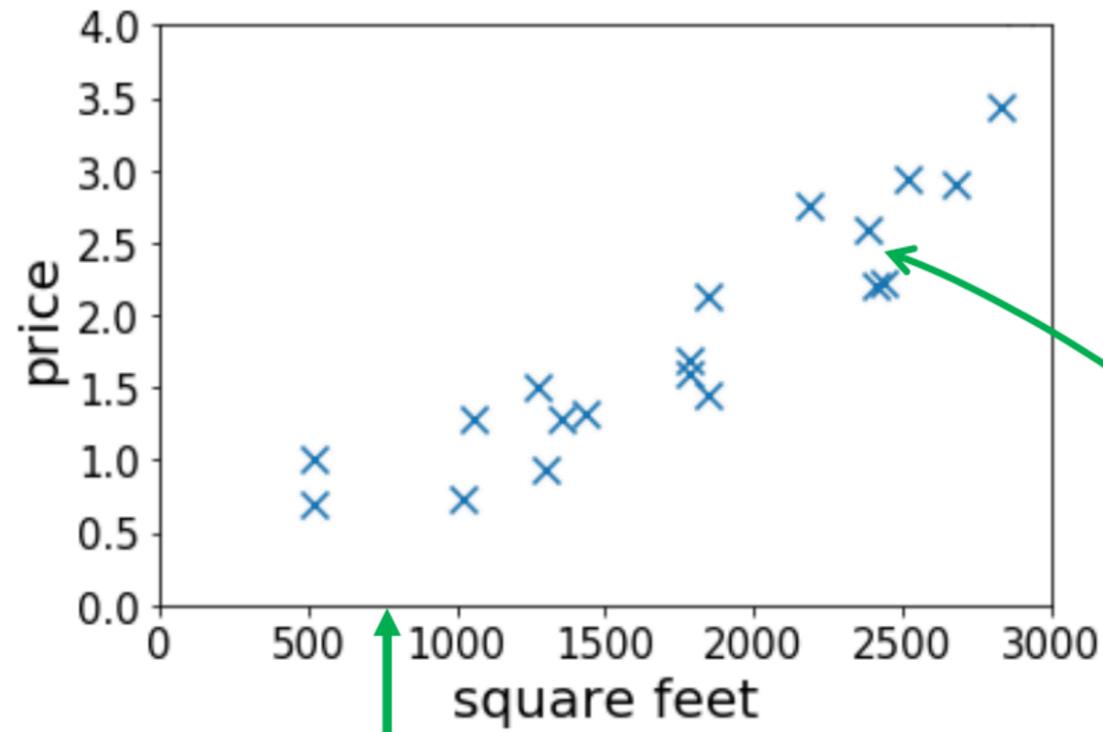
# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$



# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$



$x = 800$   
 $y = ?$

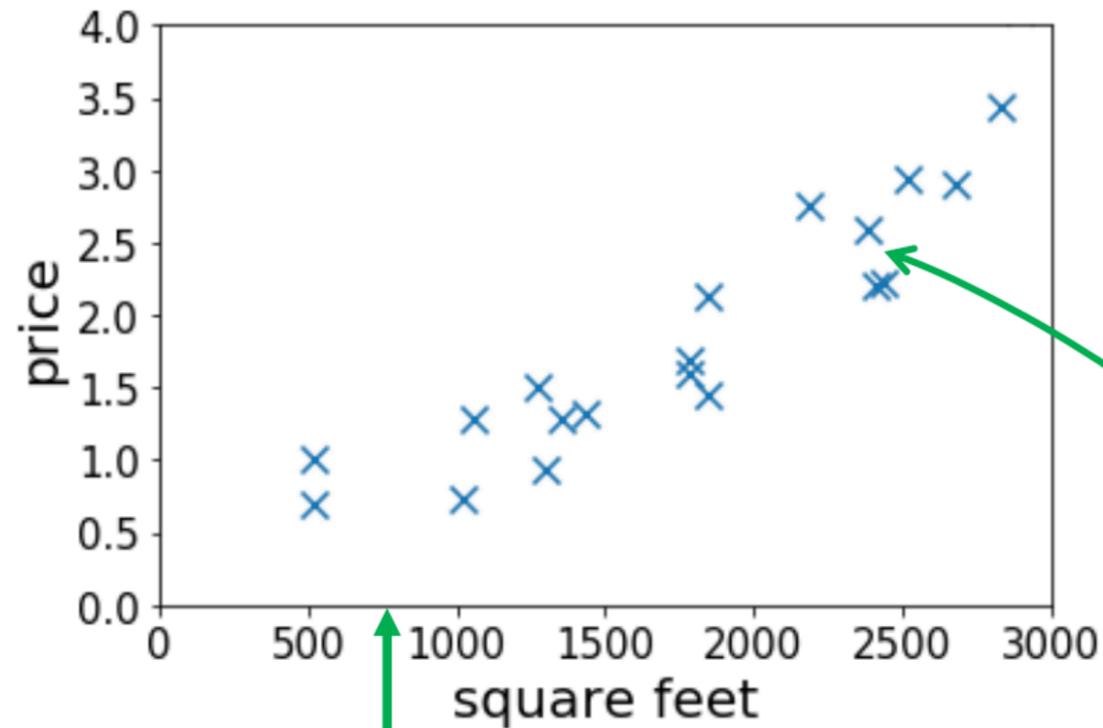
15th sample  
 $(x^{(15)}, y^{(15)})$



CAT

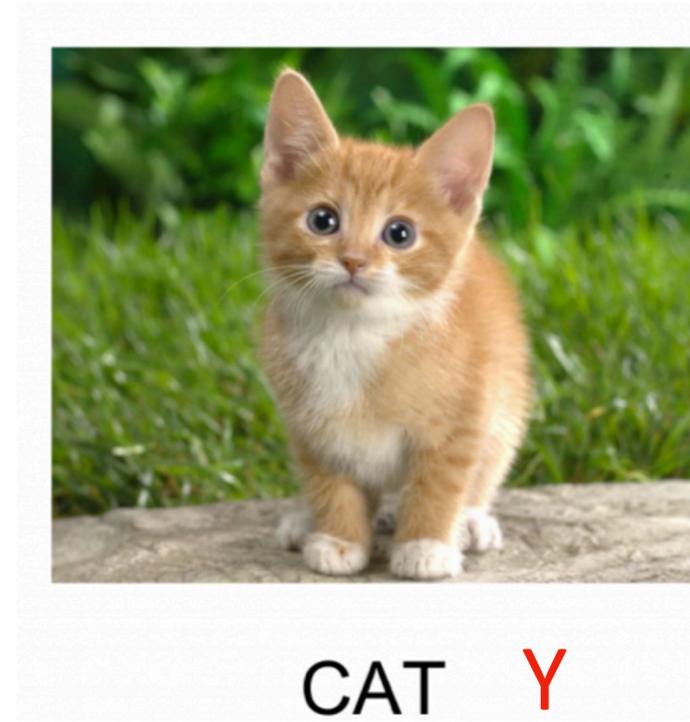
# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$



$x = 800$   
 $y = ?$

15th sample  
 $(x^{(15)}, y^{(15)})$



X

# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$

# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$
- A training set is set of pairs  $\{(\underline{x^{(1)}}), \underline{y^{(1)}}), \dots, (x^{(n)}, y^{(n)})\}$   
s.t.  $x^{(i)} \in \mathcal{X}$  and  $y^{(i)} \in \mathcal{Y}$  for  $i = 1, \dots, n$ .

# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$
- A training set is set of pairs  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$   
s.t.  $x^{(i)} \in \mathcal{X}$  and  $y^{(i)} \in \mathcal{Y}$  for  $i = 1, \dots, n$ .
- Given a training set our goal is to produce a good prediction function  $h$

# Supervised Learning

- A hypothesis or a prediction function is function  $h : \mathcal{X} \rightarrow \mathcal{Y}$
- A training set is set of pairs  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$   
s.t.  $x^{(i)} \in \mathcal{X}$  and  $y^{(i)} \in \mathcal{Y}$  for  $i = 1, \dots, n$ .
- Given a training set our goal is to produce a good prediction function  $h$
- If  $\mathcal{Y}$  is continuous, then called a regression problem
- If  $\mathcal{Y}$  is discrete, then called a classification problem

*→ GPT →*

# Supervised Learning

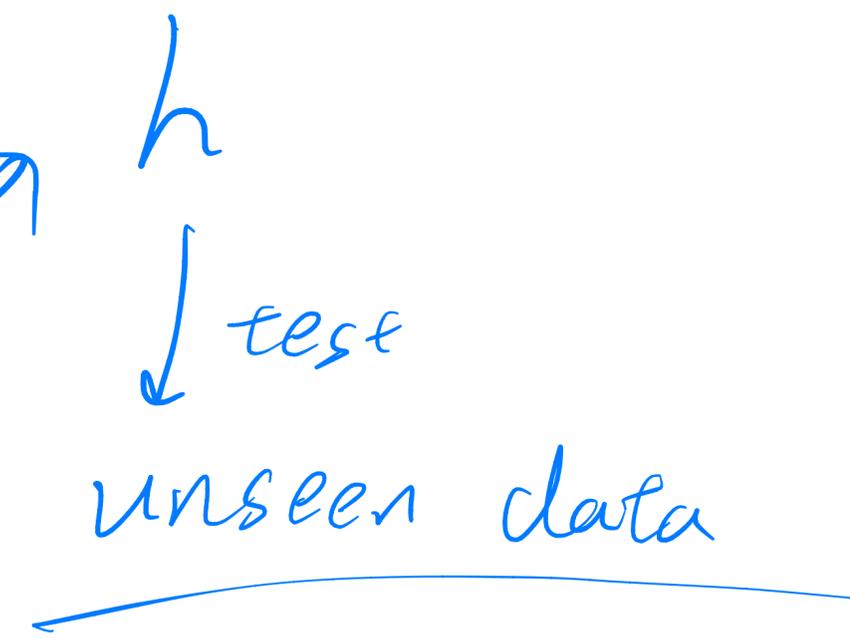
- How to define “good” for a prediction function?
  - Metrics / performance
  - Good on unseen data

Validation dataset is another set of pairs  $\{(\hat{x}^{(1)}, \hat{y}^{(1)}), \dots, (\hat{x}^{(m)}, \hat{y}^{(m)})\}$

Does not overlap with training dataset

# Supervised Learning

- How to define “good” for a prediction function?
  - Metrics / performance
  - Good on unseen data



Validation dataset is another set of pairs  $\{(\hat{x}^{(1)}, \hat{y}^{(1)}), \dots, (\hat{x}^{(m)}, \hat{y}^{(m)})\}$

Does not overlap with training dataset

develop

Test dataset is another set of pairs  $\{(\tilde{x}^{(1)}, \tilde{y}^{(1)}), \dots, (\tilde{x}^{(L)}, \tilde{y}^{(L)})\}$

Does not overlap with training and validation dataset

# Supervised Learning

● How to define “good” for a prediction function?

- Metrics / performance
- Good on unseen data

Validation dataset is another set of pairs  $\{(\hat{x}^{(1)}, \hat{y}^{(1)}), \dots, (\hat{x}^{(m)}, \hat{y}^{(m)})\}$

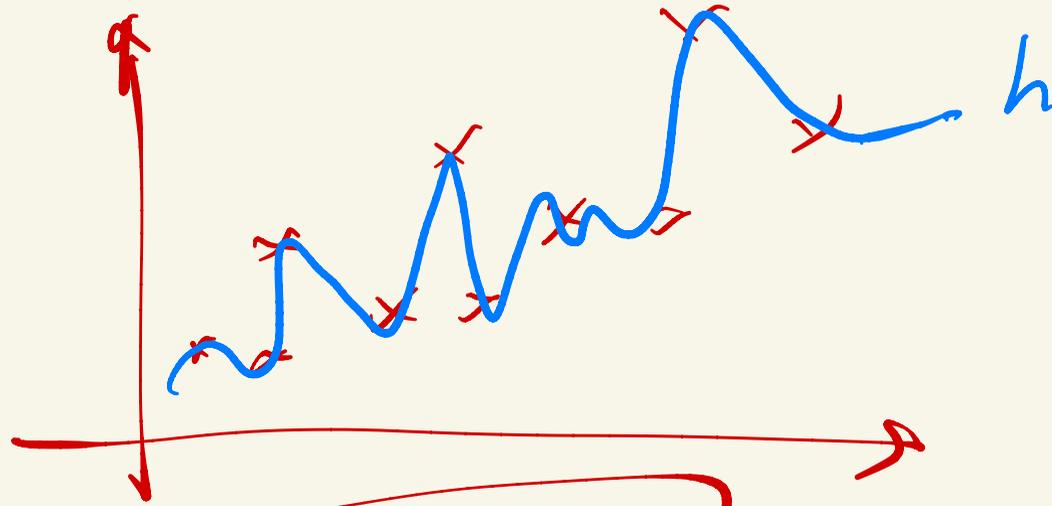
Does not overlap with training dataset

Test dataset is another set of pairs  $\{(\tilde{x}^{(1)}, \tilde{y}^{(1)}), \dots, (\tilde{x}^{(L)}, \tilde{y}^{(L)})\}$

Does not overlap with training and validation dataset

Completely unseen before deployment

Realistic setting



$$y = h(x)$$

linear?

$x^2, x^3, \frac{x^4}{\dots}$   
polynomial. what order

# Supervised Learning

- How to define “good” for a prediction function?
  - Metrics / performance
  - Good on unseen data

Validation dataset is another set of pairs  $\{(\hat{x}^{(1)}, \hat{y}^{(1)}), \dots, (\hat{x}^{(m)}, \hat{y}^{(m)})\}$

Does not overlap with training dataset

Test dataset is another set of pairs  $\{(\tilde{x}^{(1)}, \tilde{y}^{(1)}), \dots, (\tilde{x}^{(L)}, \tilde{y}^{(L)})\}$

Does not overlap with training and validation dataset

Completely unseen before deployment

*parameter*  
**Hyperparameter tuning is a form of training**

Realistic setting

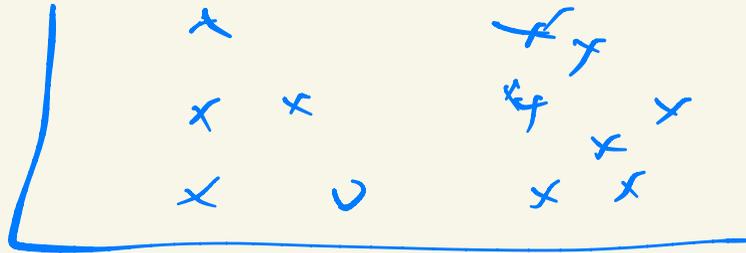
*test loss*

*optimizer*

*human*

*brain space*

# 1. unsupervised training



kmeans ?



## 2. don't care about unseen data

1. Solve a mathematical hypothesis

2. algorithm operator

→ kmeans *improved*  
box

# Supervised Training

Generalization



Train



Validation



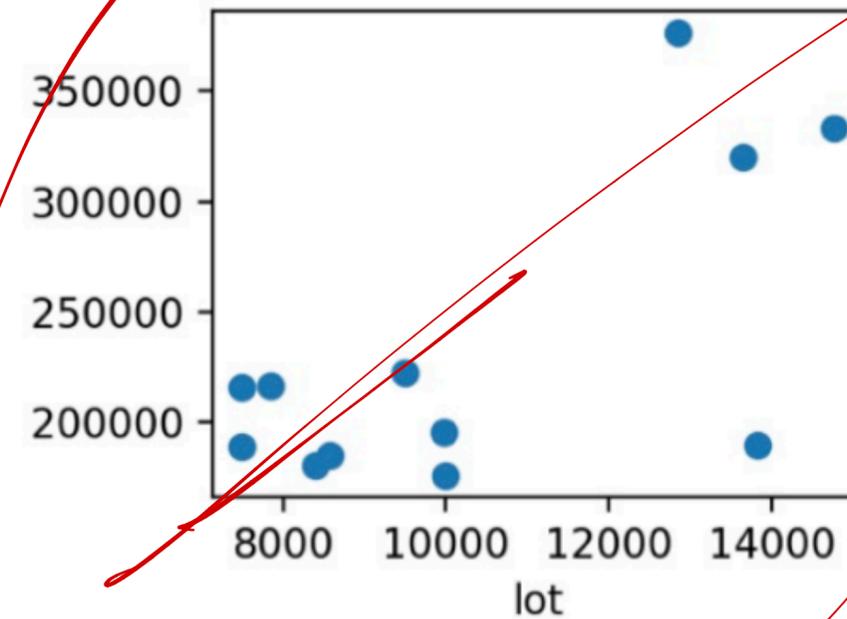
Test

Not only for supervised learning

# Example: Regression using Housing Data

# Example Housing Data

	SalePrice	Lot.Area
4	189900	13830
5	195500	9978
9	189000	7500
10	175900	10000
12	180400	8402
22	216000	7500
36	376162	12858
47	320000	13650
55	216500	7851
56	185088	8577



# Represent $h$ as a Linear Function

$h(x) = \theta_0 + \theta_1 x_1$  is an *affine function*

↗ affine

Popular choice

$$h(x) = \theta_1 x_1 + \theta_0$$

# Represent $h$ as a Linear Function

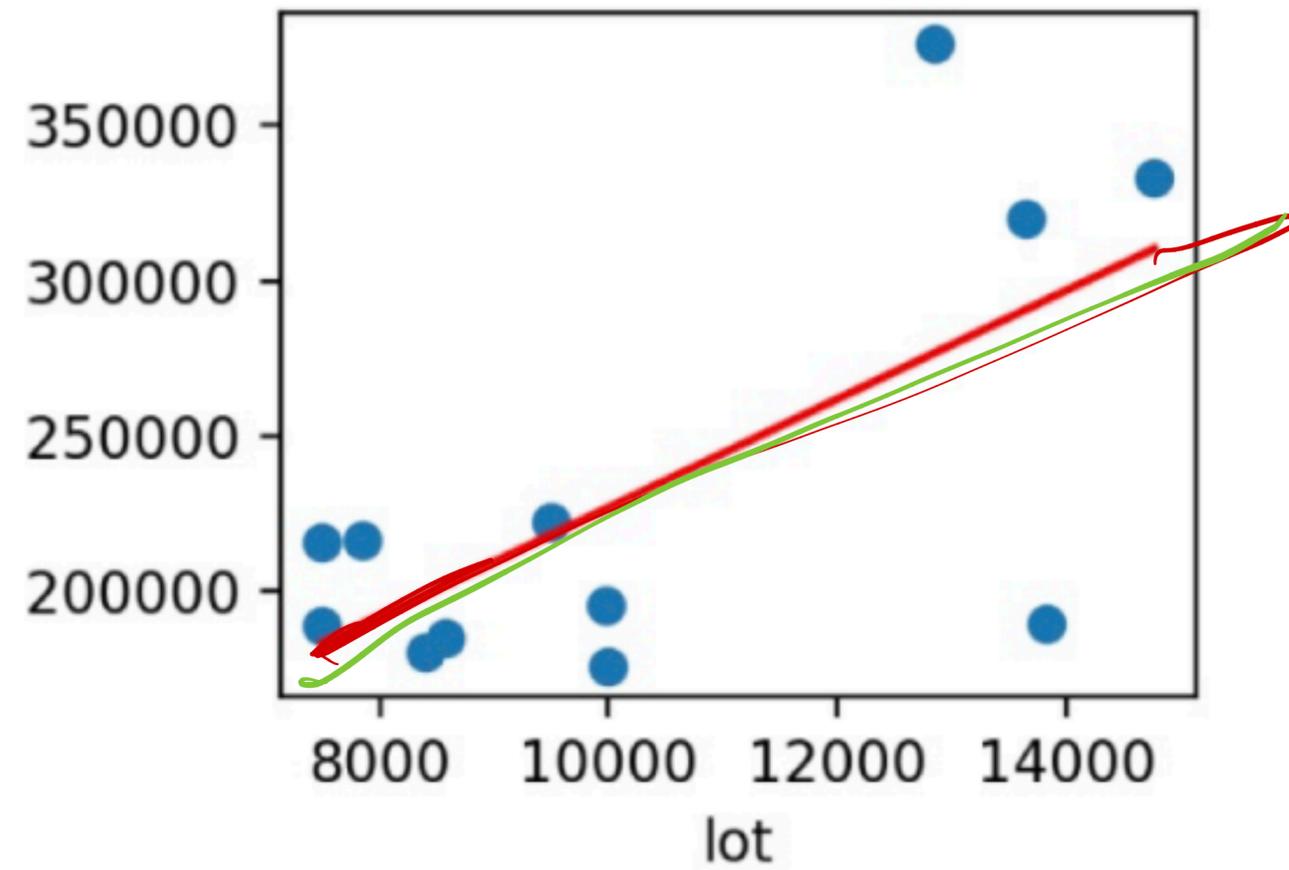
$h(x) = \theta_0 + \theta_1 x_1$  is an affine function

Popular choice

The function is defined by **parameters**  $\theta_0$  and  $\theta_1$ , the function space is greatly reduced

learning

# Simple Line Fit



# More Features

	size	bedrooms	lot size		Price
$x^{(1)}$	2104	4	45k	$y^{(1)}$	400
$x^{(2)}$	2500	3	30k	$y^{(2)}$	900



# More Features

	size	bedrooms	lot size		Price
$x^{(1)}$	2104	4	45k	$y^{(1)}$	400
$x^{(2)}$	2500	3	30k	$y^{(2)}$	900

What's a prediction here?

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3.$$

# More Features

	size	bedrooms	lot size		Price
$x^{(1)}$	2104	4	45k	$y^{(1)}$	400
$x^{(2)}$	2500	3	30k	$y^{(2)}$	900

What's a prediction here?

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3.$$

With the convention that  $x_0 = 1$  we can write:

$$h(x) = \sum_{j=0}^3 \theta_j x_j$$

$x_0 = 1$

$$h(x) = \sum_{j=0}^3 \theta_j x_j$$

# Vector Notations

# Vector Notations

	size	bedrooms	lot size		Price
$x^{(1)}$	2104	4	45k	$y^{(1)}$	400
$x^{(2)}$	2500	3	30k	$y^{(2)}$	900

We write the vectors as (important notation)

$$\vec{\theta} = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{pmatrix} \text{ and } \vec{x}^{(1)} = \begin{pmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{pmatrix} = \begin{pmatrix} 1 \\ 2104 \\ 4 \\ 45 \end{pmatrix} \text{ and } \underline{y^{(1)} = 400}$$

# Vector Notations

	size	bedrooms	lot size		Price
$x^{(1)}$	2104	4	45k	$y^{(1)}$	400
$x^{(2)}$	2500	3	30k	$y^{(2)}$	900

We write the vectors as (important notation)

$$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{pmatrix} \text{ and } x^{(1)} = \begin{pmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{pmatrix} = \begin{pmatrix} 1 \\ 2104 \\ 4 \\ 45 \end{pmatrix} \text{ and } y^{(1)} = 400$$

We call  $\theta$  **parameters**,  $x^{(i)}$  is the input or the **features**, and the output or **target** is  $y^{(i)}$ . To be clear,

$(x, y)$  is a training example and  $(x^{(i)}, y^{(i)})$  is the  $i^{\text{th}}$  example.

# Vector Notations

	size	bedrooms	lot size		Price
$x^{(1)}$	2104	4	45k	$y^{(1)}$	400
$x^{(2)}$	2500	3	30k	$y^{(2)}$	900

We write the vectors as (important notation)

$$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{pmatrix} \text{ and } x^{(1)} = \begin{pmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{pmatrix} = \begin{pmatrix} 1 \\ 2104 \\ 4 \\ 45 \end{pmatrix} \text{ and } y^{(1)} = 400$$

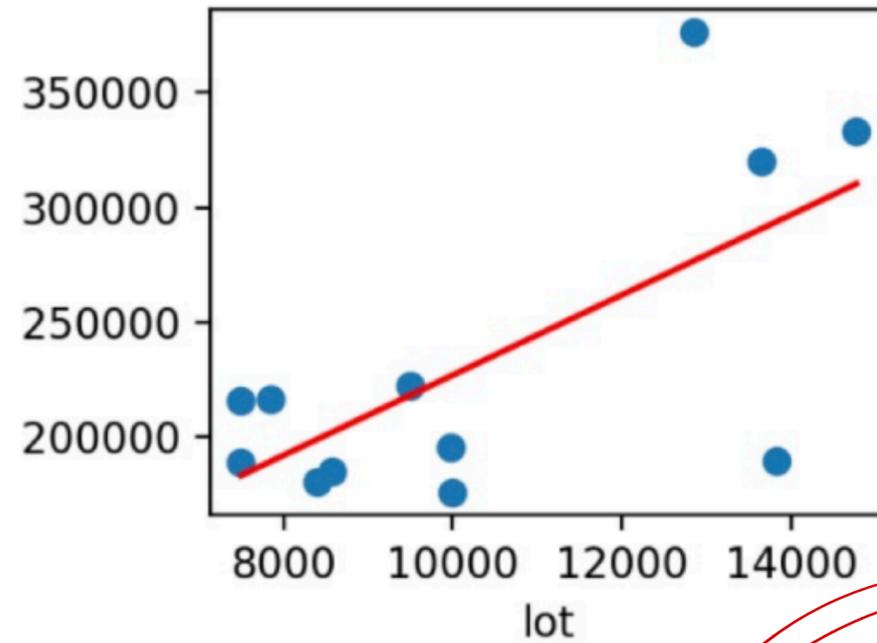
We call  $\theta$  **parameters**,  $x^{(i)}$  is the input or the **features**, and the output or **target** is  $y^{(i)}$ . To be clear,

$(x, y)$  is a training example and  $(x^{(i)}, y^{(i)})$  is the  $i^{\text{th}}$  example.

We have  $n$  examples. There are  $d$  features.  $x^{(i)}$  and  $\theta$  are  $d+1$  dimensional (since  $x_0 = 1$ )

$\theta_0$   
 $x_0 = 1$

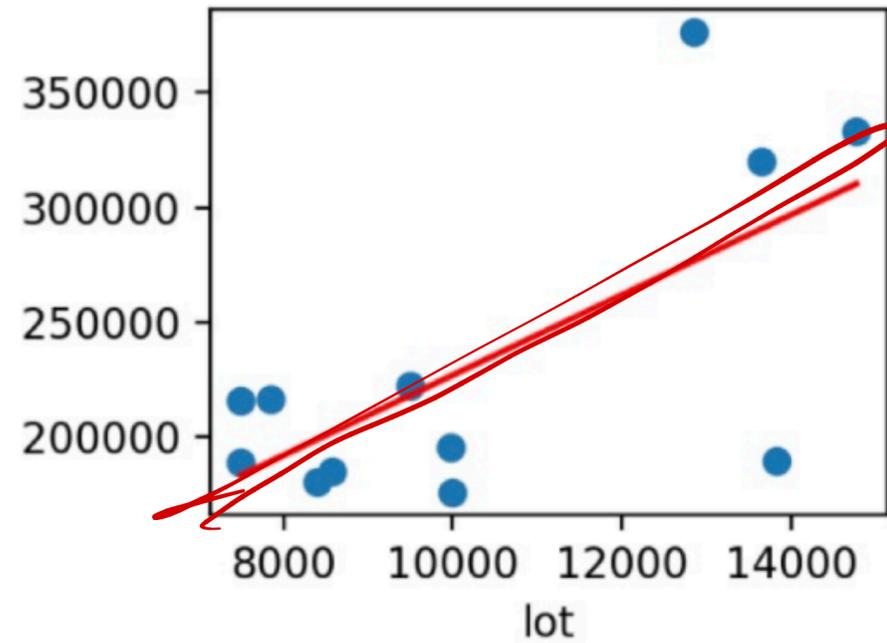
# Vector Notation of Prediction



$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$



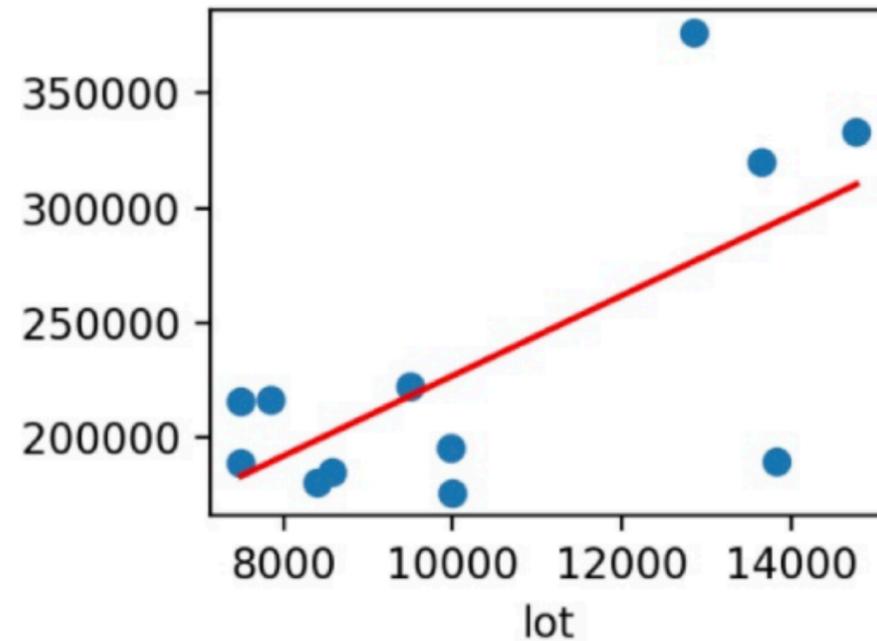
# Vector Notation of Prediction



$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

We want to choose  $\theta$  so that  $h_{\theta}(x) \approx y$

# Loss Function



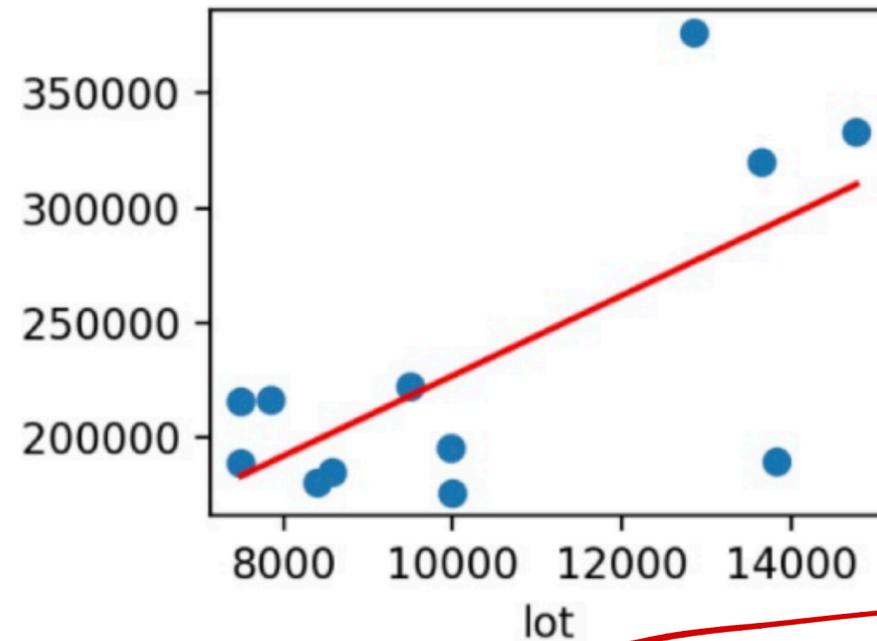
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

We want to choose  $\theta$  so that  $h_{\theta}(x) \approx y$

How to quantify the deviation of  $h_{\theta}(x)$  from  $y$

*metric*

# Least Squares

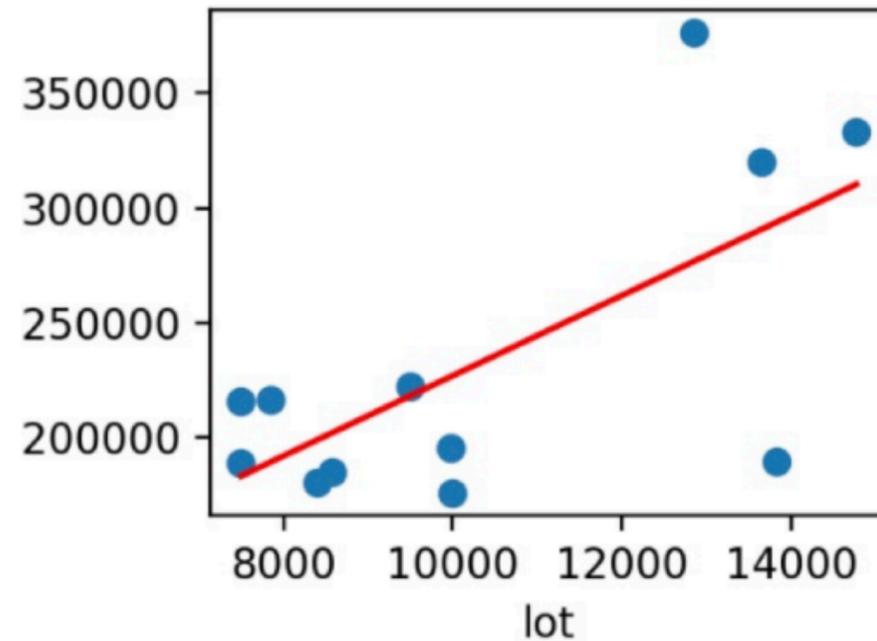


$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( \underbrace{h_{\theta}(x^{(i)})}_{\text{red underline}} - \underbrace{y^{(i)}}_{\text{red underline}} \right)^2$$

$$= \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

# Least Squares



$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Choose

$$\theta = \underset{\theta}{\operatorname{argmin}} J(\theta).$$

# Solving Least Square Problem

Direct Minimization

$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( \underbrace{h_{\theta}(x^{(i)})}_{\theta} - \underbrace{y^{(i)}} \right)^2$$

Choose

$$\theta = \underset{\theta}{\operatorname{argmin}} J(\theta).$$

$$\nabla_{\theta} J(\theta) = 0$$

$$\frac{d J(\theta)}{d \theta} = 0$$

# Solving Least Square Problem

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\ &= \frac{1}{2} \nabla_{\theta} ((X\theta)^T X\theta - (X\theta)^T \vec{y} - \vec{y}^T (X\theta) + \vec{y}^T \vec{y}) \\ &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X)\theta - \vec{y}^T (X\theta) - \vec{y}^T (X\theta)) \\ &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X)\theta - 2(X^T \vec{y})^T \theta) \\ &= \frac{1}{2} (2X^T X\theta - 2X^T \vec{y}) \\ &= X^T X\theta - X^T \vec{y} \quad \Rightarrow 0\end{aligned}$$

$$X^T X\theta = X^T \vec{y}$$

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

# Solving Least Square Problem

$$\begin{aligned}\nabla_{\theta} J(\theta) &= \nabla_{\theta} \frac{1}{2} (X\theta - \vec{y})^T (X\theta - \vec{y}) \\ &= \frac{1}{2} \nabla_{\theta} ((X\theta)^T X\theta - (X\theta)^T \vec{y} - \vec{y}^T (X\theta) + \vec{y}^T \vec{y}) \\ &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X)\theta - \vec{y}^T (X\theta) - \vec{y}^T (X\theta)) \\ &= \frac{1}{2} \nabla_{\theta} (\theta^T (X^T X)\theta - 2(X^T \vec{y})^T \theta) \\ &= \frac{1}{2} (2X^T X\theta - 2X^T \vec{y}) \\ &= X^T X\theta - X^T \vec{y}\end{aligned}$$

Normal equations

$$X^T X\theta = X^T \vec{y}$$

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

$\text{rank}(X^T X) \leq \min(\text{rank}(X^T), \text{rank}(X)) = \text{rank}(X)$

# Solving Least Square Problem

$d \times d$

necessary condition:  
 $\text{rank}(X) \geq d$

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

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

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

$\text{rank}(X) \leq \min(n, d)$

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

$[n < d]$

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

invertible  $\iff$  full rank

$$= X^T X\theta - X^T \vec{y}$$

rank

$n$ : # samples

$d$ : # features

Normal equations  $X^T X\theta = X^T \vec{y}$

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

When is  $X^T X$  invertible? What if it is not invertible?

$d+1$

$\theta_1 + \theta_2 + \dots + \theta_{d+1} = 0$

$\theta_1 + \dots + \theta_d = 0$

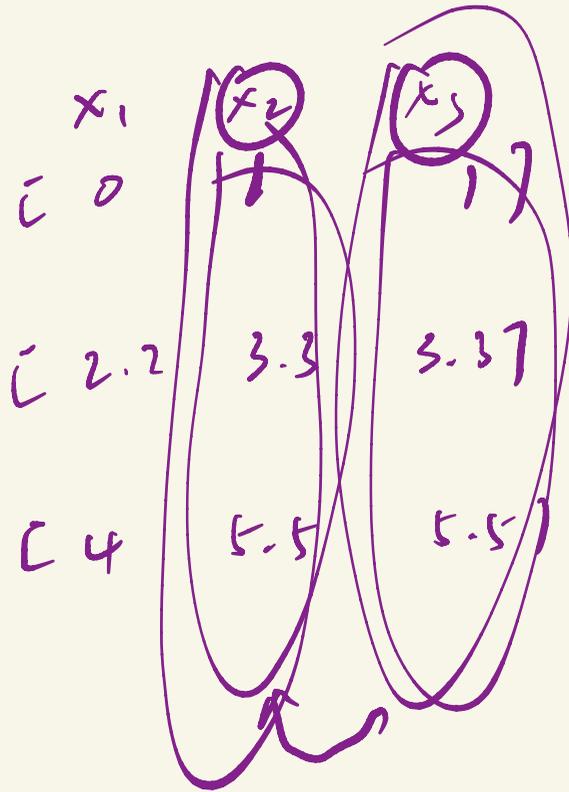
$\theta_1 + \dots + \theta_{d-1} = 0$

$n$

sufficient

$\text{rank}(x) < d+1$

$n < d+1$



$$\boxed{X \in \mathbb{R}^{n \times d}}$$

↓

# Why Least-Square Loss Function?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

$$\sum_{i=1}^n |h_{\theta}(x^{(i)}) - y^{(i)}|$$

$$\sum_{i=1}^n |h_{\theta}(x^{(i)}) - y^{(i)}|^q$$

Maximum Likelihood

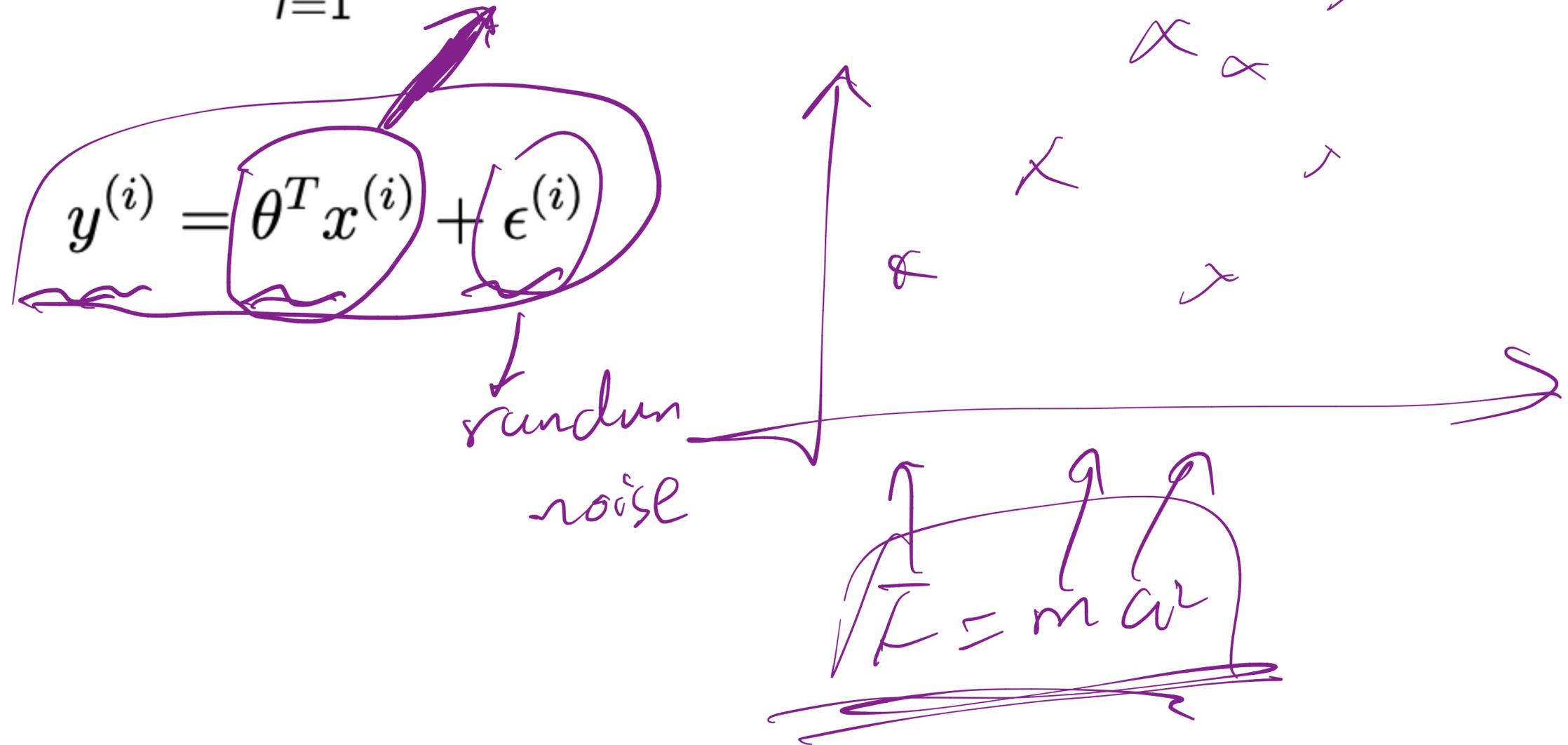
estimation

MLE

# Why Least-Square Loss Function?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Assume



# Why Least-Square Loss Function?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Assume

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$x, y$ : random variable

# Why Least-Square Loss Function?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Assume

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$x, y$ : random variable

$\epsilon$ : deviation of prediction from the truth, Gaussian random variable

# Why Least-Square Loss Function?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Assume

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$x, y$ : random variable

$\epsilon$ : deviation of prediction from the truth, Gaussian random variable

$x^{(i)}, y^{(i)}$ : observations, or the data



# Why Least-Square Loss Function?

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Assume

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$x, y$ : random variable

$\epsilon$ : deviation of prediction from the truth, Gaussian random variable

$x^{(i)}, y^{(i)}$ : observations, or the data

$\epsilon^{(i)}$ : the actual prediction error of the  $i_{th}$  example, sampled from the Gaussian distribution, IID (independently and identically distributed)

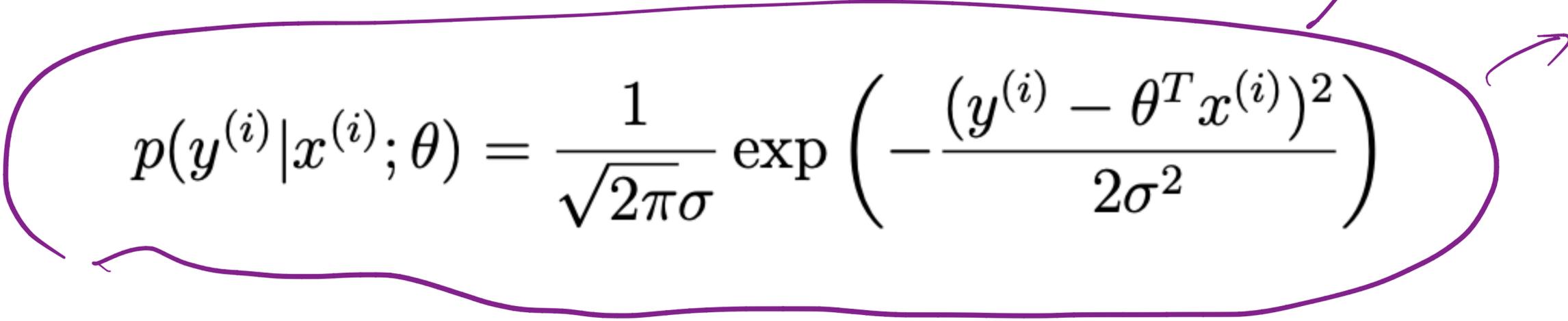
# Why Least-Square Loss Function?

# Why Least-Square Loss Function?

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

# Why Least-Square Loss Function?

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$


$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

# Why Least-Square Loss Function?

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$\begin{aligned} p(\vec{y} | X; \theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

# Why Least-Square Loss Function?

$$p(\epsilon^{(i)}) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(\epsilon^{(i)})^2}{2\sigma^2}\right)$$

$$p(y^{(i)} | x^{(i)}; \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$p(\vec{y} | X; \theta) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$$

Function of  $\theta$  =  $\prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$

# Why Least-Square Loss Function?

# Why Least-Square Loss Function?

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

# Why Least-Square Loss Function?

$$L(\theta) = \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta)$$

$$= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

*Gaussian*

Likelihood Function

# Why Least-Square Loss Function?

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned} \quad \text{Likelihood Function}$$

What is a reasonable guess of  $\theta$ ?



# Why Least-Square Loss Function?

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned}$$

Likelihood Function

What is a reasonable guess of  $\theta$ ?

Maximize the probability of Y's happening!

$\underset{\theta}{\operatorname{argmax}} L(\theta) = \underset{\theta}{\operatorname{argmax}} \log L(\theta)$

# Maximum Likelihood Estimation (MLE)

$$\begin{aligned}\ell(\theta) &= \log L(\theta) \\ &= \log \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= \sum_{i=1}^n \log \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \\ &= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{\sigma^2} \cdot \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2.\end{aligned}$$

↓ constant

↓ least square

# Why MLE?

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned} \quad \text{Likelihood Function}$$

What is a reasonable guess of  $\theta$ ?

Maximize the probability of Y's happening?

# Why MLE?

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned} \quad \text{Likelihood Function}$$

What is a reasonable guess of  $\theta$ ?

Maximize the probability of Y's happening?

Maximizing likelihood estimation  $\rightarrow \hat{\theta}$

# Why MLE?

$$\begin{aligned} L(\theta) &= \prod_{i=1}^n p(y^{(i)} | x^{(i)}; \theta) \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right) \end{aligned} \quad \text{Likelihood Function}$$

What is a reasonable guess of  $\theta$ ?

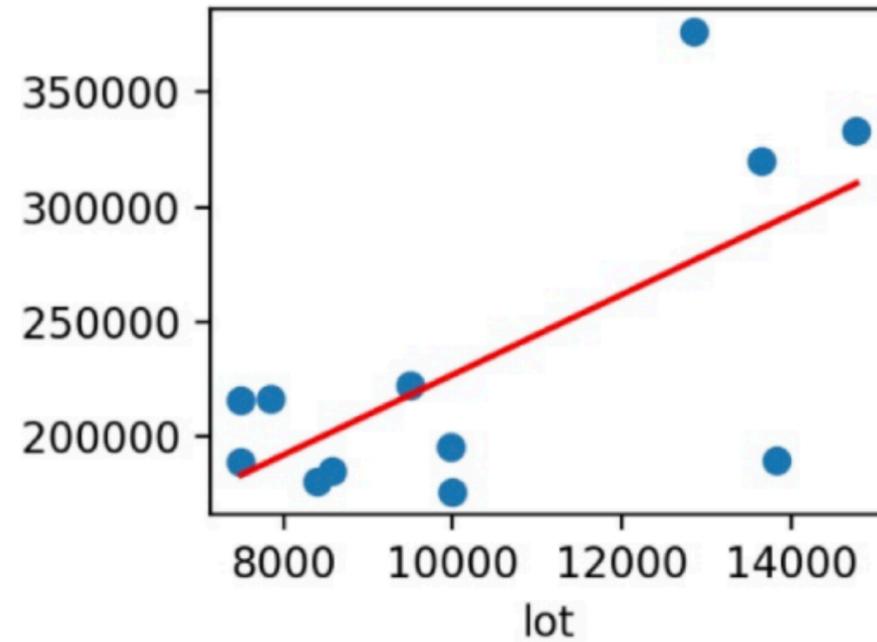
Maximize the probability of Y's happening?

Maximizing likelihood estimation  $\rightarrow \hat{\theta}$

Ground-truth  $\theta^*$



# Another Solution — Gradient Descent



$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Choose

$$\theta = \underset{\theta}{\operatorname{argmin}} J(\theta).$$

# Gradient Descent

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

# Gradient Descent

step size

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

direction

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

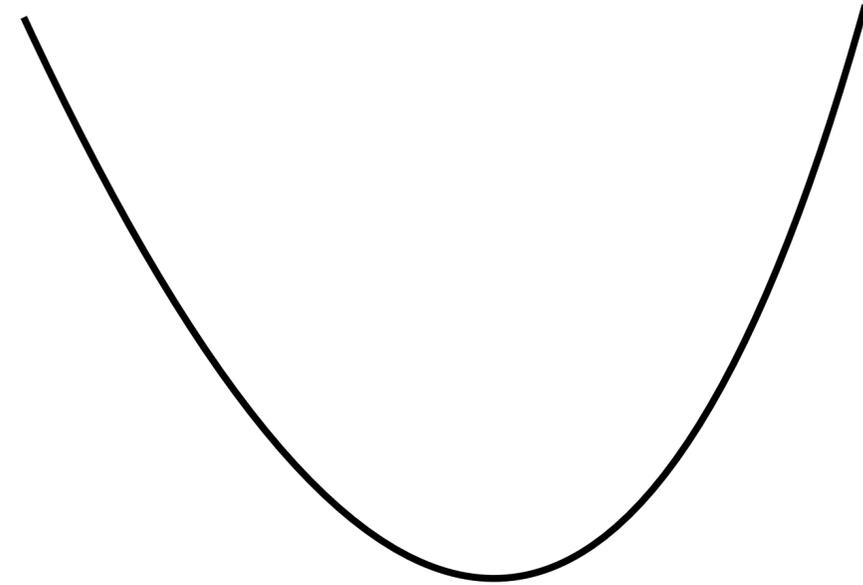
derivative

gradient

# Gradient Descent

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

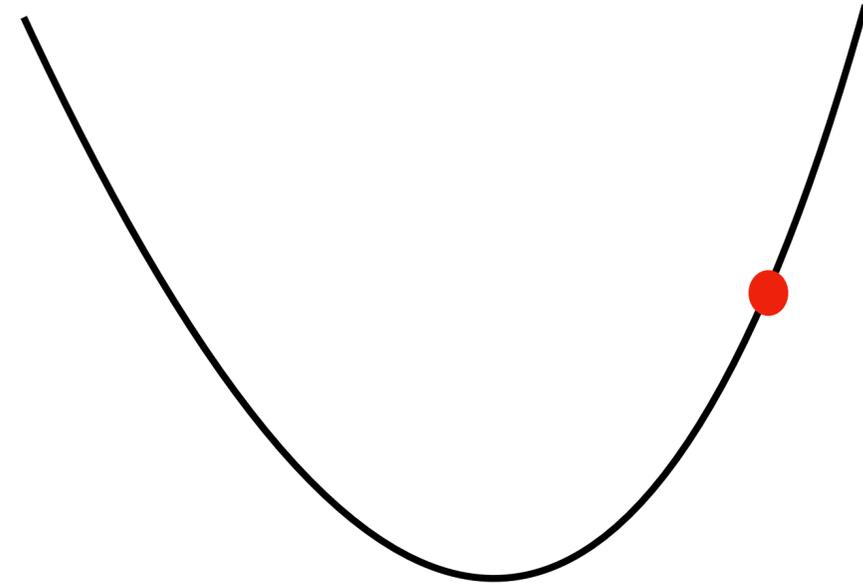
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$



# Gradient Descent

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

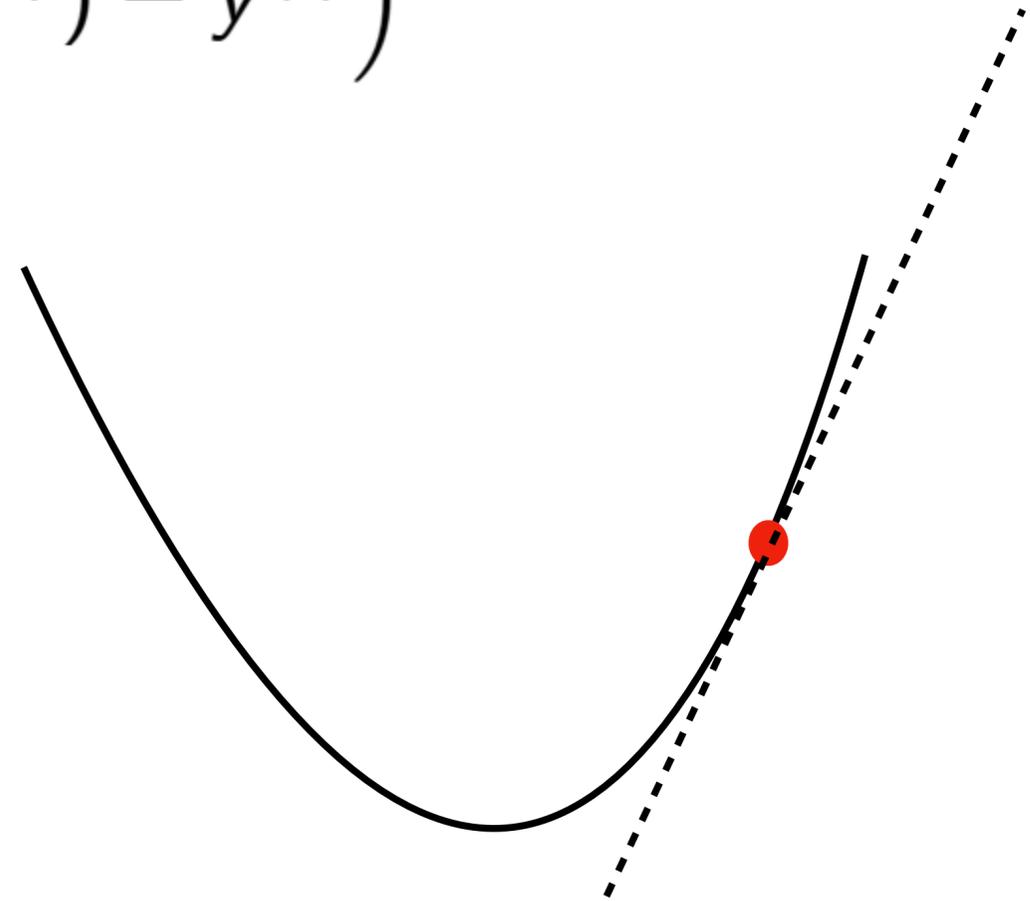
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$



# Gradient Descent

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

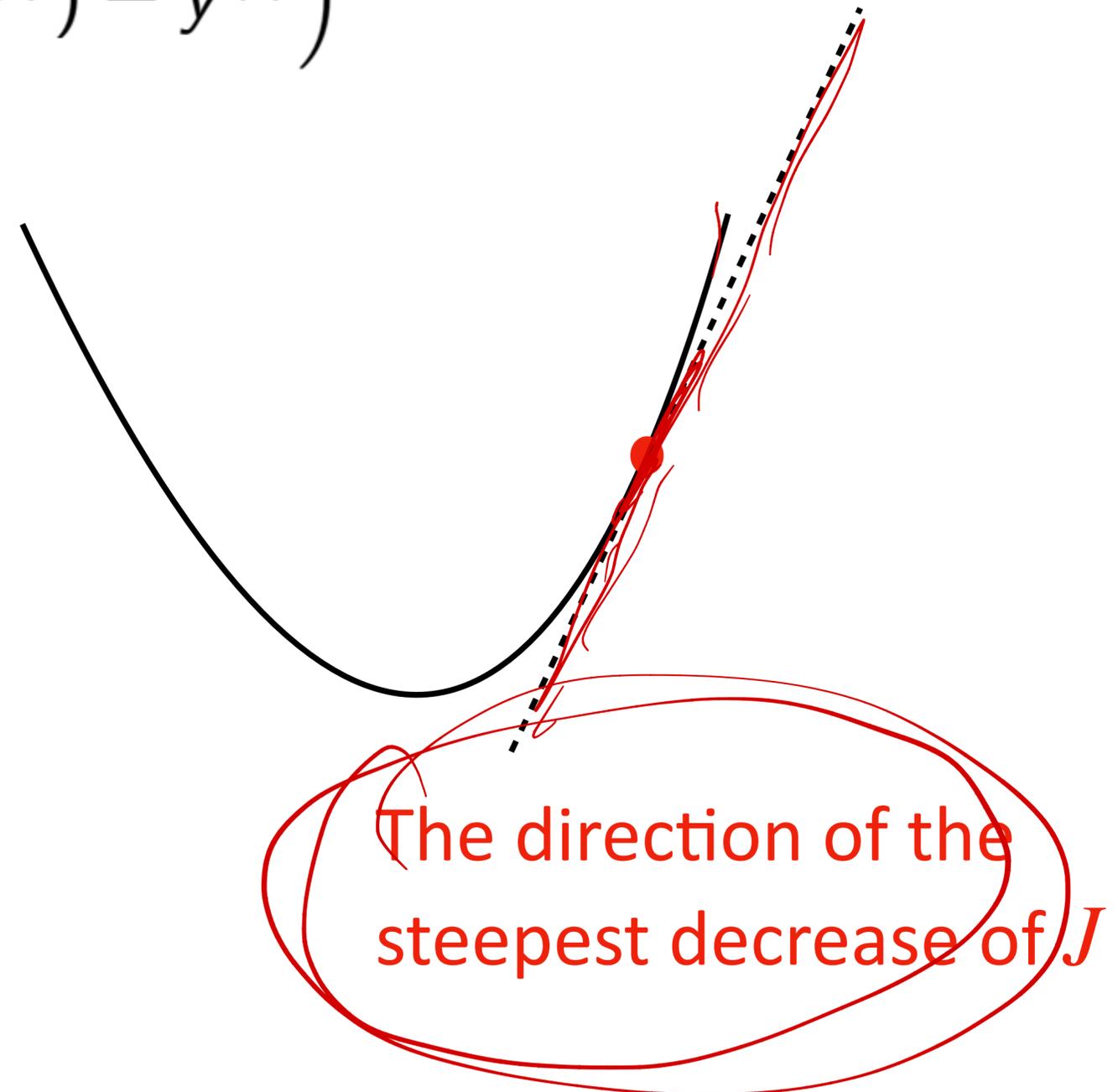
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$



# Gradient Descent

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

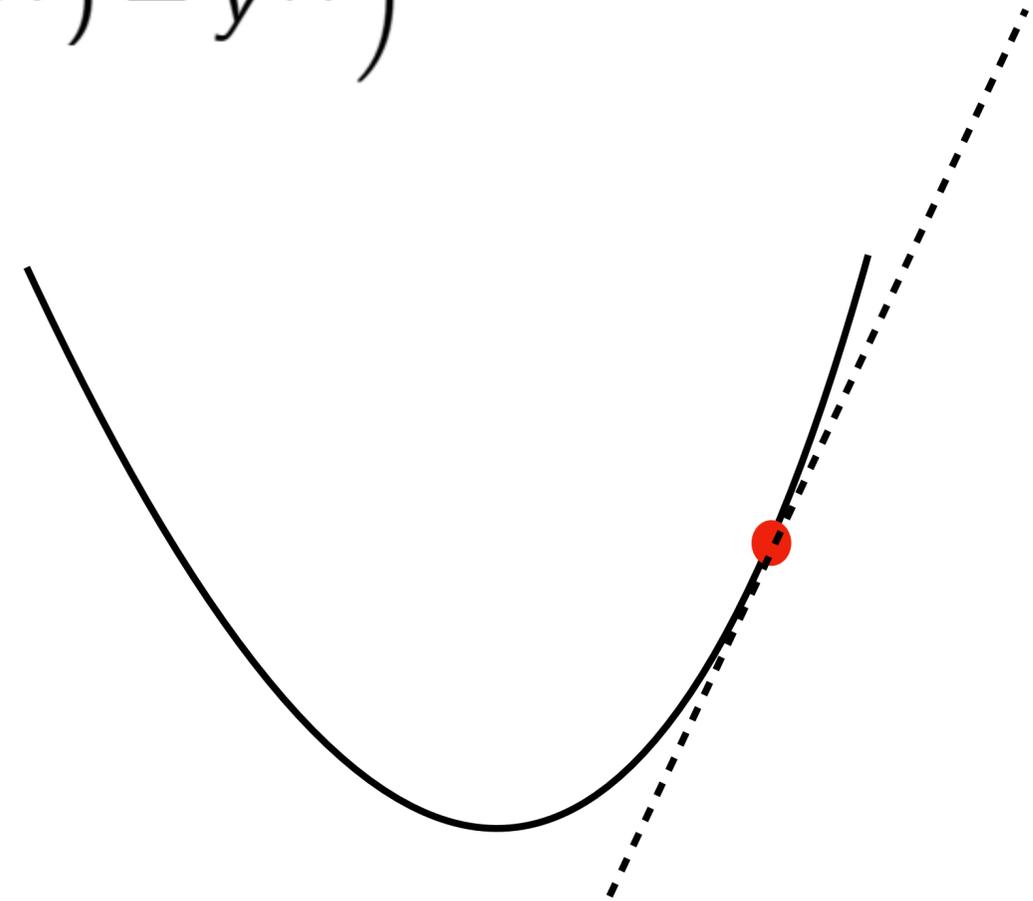


# Gradient Descent

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Learning Rate

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$



The direction of the  
steepest decrease of  $J$

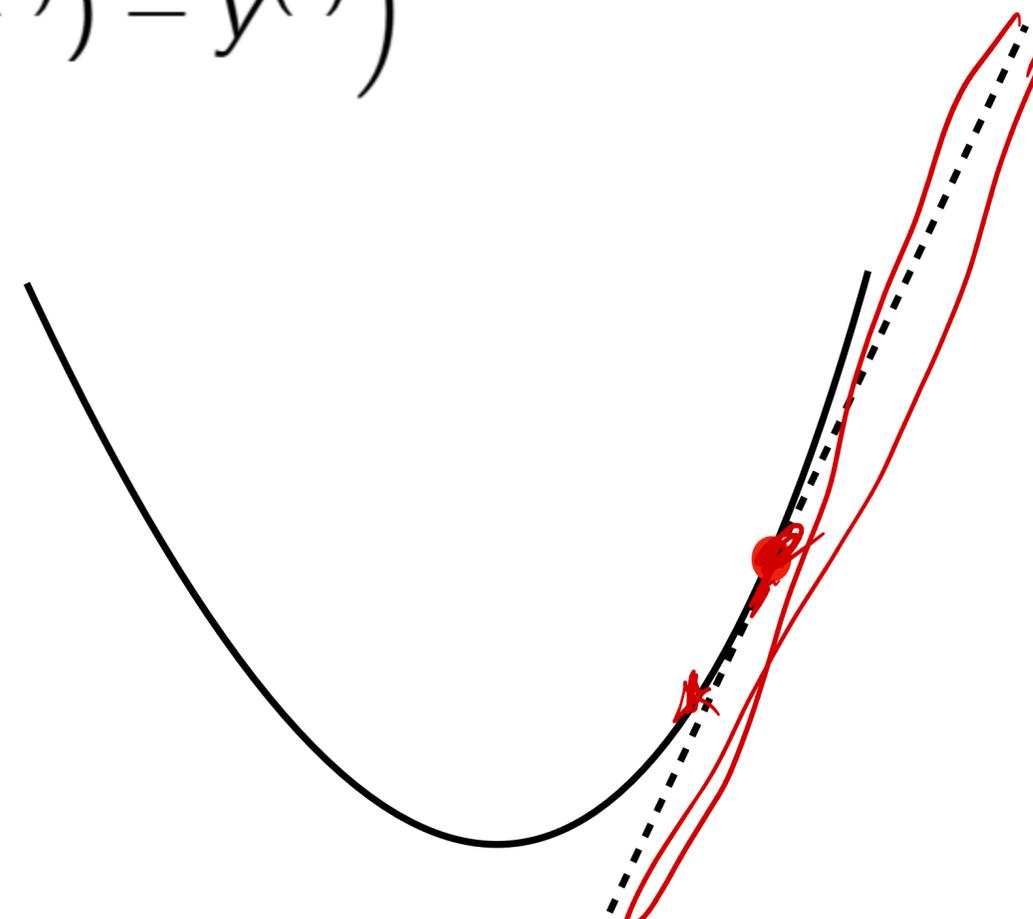
# Gradient Descent

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left( h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Learning Rate

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

This update is simultaneously performed for all values of  $j = 0, \dots, d$ .



The direction of the steepest decrease of  $J$

# Gradient Descent

For a single training example:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^d \theta_i x_i - y \right) \\ &= (h_{\theta}(x) - y) x_j\end{aligned}$$

# Gradient Descent

For a single training example:

$$\begin{aligned}\frac{\partial}{\partial \theta_j} J(\theta) &= \frac{\partial}{\partial \theta_j} \frac{1}{2} (h_{\theta}(x) - y)^2 \\ &= 2 \cdot \frac{1}{2} (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} (h_{\theta}(x) - y) \\ &= (h_{\theta}(x) - y) \cdot \frac{\partial}{\partial \theta_j} \left( \sum_{i=0}^d \theta_i x_i - y \right) \\ &= (h_{\theta}(x) - y) x_j\end{aligned}$$

LMS (Least Mean Square) Update Rule

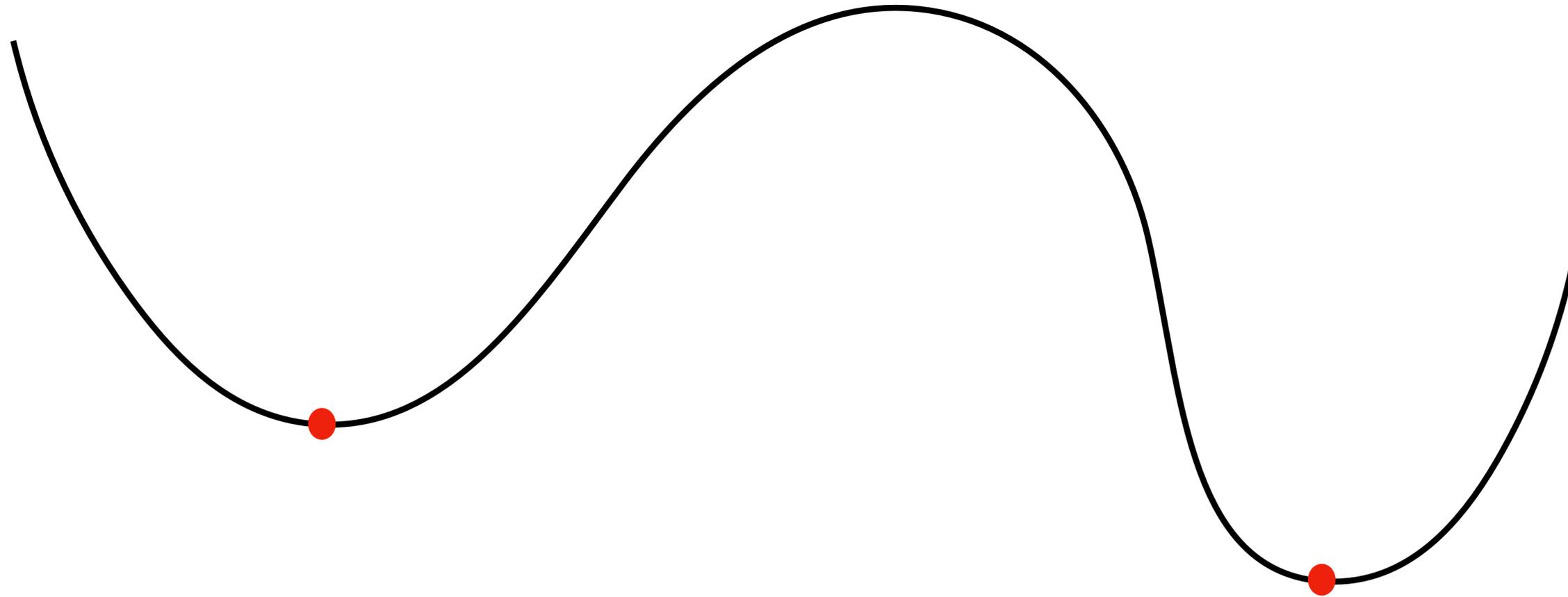
# Batch Gradient Descent

For a multiple training examples:

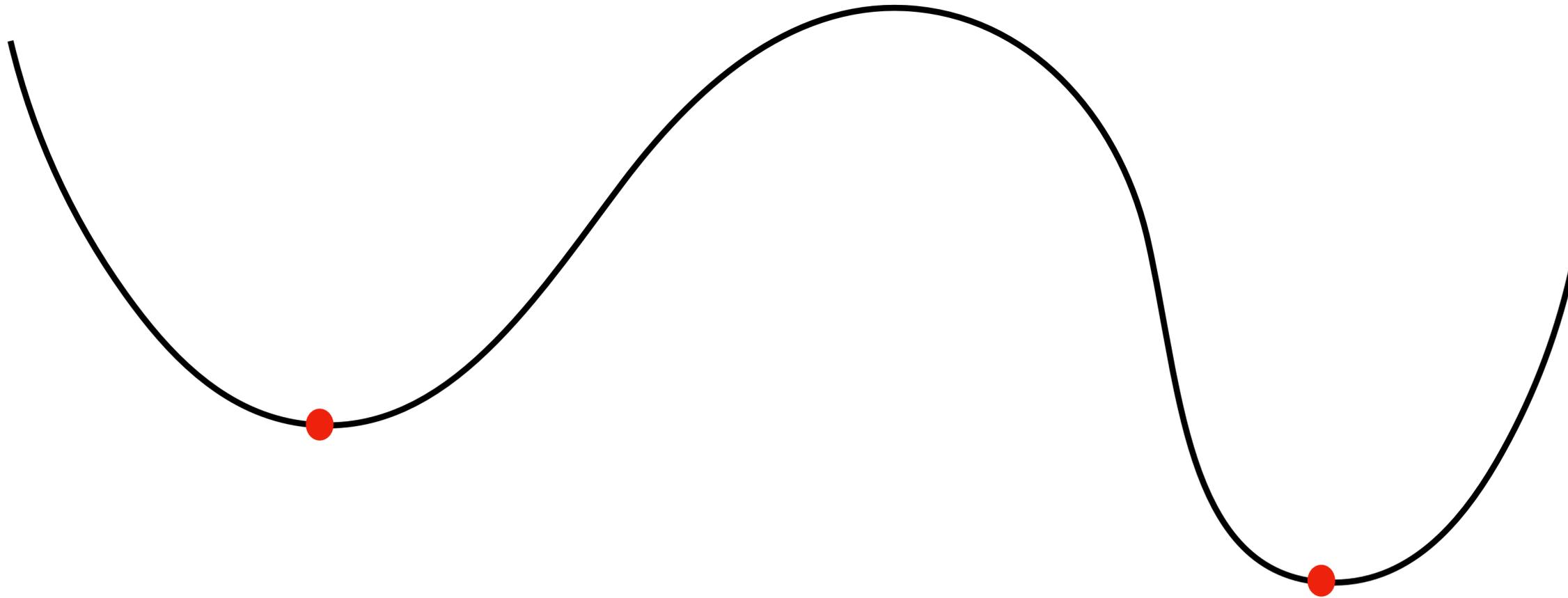
$$\theta_j := \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Repeat until convergence

# Local Minimum



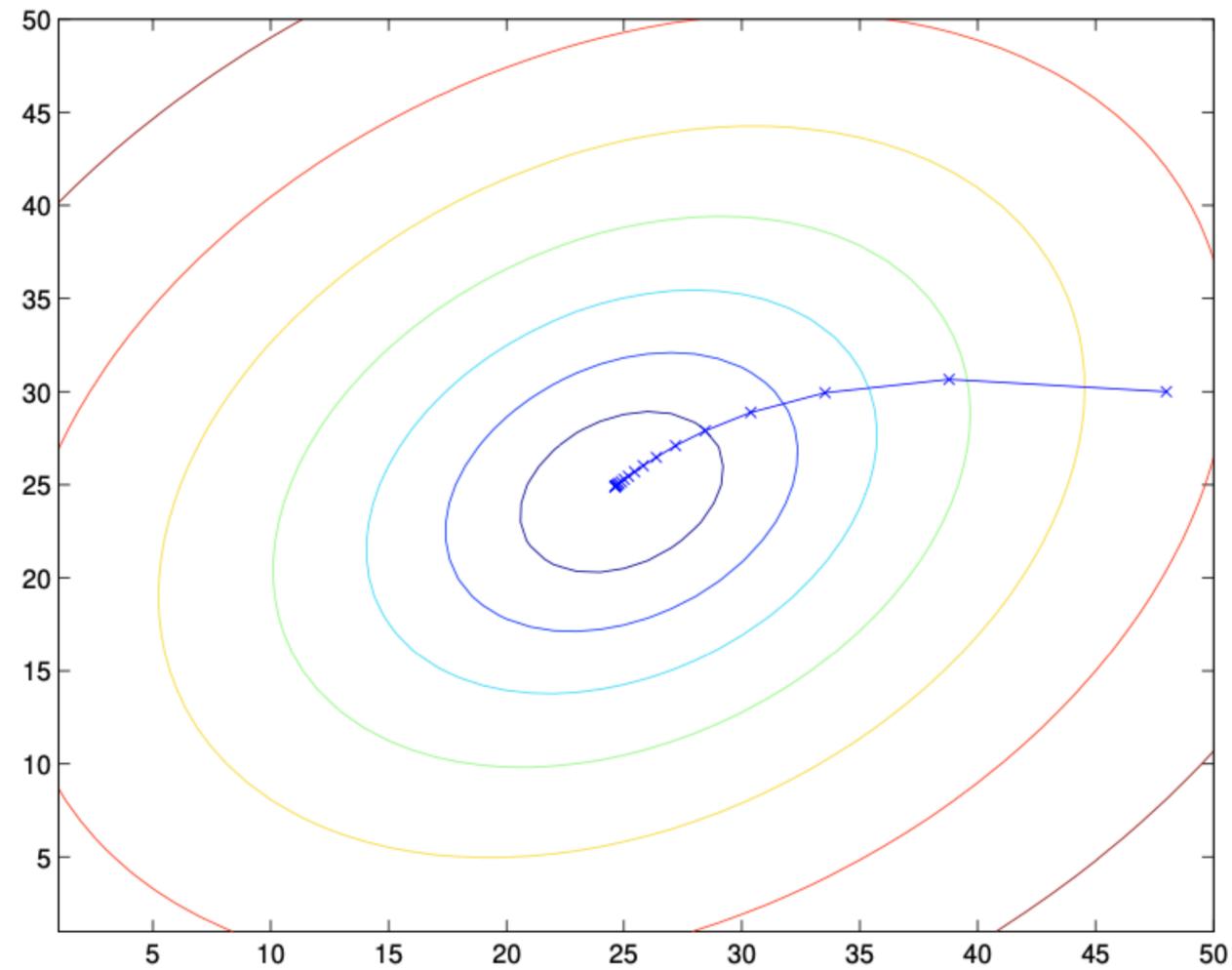
# Local Minimum



For least square optimization, are we likely to get local minima rather than the global minima through gradient descent?

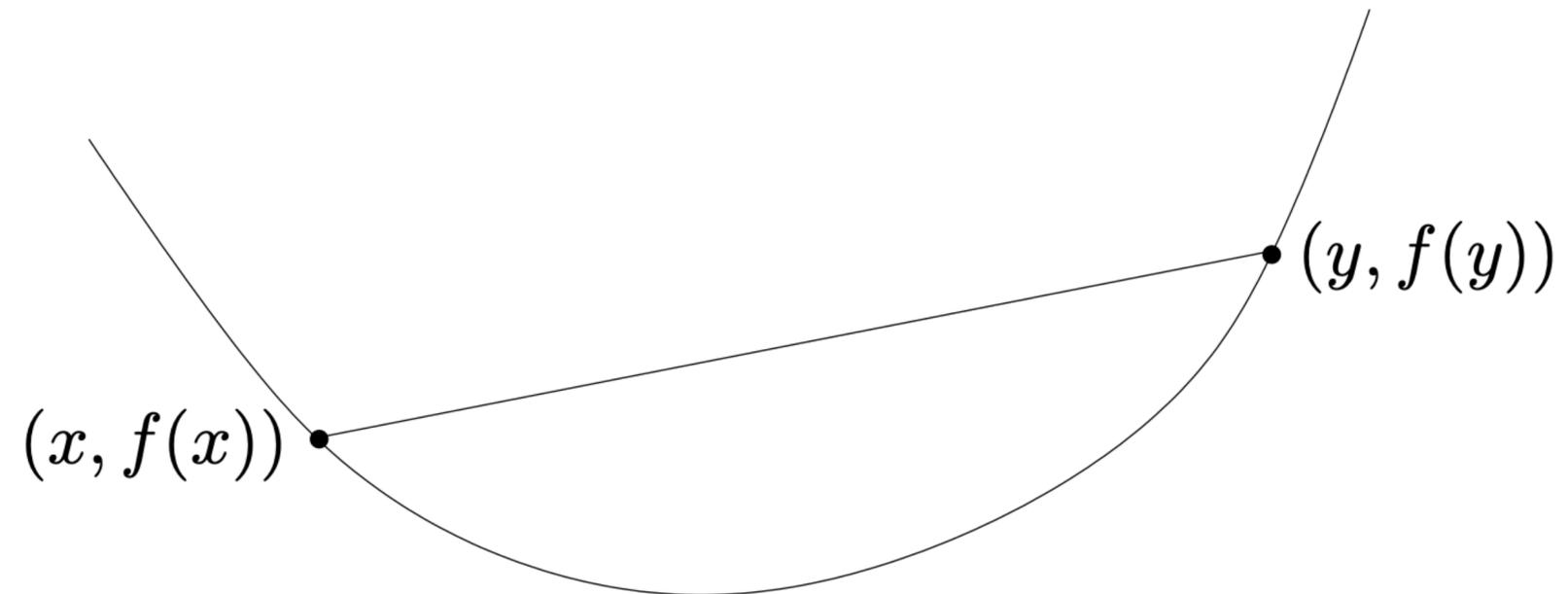
# $J$ is a convex quadratic function

There is only one local minima for  $J$



# Convex Function

$$f(tx + (1 - t)y) \leq tf(x) + (1 - t)f(y) \quad \text{for } 0 \leq t \leq 1$$



**Thank You!**  
**Q & A**