



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

COMP 5212
Machine Learning
Lecture 2

Supervised Learning: Regression

Junxian He
Sep 10, 2024

Announcement

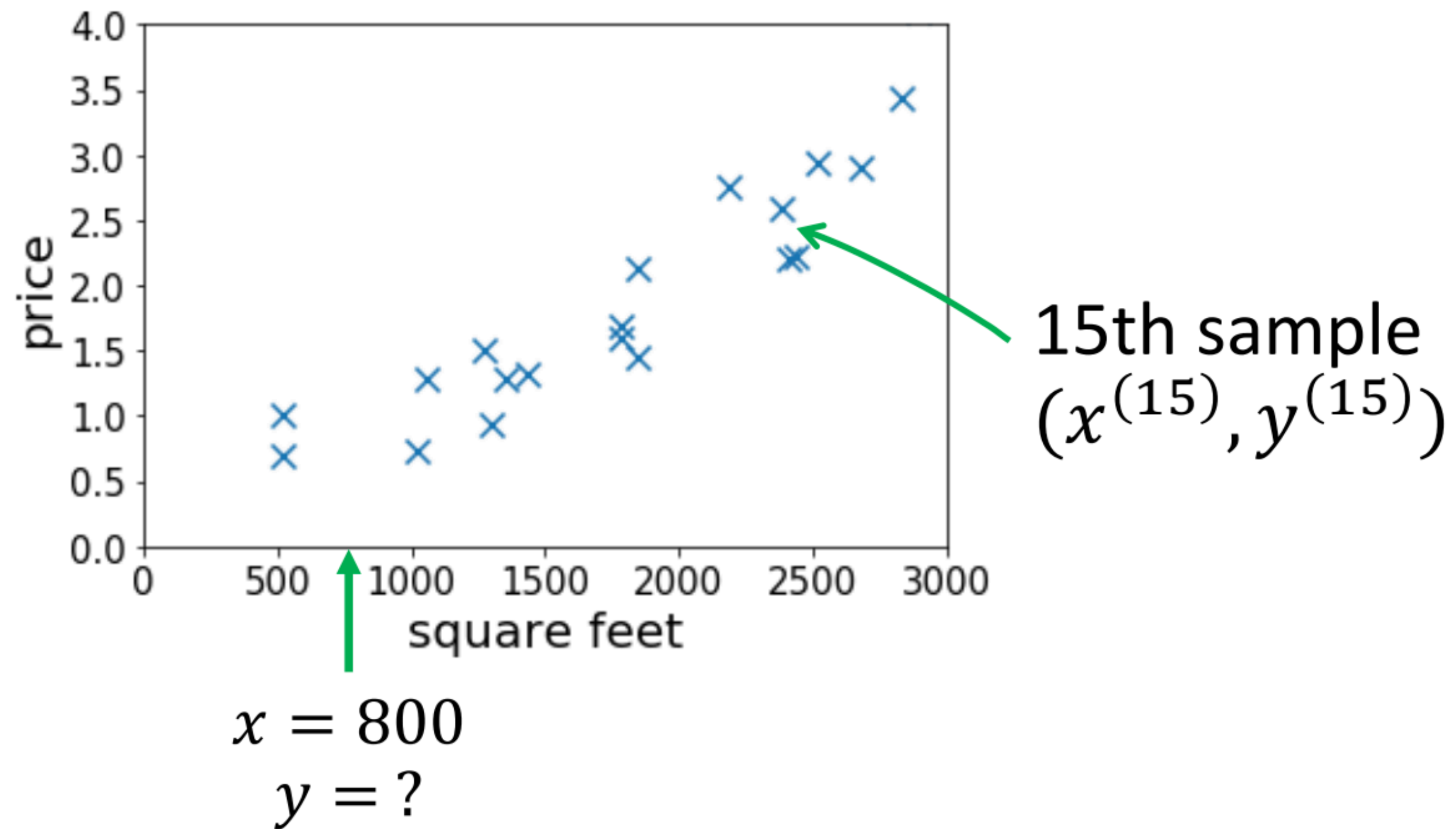
Lecture on Sep 17 (Mid-Autumn Festival) is rescheduled to Sep 23 (Monday) from 130pm - 250pm at LG3009.

Supervised Learning

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

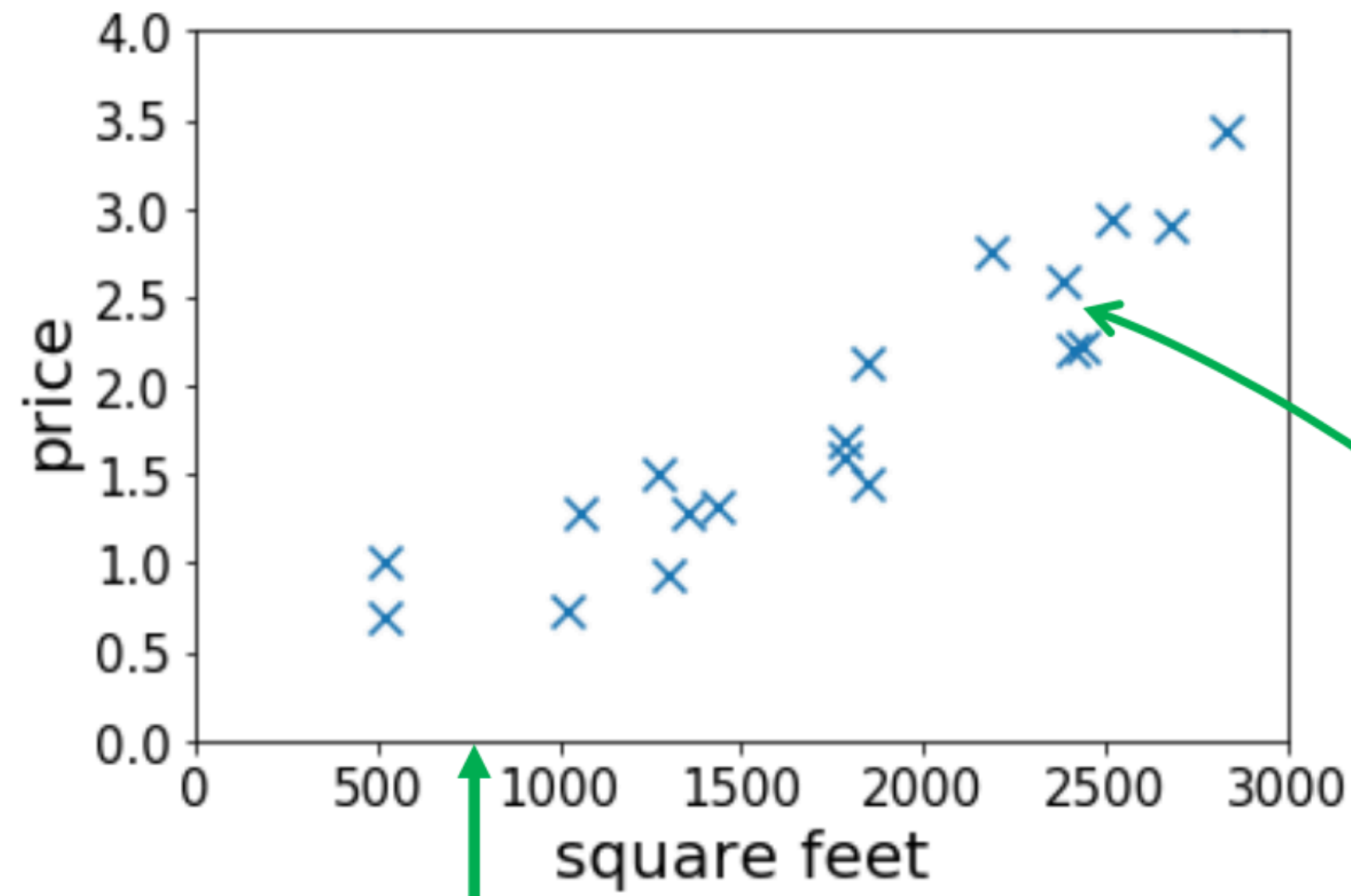
Supervised Learning

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



Supervised Learning

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



$x = 800$
 $y = ?$

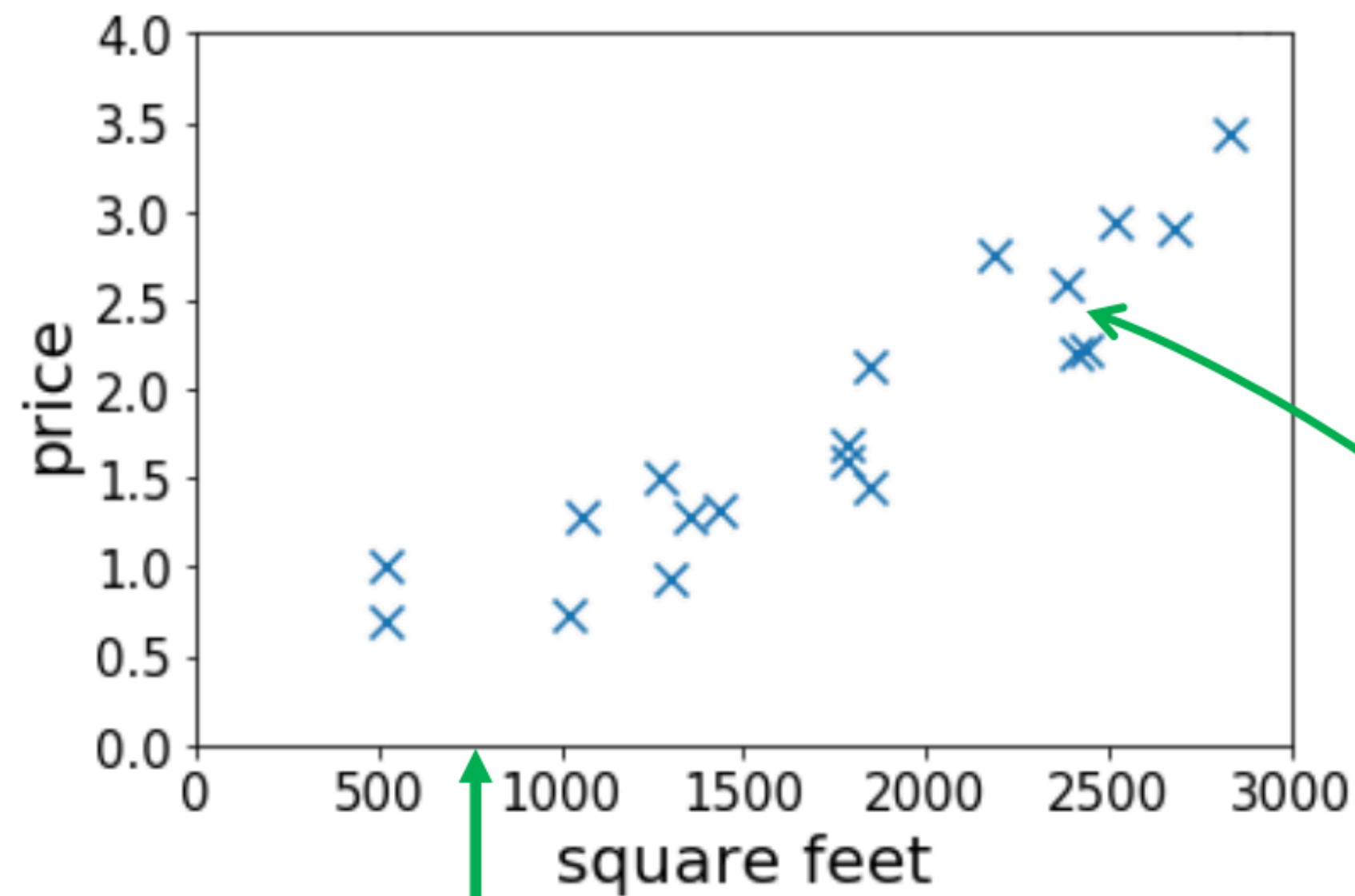
15th sample
 $(x^{(15)}, y^{(15)})$



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 $\{(x^{(1)}, y^{(1)}), \dots, (x^{(n)}, y^{(n)})\}$
s.t. $x^{(i)} \in \mathcal{X}$ and $y^{(i)} \in \mathcal{Y}$ for $i = 1, \dots, n$.

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

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

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

Supervised Learning



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

100K *90K → training* *10K → validation*
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

Hyperparameter tuning is a form of training *Realistic setting*

Supervised Training



Train



Validation



Test

Not only for supervised learning

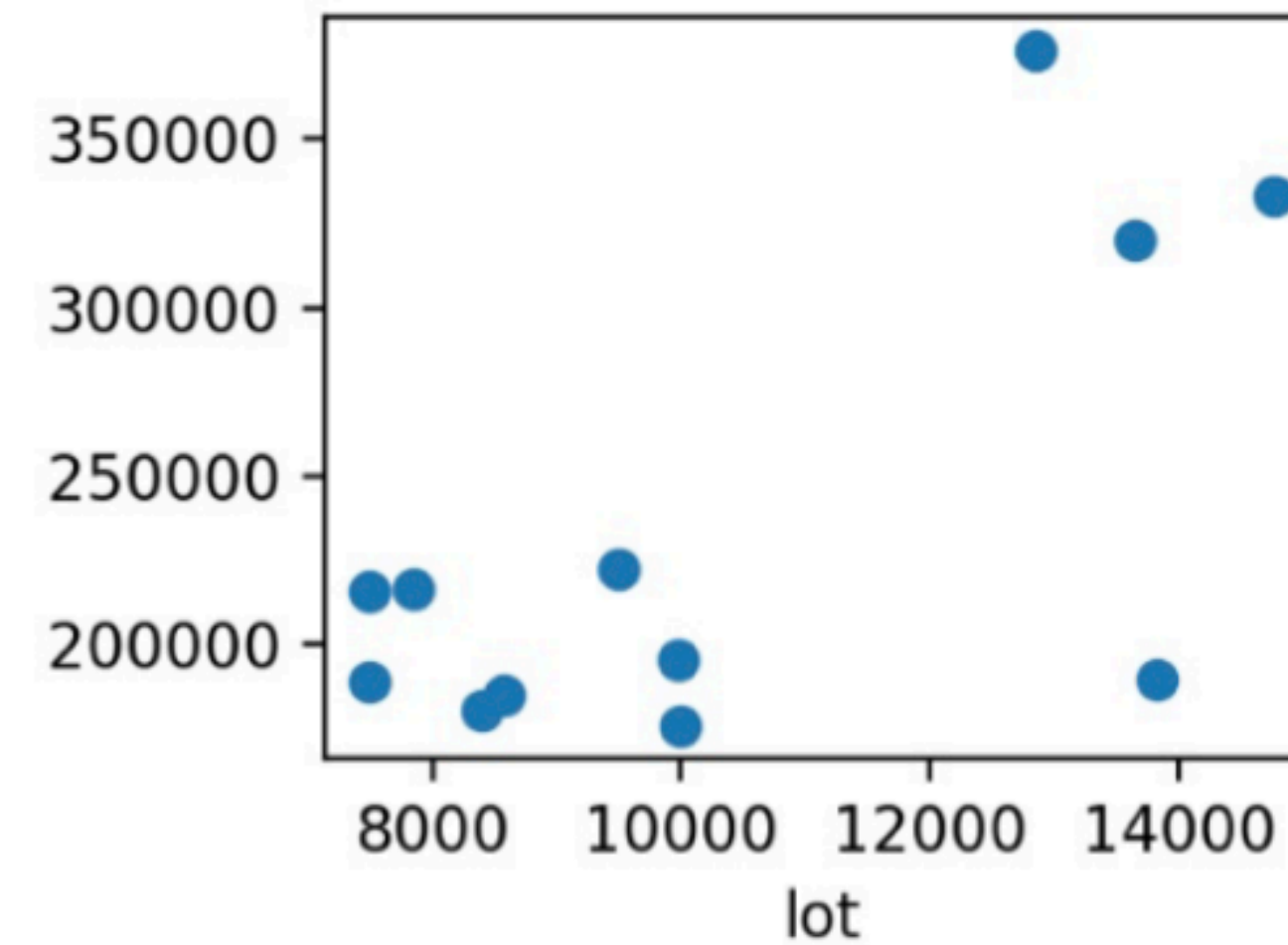
unsupervised $M \leftarrow X$
(labeled validation (X, y))

X

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*

Popular choice

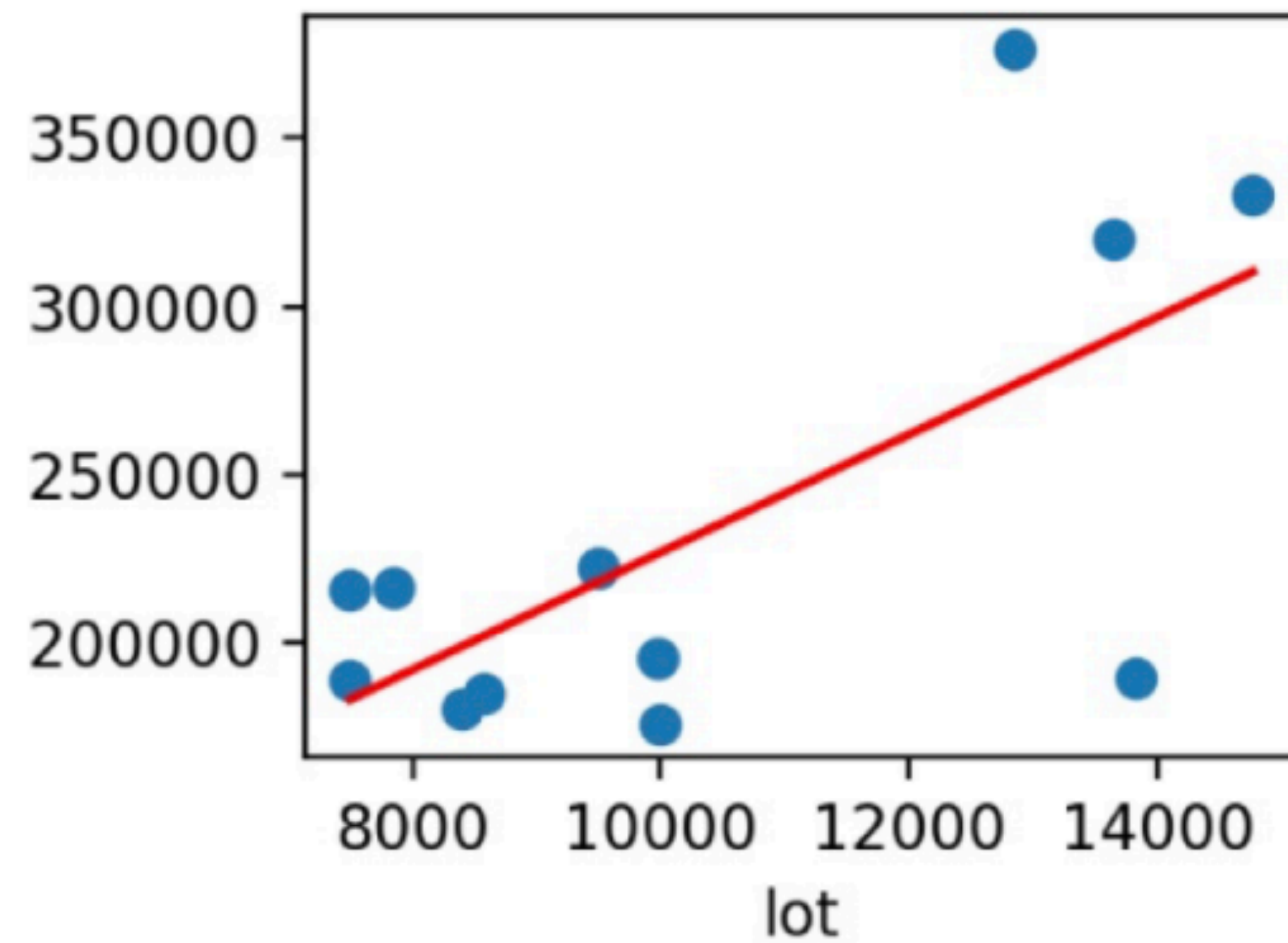
Represent h as a Linear Function

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

Popular choice

The function is defined by **parameters** θ_0 and θ_1 , the function space is greatly reduced

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?

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

$x_0 = 1$

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

3 x

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

4 terms

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)

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

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 θ **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^{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)

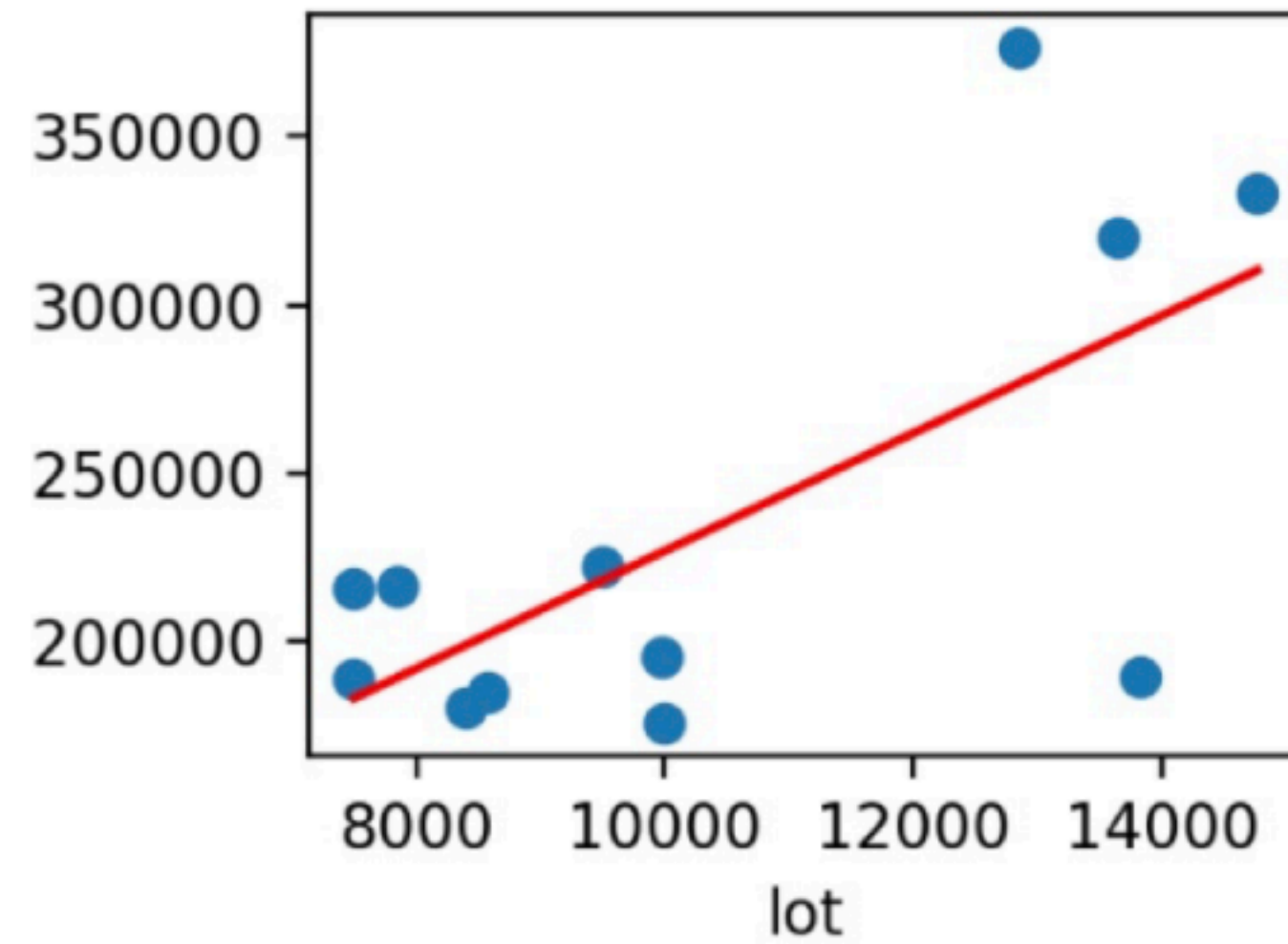
$$\theta = \begin{pmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \end{pmatrix} \text{ and } x^{(1)} = \begin{pmatrix} x_0^{(1)} \\ x_1^{(1)} \\ x_2^{(1)} \\ x_3^{(1)} \end{pmatrix} = \begin{pmatrix} 1 \\ 2104 \\ 4 \\ 45 \end{pmatrix} \text{ and } y^{(1)} = 400$$

We call θ **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^{th} example.

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

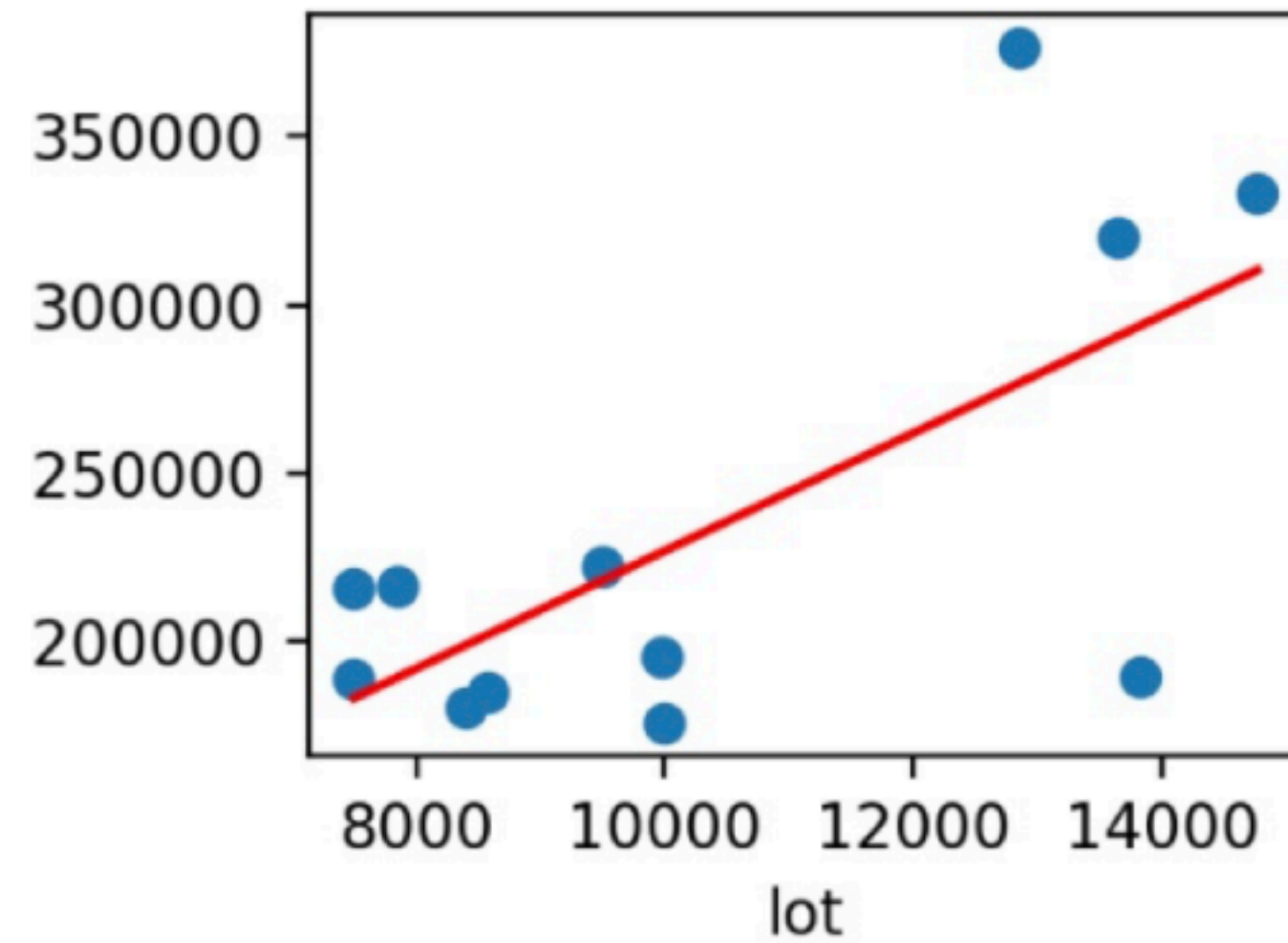
Vector Notation of Prediction



Kernel methods

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

Vector Notation of Prediction



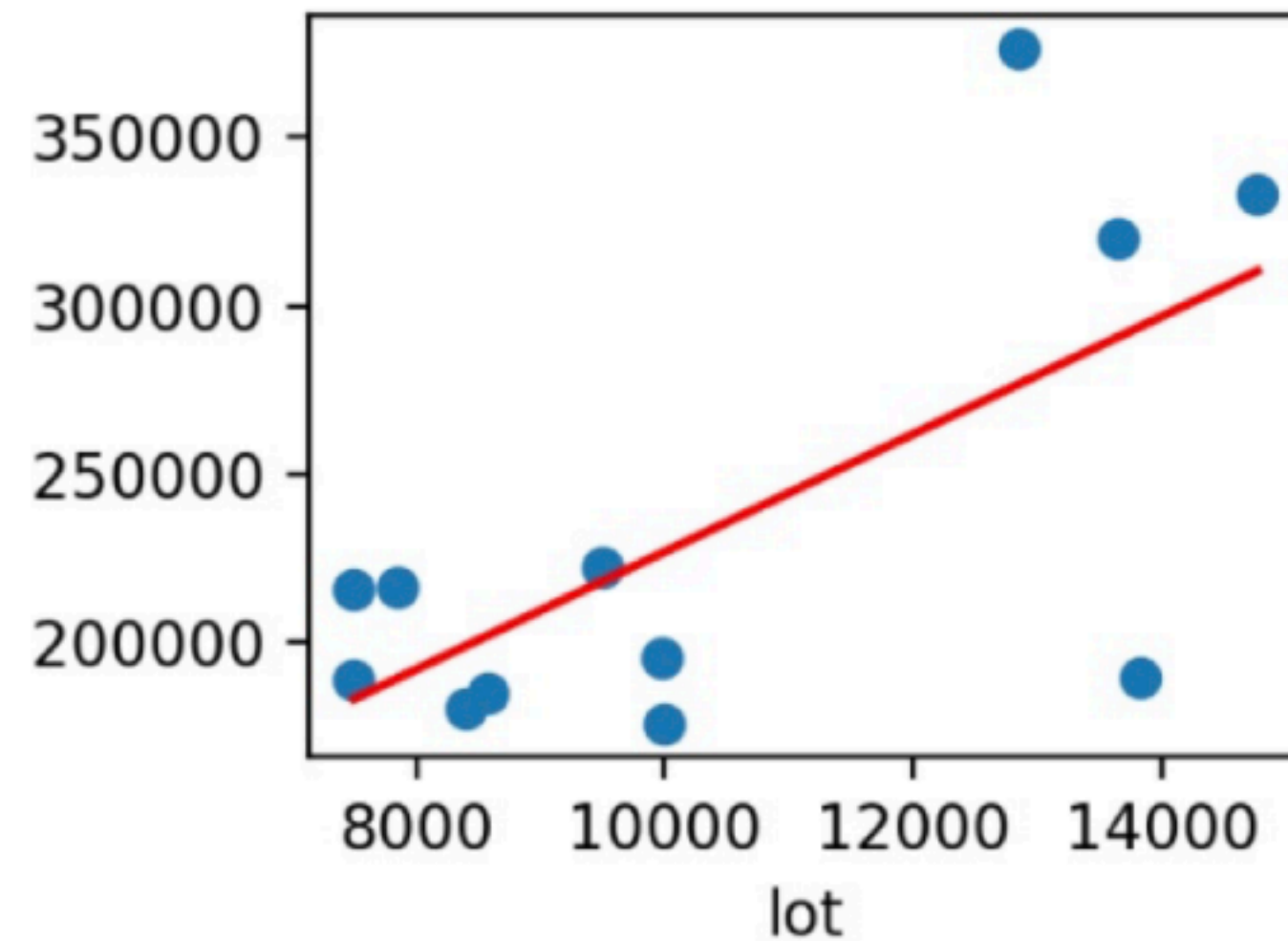
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

how to define metric

similarity

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

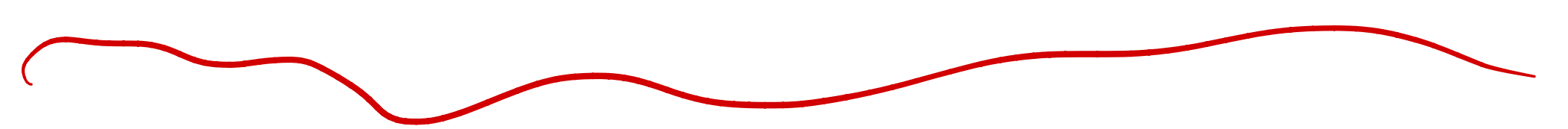
Loss Function



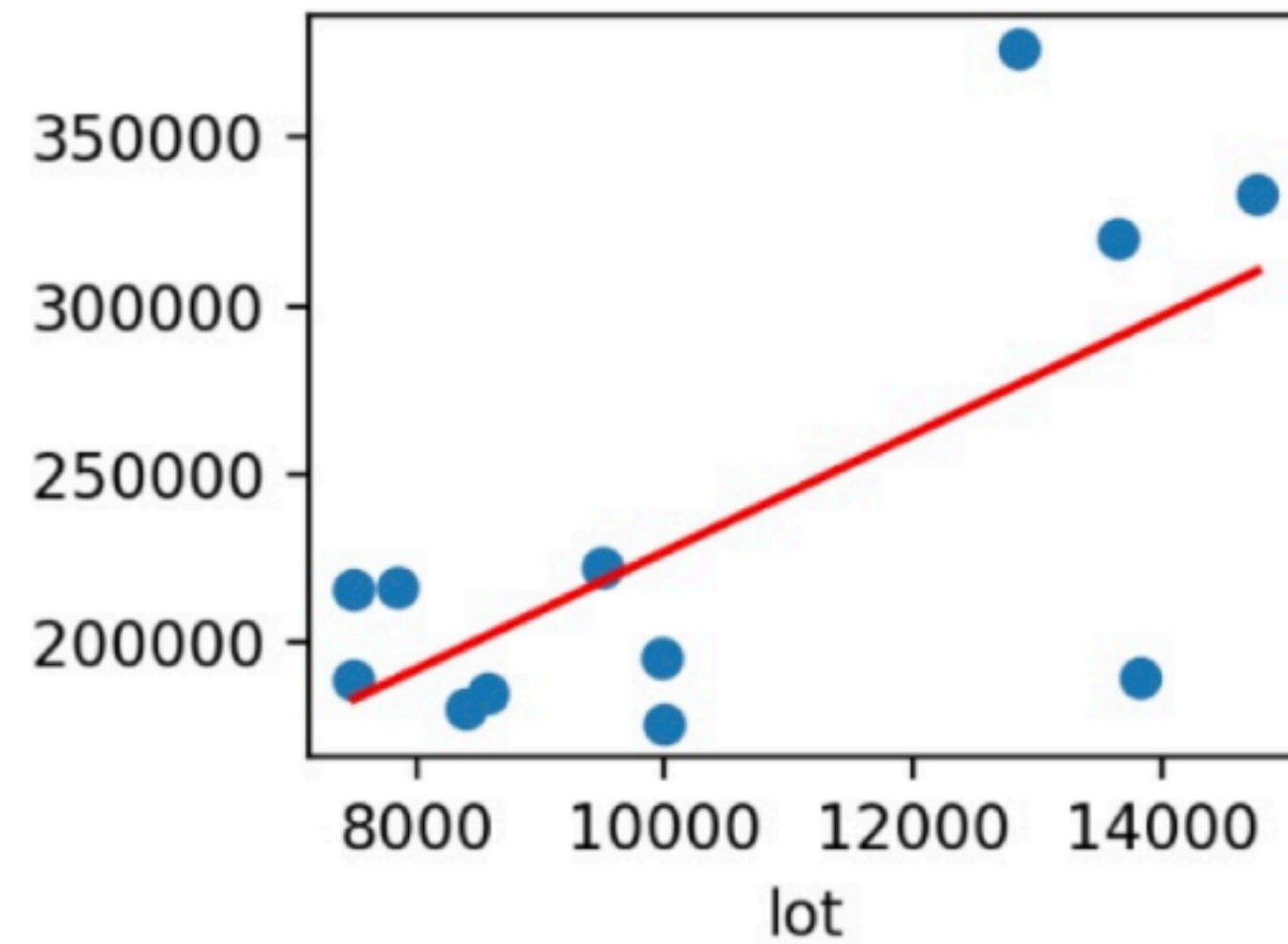
$$h_{\theta}(x) = \sum_{j=0}^d \theta_j x_j = x^T \theta$$

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

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



Least Squares



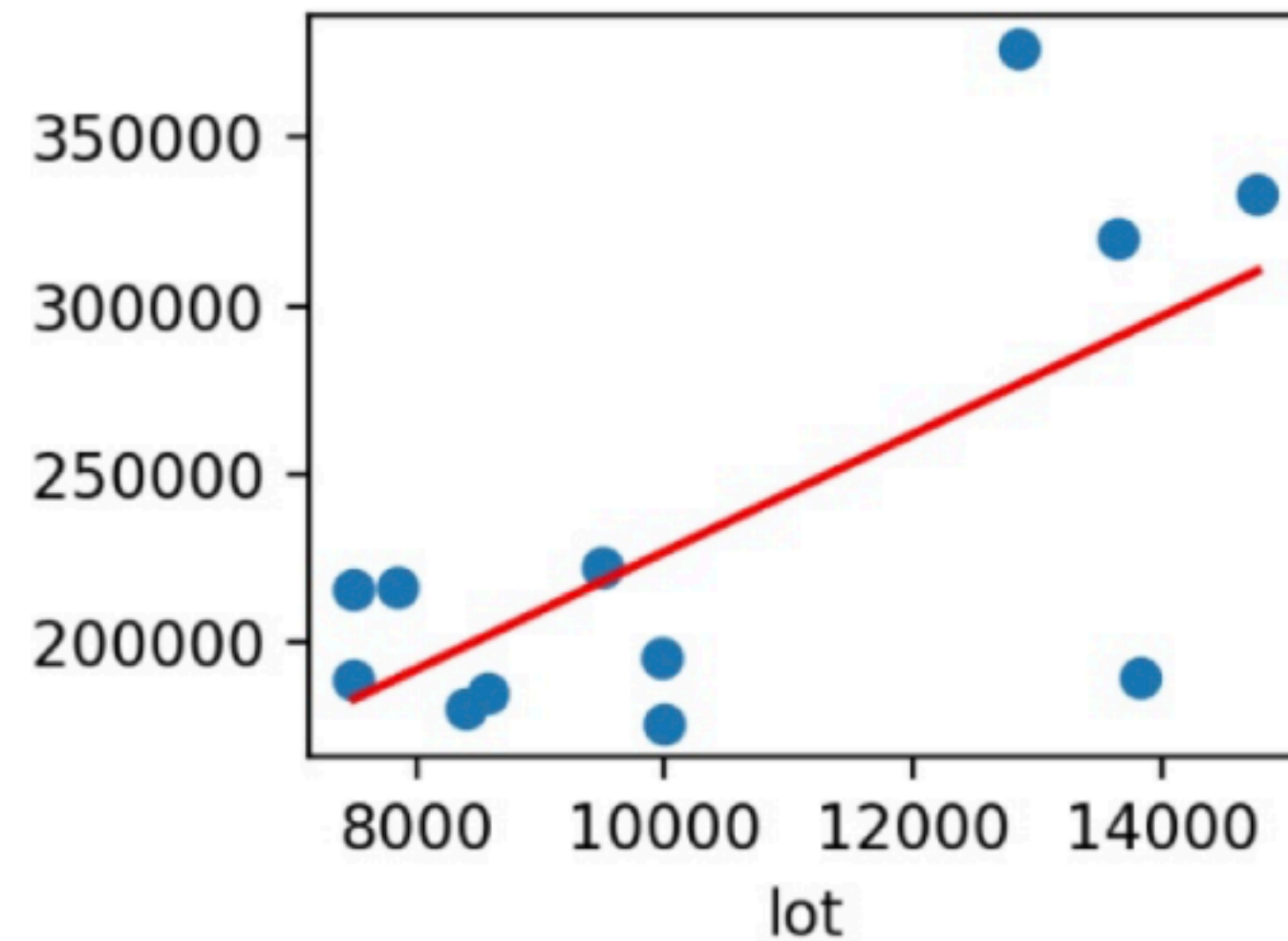
$$\|\vec{h}_\theta(x) - \vec{y}\|$$

$$\|\vec{h}_\theta(x) - \vec{y}\|^2$$

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

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

Least Squares



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

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

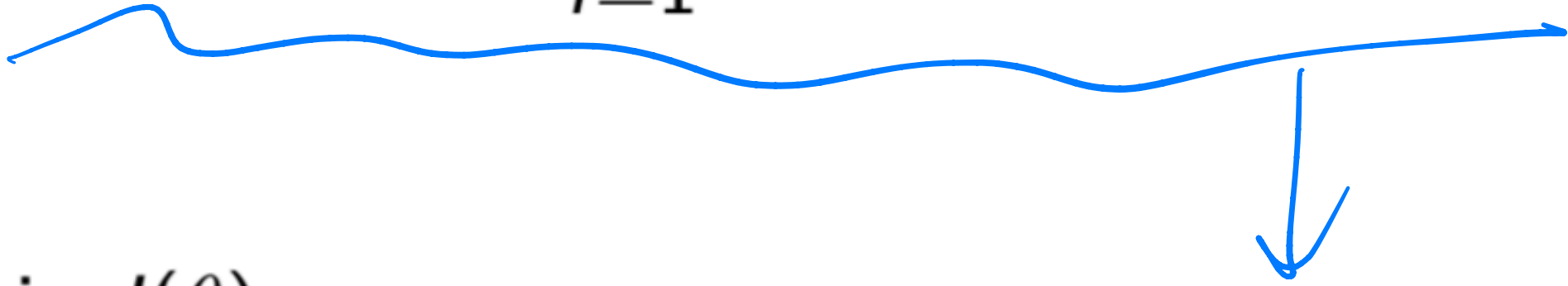
Choose

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


Solving Least Square Problem

Direct Minimization

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

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


Choose

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


Solving Least Square Problem

$$\text{rank}(X) < d$$

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

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

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

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

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

$$= X^T X \theta - X^T \vec{y} = 0$$

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

$$N > d$$

N data samples

d feature size

$X \in \mathbb{R}^{N \times d}$

$X^T X \in \mathbb{R}^{d \times d}$

$\in \mathbb{R}^{d \times d}$

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

$N < d$

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

$\text{rank}(X) < d$

$$\text{rank}(AB) \leq \min(\text{rank}(A), \text{rank}(B))$$

whether $(X^T X)^{-1}$ exists

$$\begin{array}{c}
 n \times \\
 \left[\begin{array}{c}
 V_2 \\
 2 \\
 \left[\begin{array}{c}
 1 \\
 2 \\
 5 \\
 6 \\
 12
 \end{array} \right] \\
 \hline \\
 \hline
 \end{array} \right.
 \end{array}
 \quad
 \begin{array}{c}
 d \\
 5 \\
 \left[\begin{array}{c}
 V_5 \\
 \left[\begin{array}{c}
 2 \\
 4 \\
 10 \\
 12 \\
 24
 \end{array} \right] \\
 \hline \\
 \hline
 \end{array} \right.
 \end{array}
 \end{array}$$

$n < d \rightarrow$ not full rank

$n > d$

$$V_5 = 2 \cdot V_2$$

$$\text{rank}(X) < d$$

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

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

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

Why Least-Square Loss Function?

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

Why Least-Square Loss Function?

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

Assume

maximum likelihood

$$\underline{y^{(i)}} = \theta^T \underline{x^{(i)}} + \underline{\epsilon^{(i)}} \quad \text{noise}$$

Why Least-Square Loss Function?

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

Assume

$$\underbrace{y^{(i)}} = \theta^T x^{(i)} + \underbrace{\epsilon^{(i)}}_{\text{noise}}$$

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

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)} \quad \text{random noise}$$

x, y : random variable

ϵ : 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

ϵ : 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

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

$\epsilon^{(i)}$

$\epsilon^{(i+1)}$

x, y : random variable

ϵ : 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?

Page

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

Σ

Why Least-Square Loss Function?

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

$\xi \sim N(0, 1)$

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

$x \sim N(\mu, \sigma^2)$

$\xi = y - \theta^T x$

$ax + b \sim N(ax + b, a^2)$

Why Least-Square Loss Function?

$$P(A_1, A_2, \dots, A_n) = P(A_1) P(A_2 | A_1) P(A_3 | A_1, A_2) \dots$$

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

one data example

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

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

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

Why Least-Square Loss Function?

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

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

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

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

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?

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

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 θ ?



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 θ ?

Maximize the probability of Y's happening!

$$\arg \max_{\theta} L(\theta) = \arg \max_{\theta} \log L(\theta)$$

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

$$\operatorname{argmax}_{\theta} \ell(\theta) = \operatorname{argmax}_{\theta} - \frac{1}{\sigma^2} \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^T x^{(i)})^2$$

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 θ ?

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 θ ?

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 θ ?

Maximize the probability of Y's happening?

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

Ground-truth θ^*

θ^*



$$E_{x_{np} \text{ data}}[\hat{\theta}] = \theta^*$$

unbiased estimator

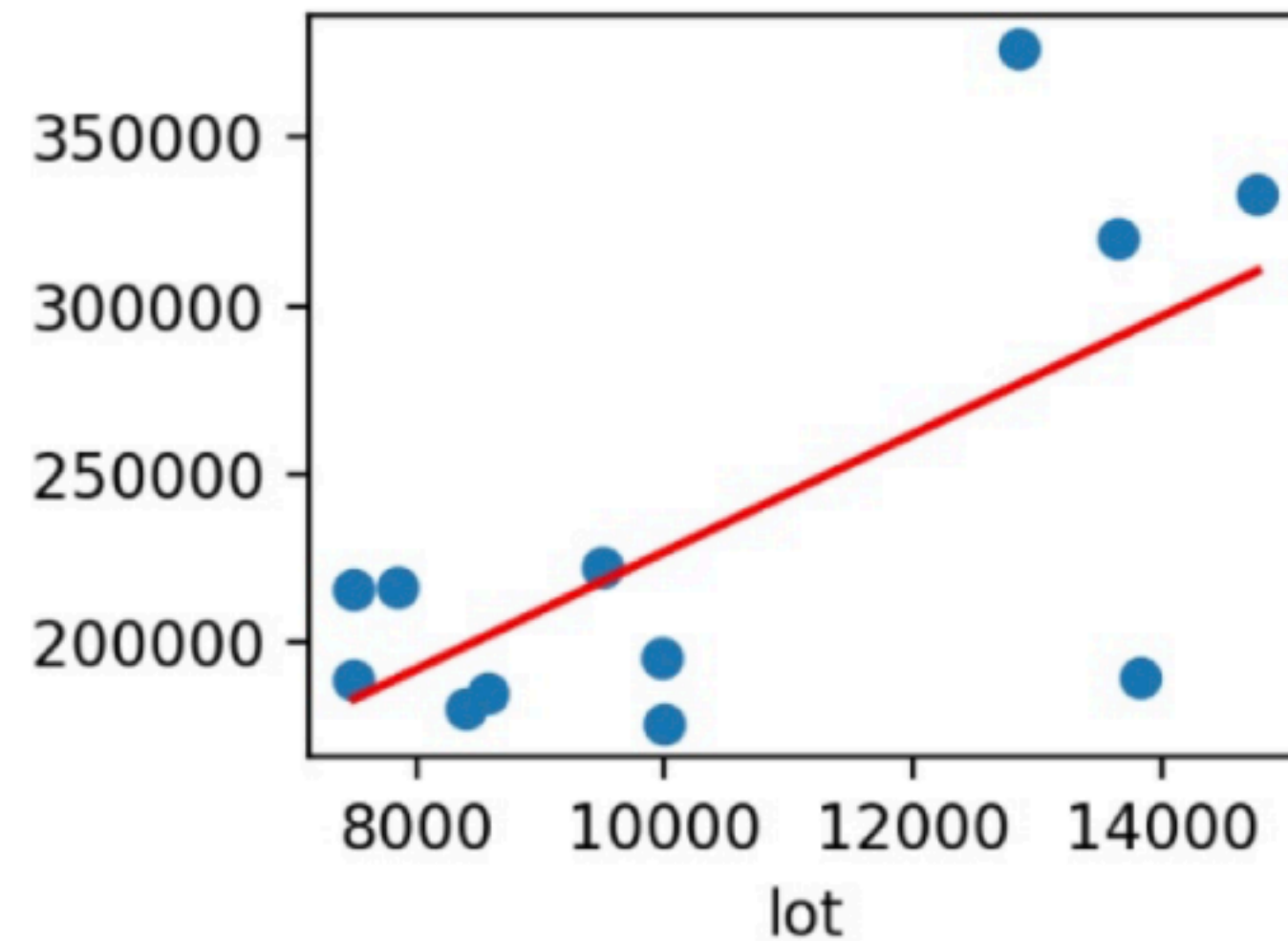
biased estimator

$x_1, \dots, x_n \sim D$ mean μ , σ^2

$$E\left[\hat{\mu} = \frac{x_1 + \dots + x_n}{n}\right] = \mu \quad \text{Var}(\mu) = E[(x - \mu)^2]$$

$$E\left[\hat{\sigma}^2 = \frac{\sum_i (x_i - \hat{\mu})^2}{n}\right] = \frac{\sum_i (x_i - \hat{\mu})^2}{n-1} \rightarrow \sigma^2$$

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

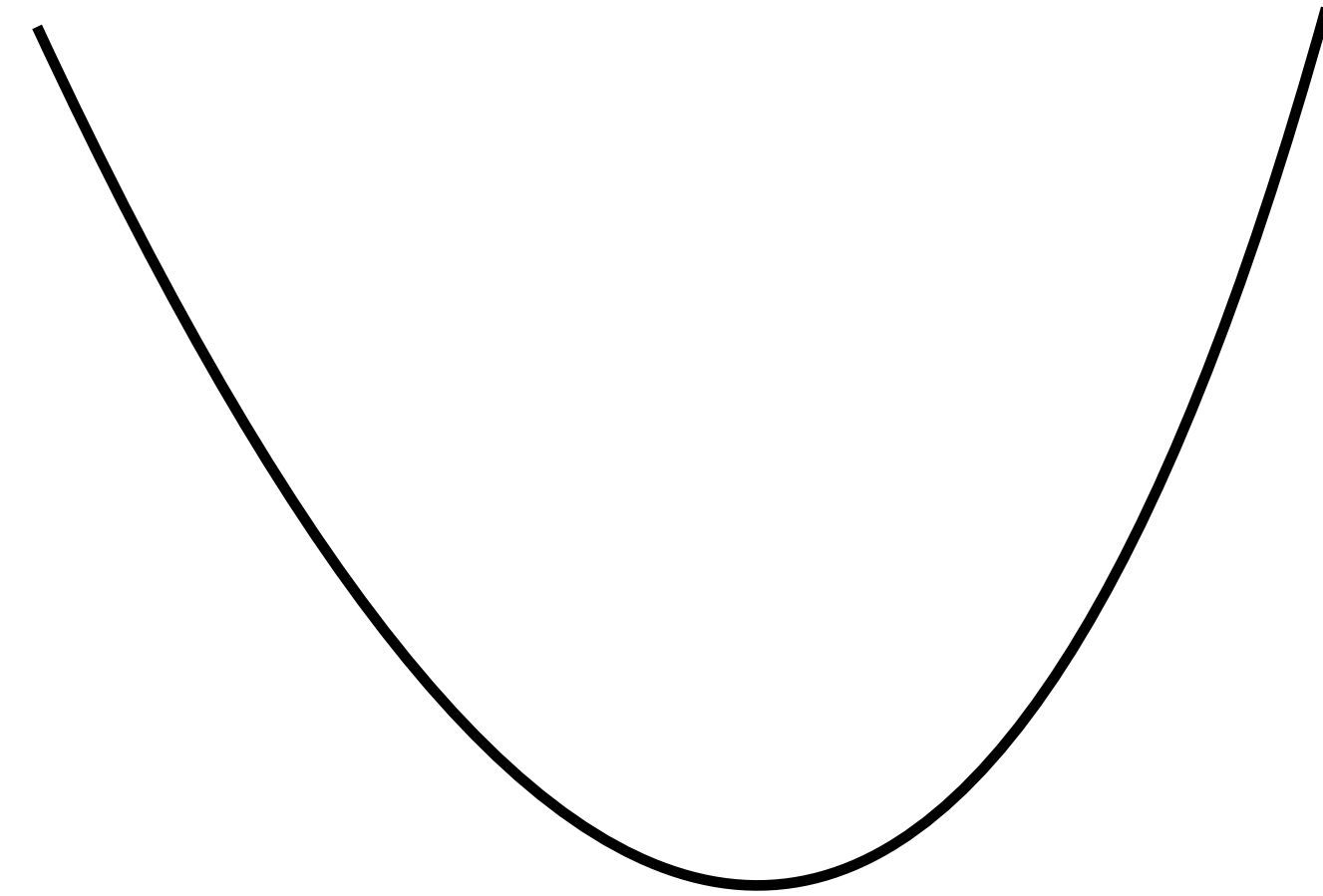
$$J(\theta) = \frac{1}{2} \sum_{i=1}^n \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

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

Gradient Descent

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

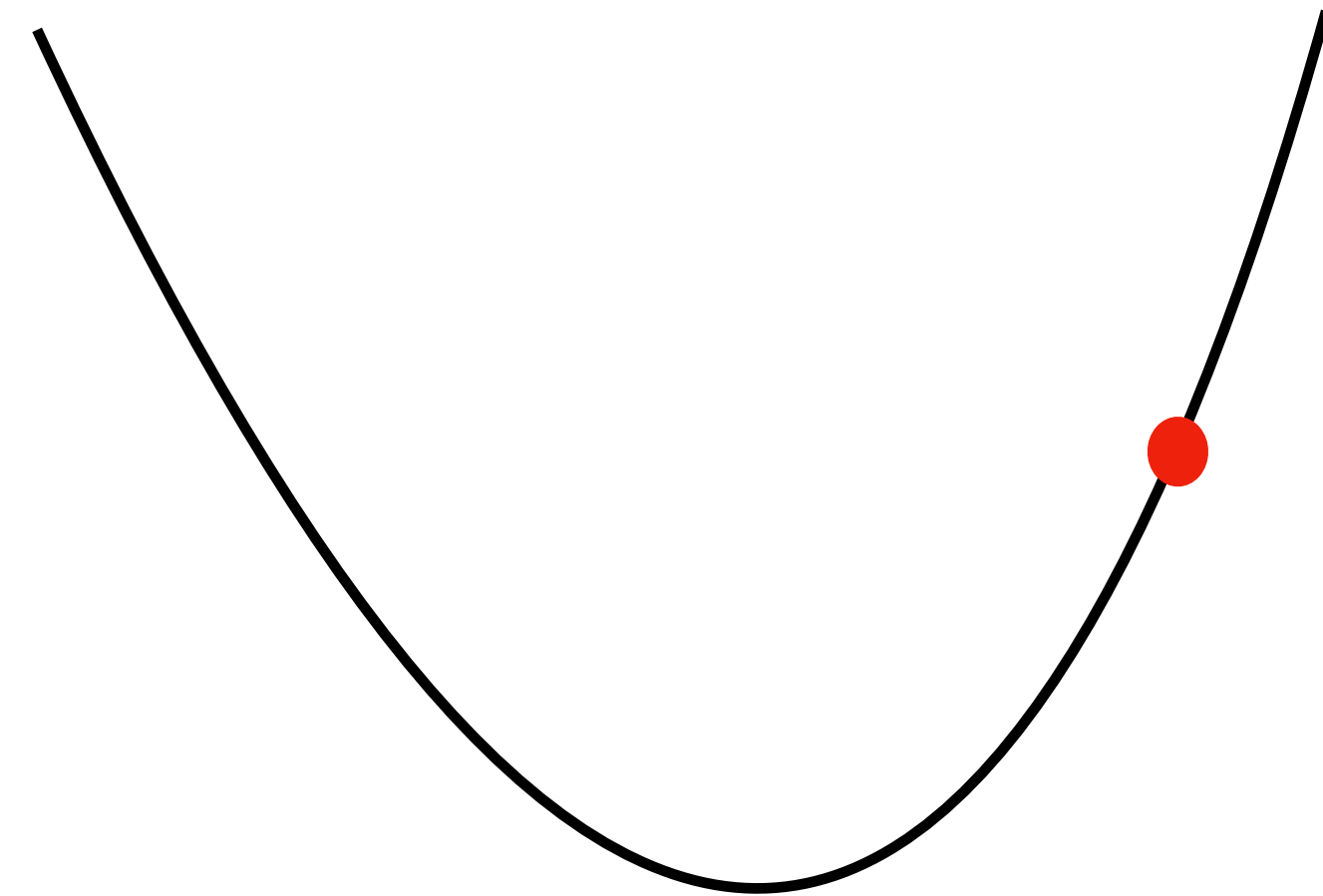
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$



Gradient Descent

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

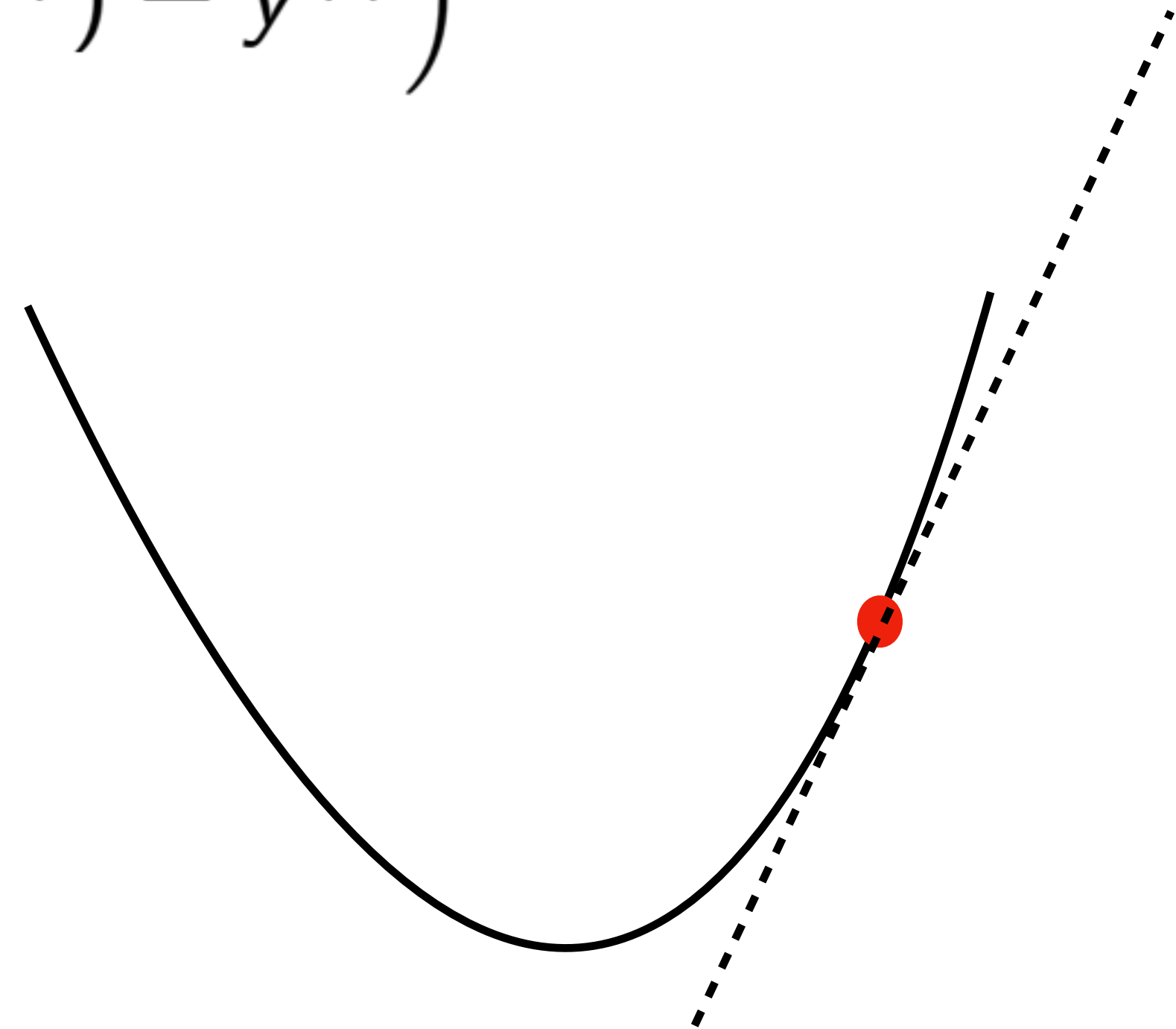
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$



Gradient Descent

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

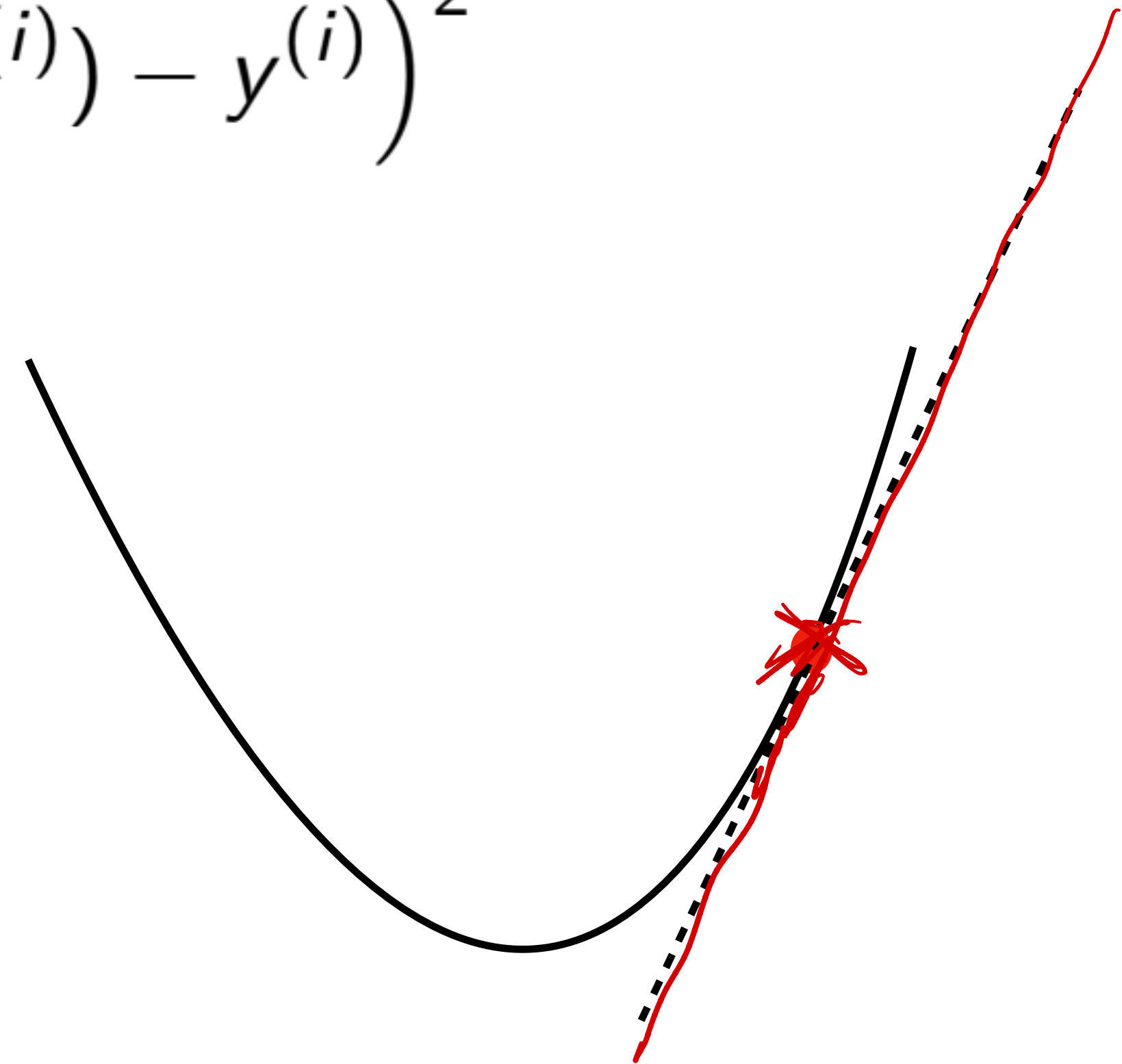
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$



Gradient Descent

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

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



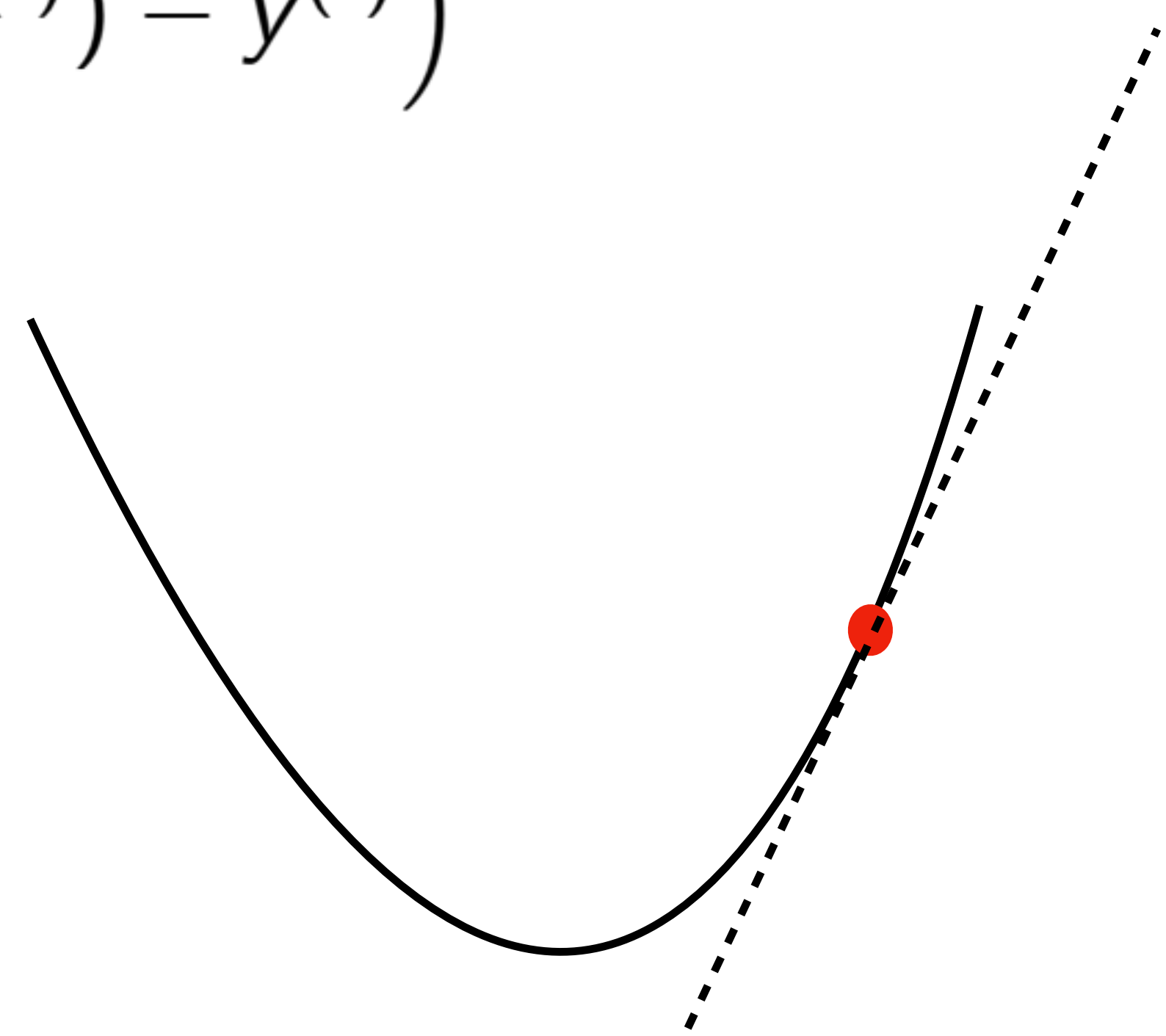
The direction of the steepest decrease of J

Gradient Descent

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

Learning Rate

$$\theta_j := \theta_j - \underbrace{\alpha}_{\text{Learning Rate}} \frac{\partial}{\partial \theta_j} J(\theta)$$



The direction of the steepest decrease of J

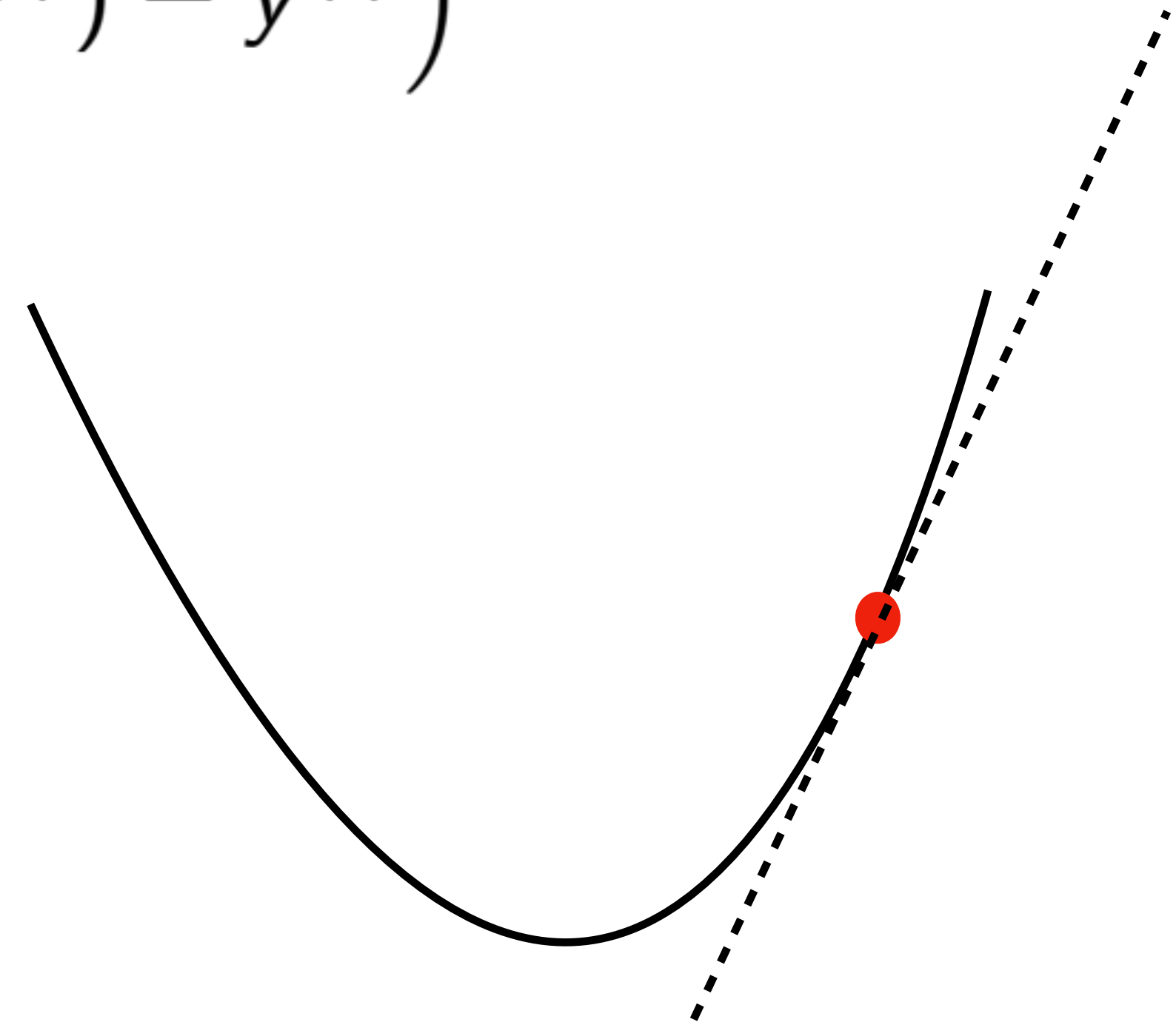
Gradient Descent

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

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:

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

Gradient Descent

For a single training example:

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

LMS (Least Mean Square) Update Rule

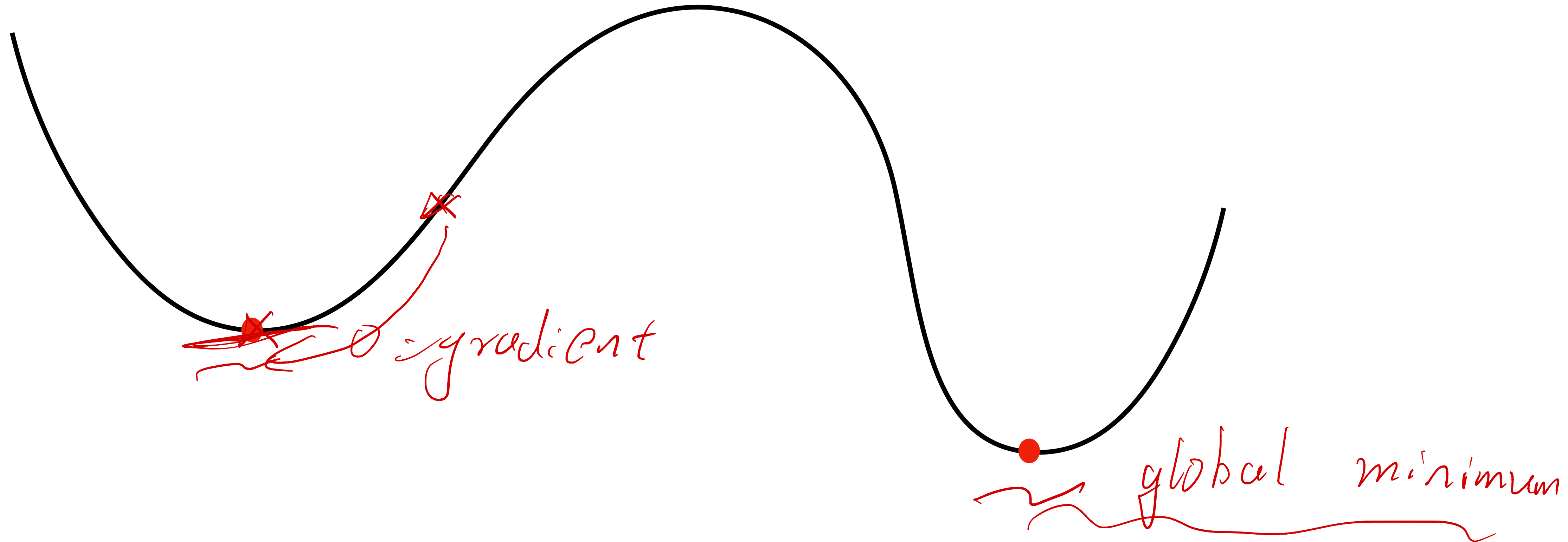
Batch Gradient Descent

For a multiple training examples:

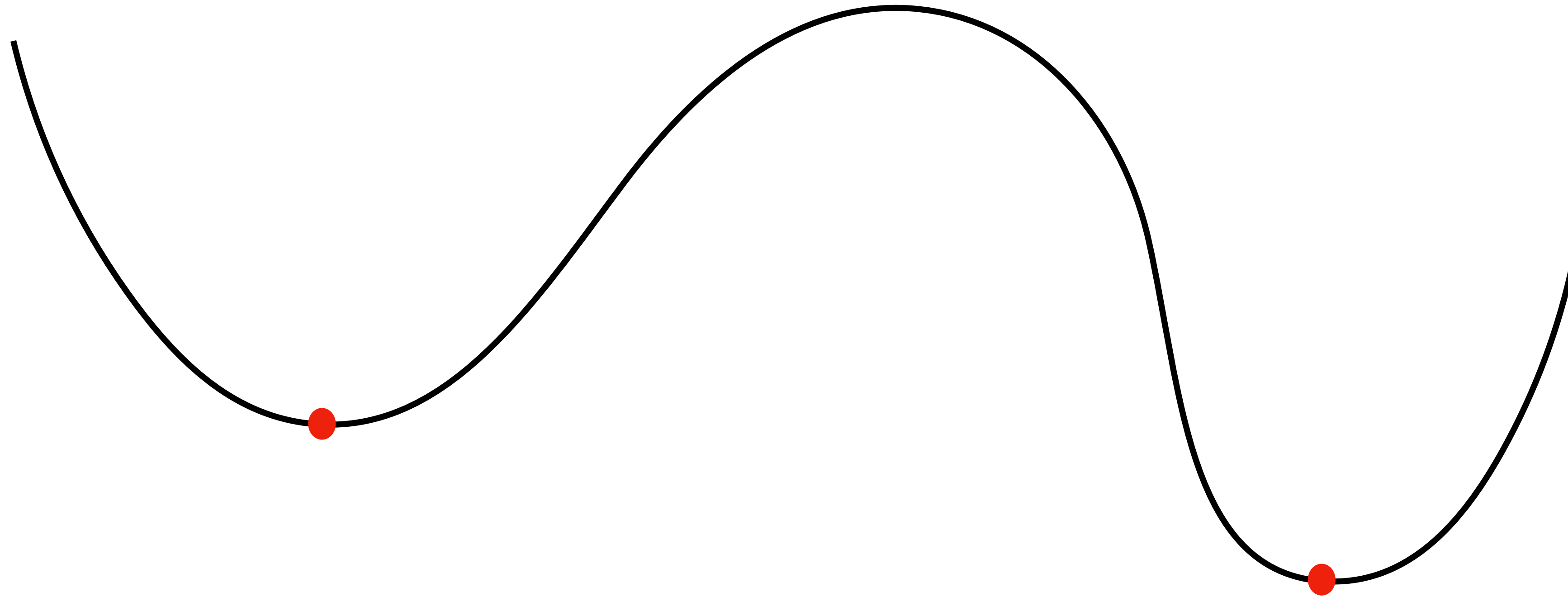
$$\theta_j := \theta_j + \alpha \sum_{i=1}^n (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

Repeat until convergence

Local Minimum



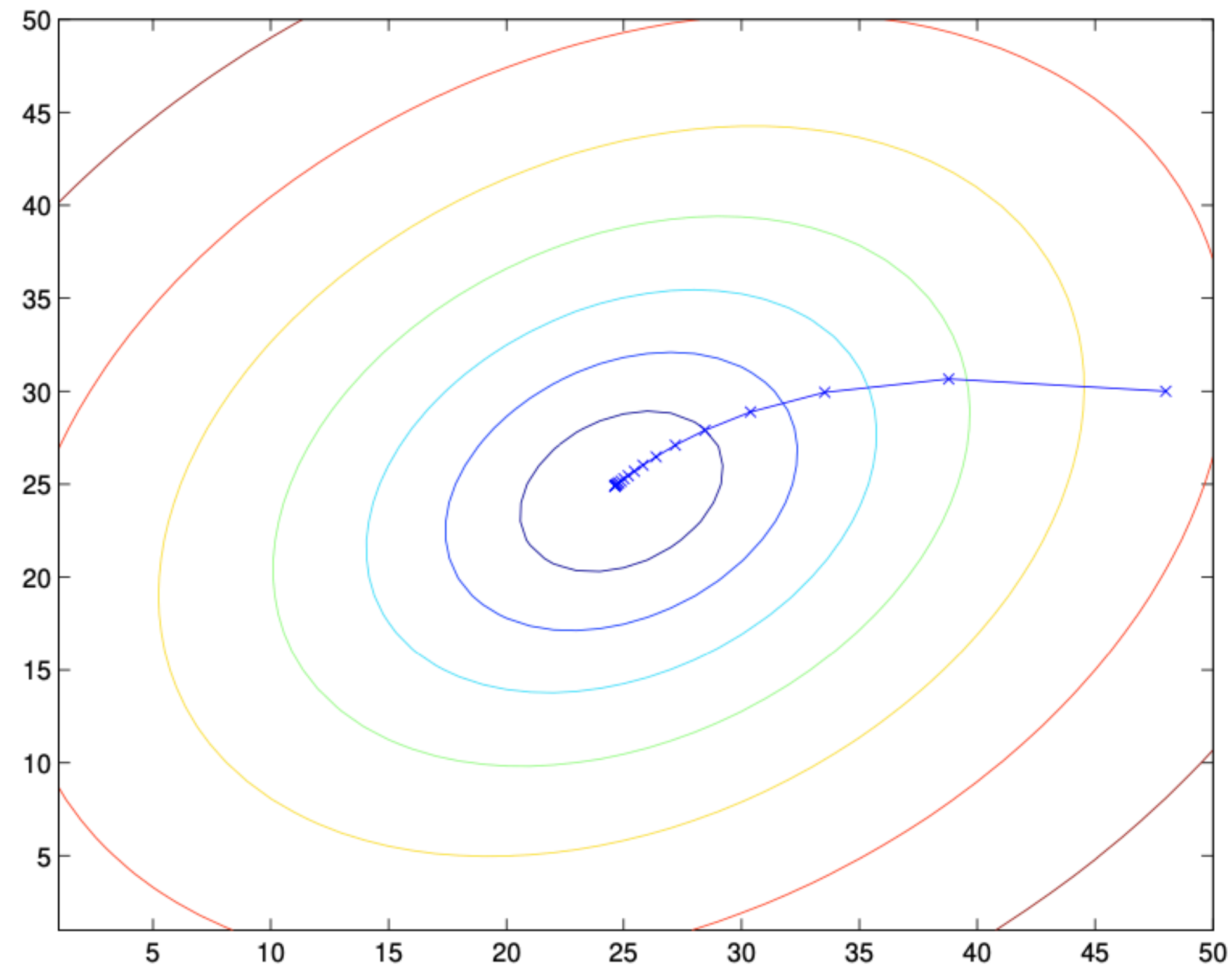
Local Minimum



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

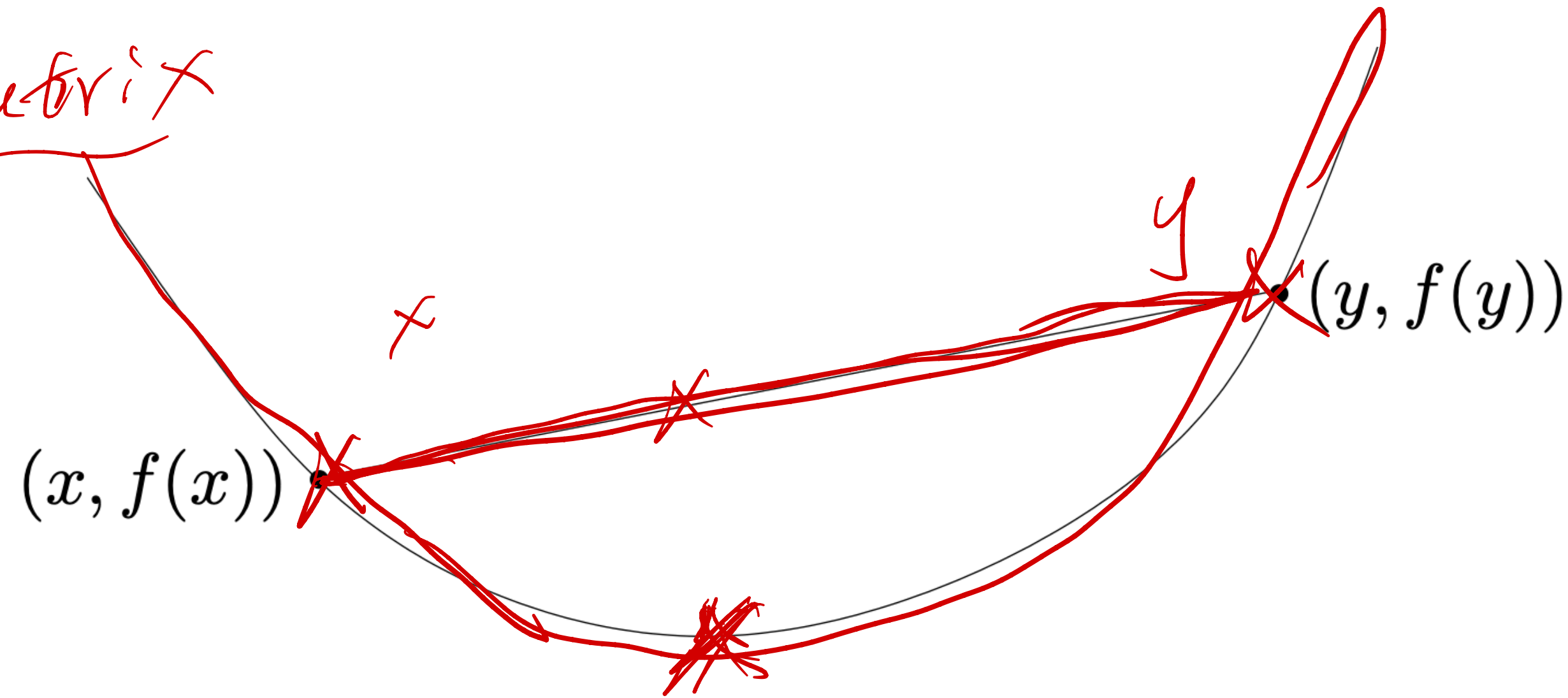


Evidence lower bound (ELBO) Convex Function

Jensen inequality

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

Hessian matrix



$$f(t_1 x_1 + t_2 x_2 + \dots + t_n x_n) \leq t_1 f(x_1) + t_2 f(x_2) + \dots + t_n f(x_n)$$

Thank You!
Q & A

Thank You!
Q & A