

Linear Regression

Prepared by: Joseph Bakarji

Data

Given a table of numbers, what can you do?

Living area (feet ²)	#bedrooms	Price (1000\$s)		x_1	x_2	x_3	x_4
2104	3	400		Apt. 1	$x_1^{(1)}$	$x_2^{(1)}$	$x_3^{(1)}$	$x_4^{(1)}$	
1600	3	330							
2400	3	369							
1416	2	232	→	Apt. 2	$x_1^{(2)}$	$x_2^{(2)}$	$x_3^{(2)}$	$x_4^{(2)}$	
3000	4	540							
⋮	⋮	⋮		Apt. 3	$x_1^{(3)}$	$x_2^{(3)}$	$x_3^{(3)}$	$x_4^{(3)}$	⋮
				Apt. 4	$x_1^{(4)}$	$x_2^{(4)}$	$x_3^{(4)}$	$x_4^{(4)}$	
					⋮	⋮	⋮	⋮	

- Visualization: Look at it!
- Find statistical features: Mean, median, outliers etc.
- Clean it: missing values, ...

What are the **inputs** and **outputs**?

- **Inputs:** quantities that are typically **given**
- **Outputs:** quantities we want to **predict**

Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400
1600	3	330
2400	3	369
1416	2	232
3000	4	540
⋮	⋮	⋮



Apt. 1

Apt. 2

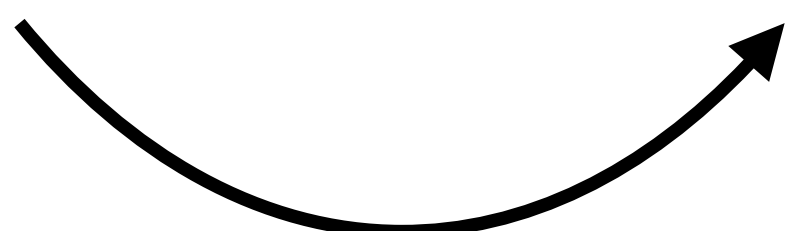
Apt. 3

Apt. 4

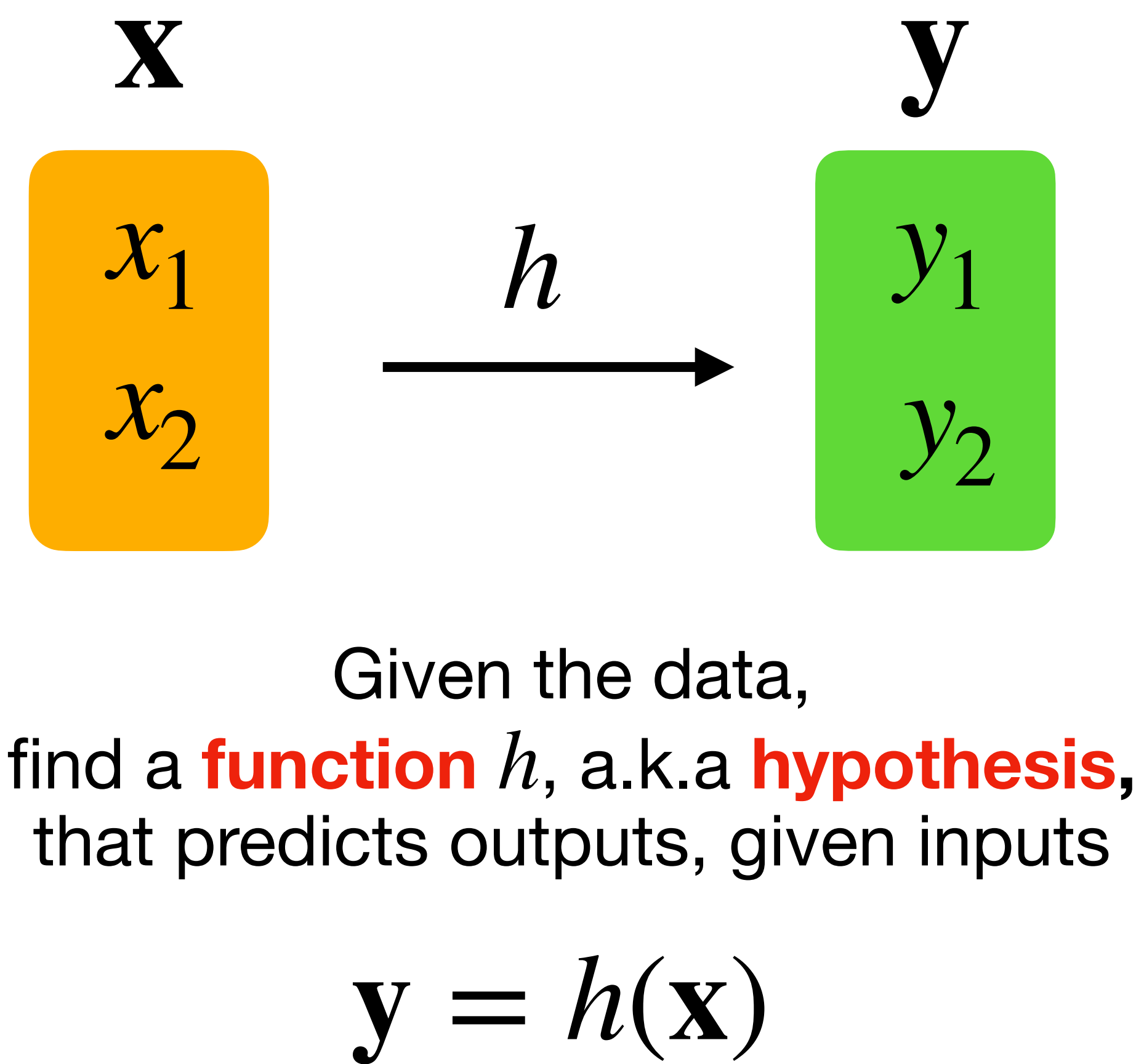
Inputs		Outputs	
x_1	x_2	x_3	x_4	
$x_1^{(1)}$	$x_2^{(1)}$	$x_3^{(1)}$	$x_4^{(1)}$	
$x_1^{(2)}$	$x_2^{(2)}$	$x_3^{(2)}$	$x_4^{(2)}$	
$x_1^{(3)}$	$x_2^{(3)}$	$x_3^{(3)}$	$x_4^{(3)}$	
$x_1^{(4)}$	$x_2^{(4)}$	$x_3^{(4)}$	$x_4^{(4)}$	
⋮	⋮	⋮	⋮	⋮

Given inputs, predict outputs

	Inputs		Outputs	
	x_1	x_2	y_1	y_2
Apt. 1	$x_1^{(1)}$	$x_2^{(1)}$	$y_1^{(1)}$	$y_2^{(1)}$
Apt. 2	$x_1^{(2)}$	$x_2^{(2)}$	$y_1^{(2)}$	$y_2^{(2)}$
Apt. 3	$x_1^{(3)}$	$x_2^{(3)}$	$y_1^{(3)}$	$y_2^{(3)}$
Apt. 4	$x_1^{(4)}$	$x_2^{(4)}$	$y_1^{(4)}$	$y_2^{(4)}$
	\vdots	\vdots	\vdots	\vdots
Given Input	$x_1^{(n)}$ $x_2^{(n)}$??	Unknown output



Supervised Learning



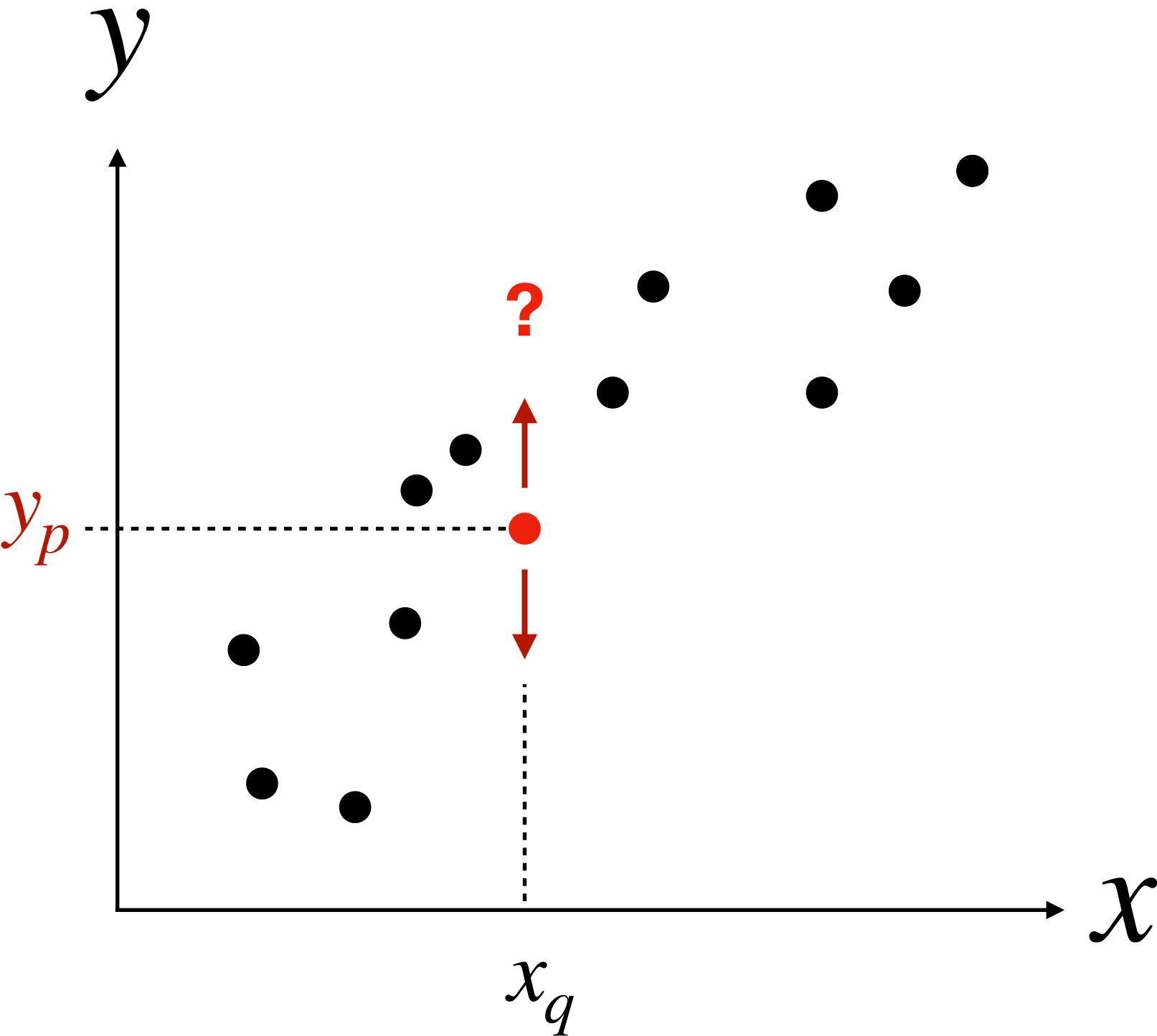
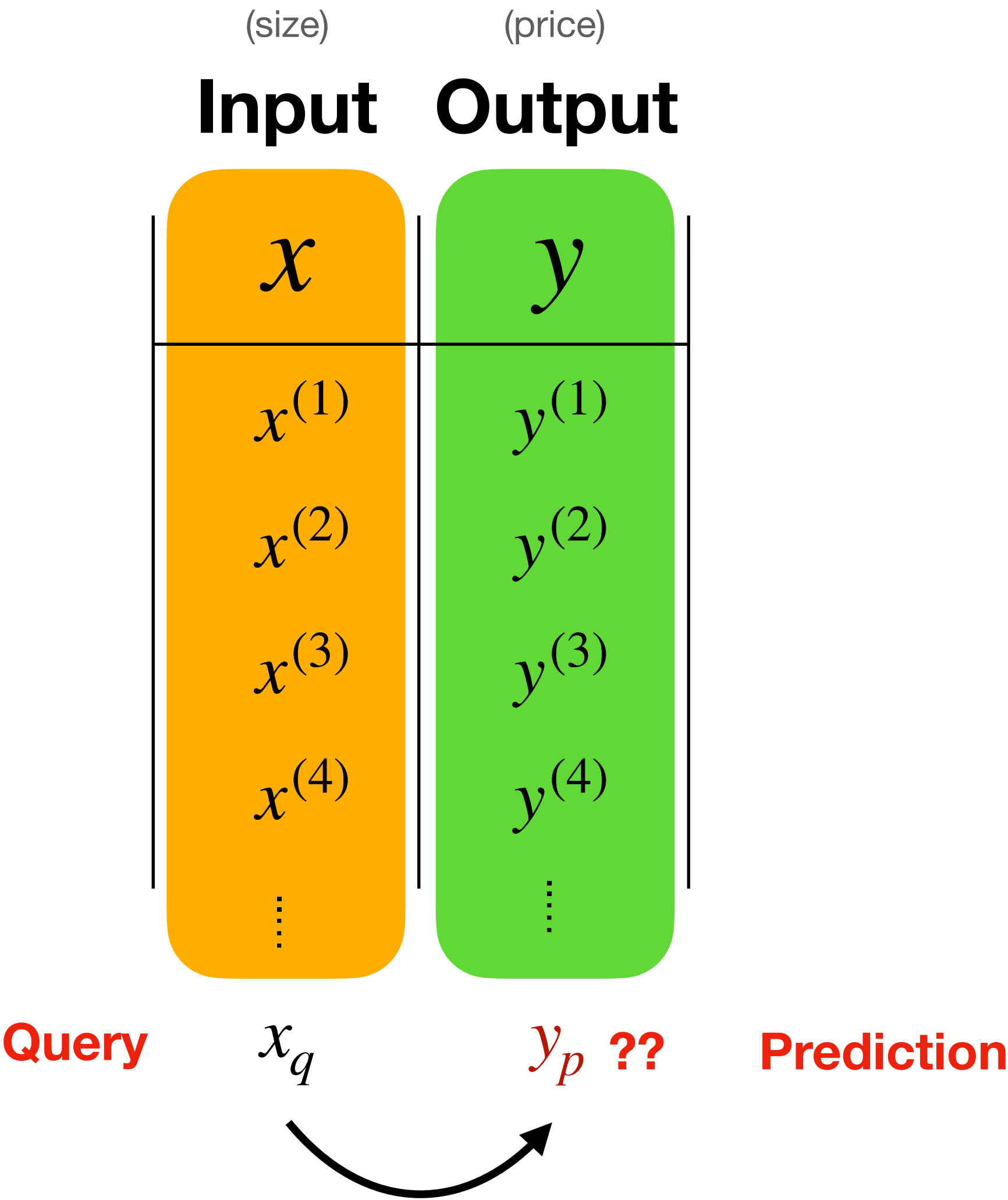
Assume multiple inputs, 1 output

x_1	x_2	y	
Living area (feet ²)	#bedrooms	Price (1000\$s)
2104	3	400	
1600	3	330	
2400	3	369	
1416	2	232	
3000	4	540	
\vdots	\vdots	\vdots	

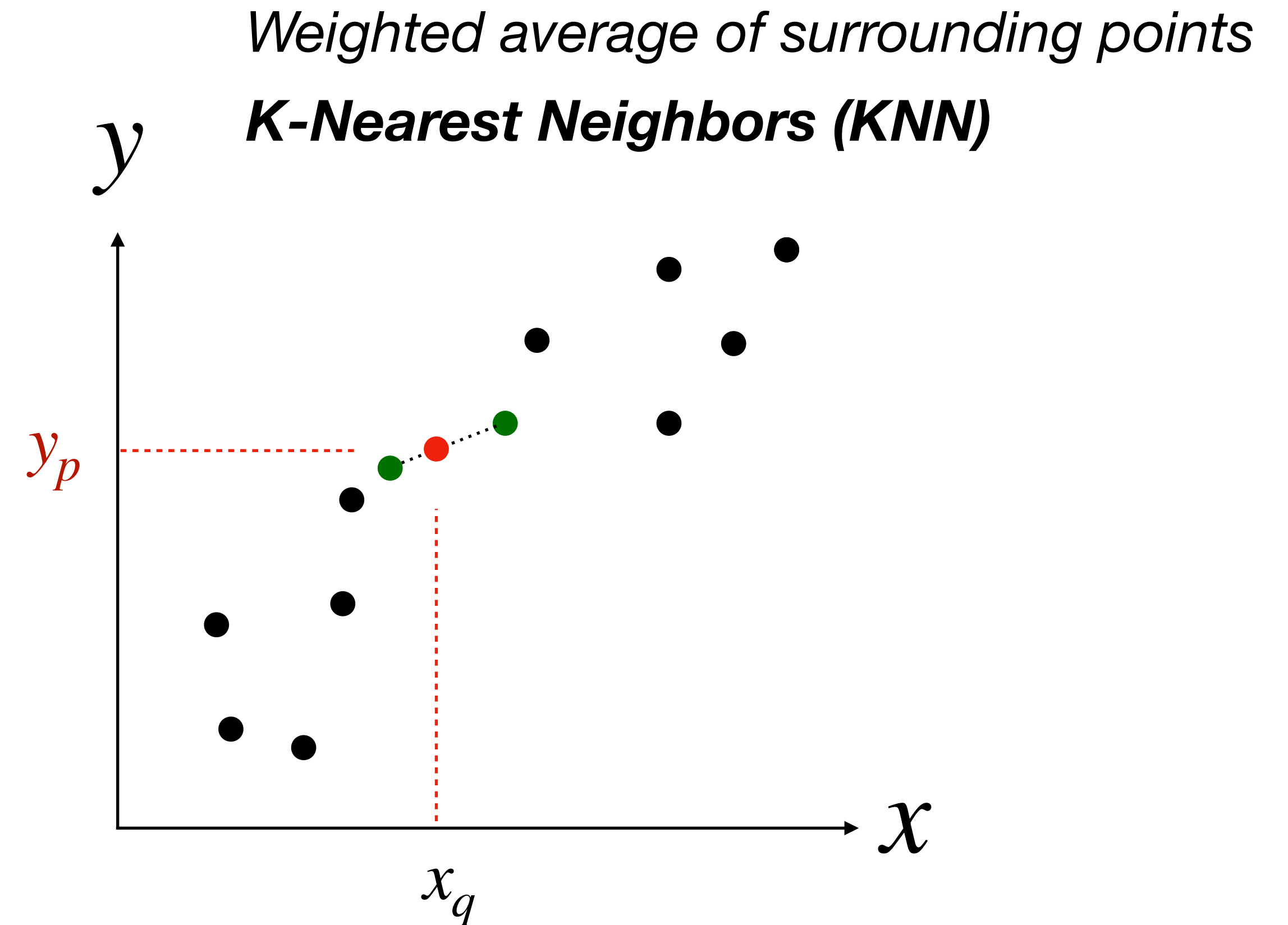
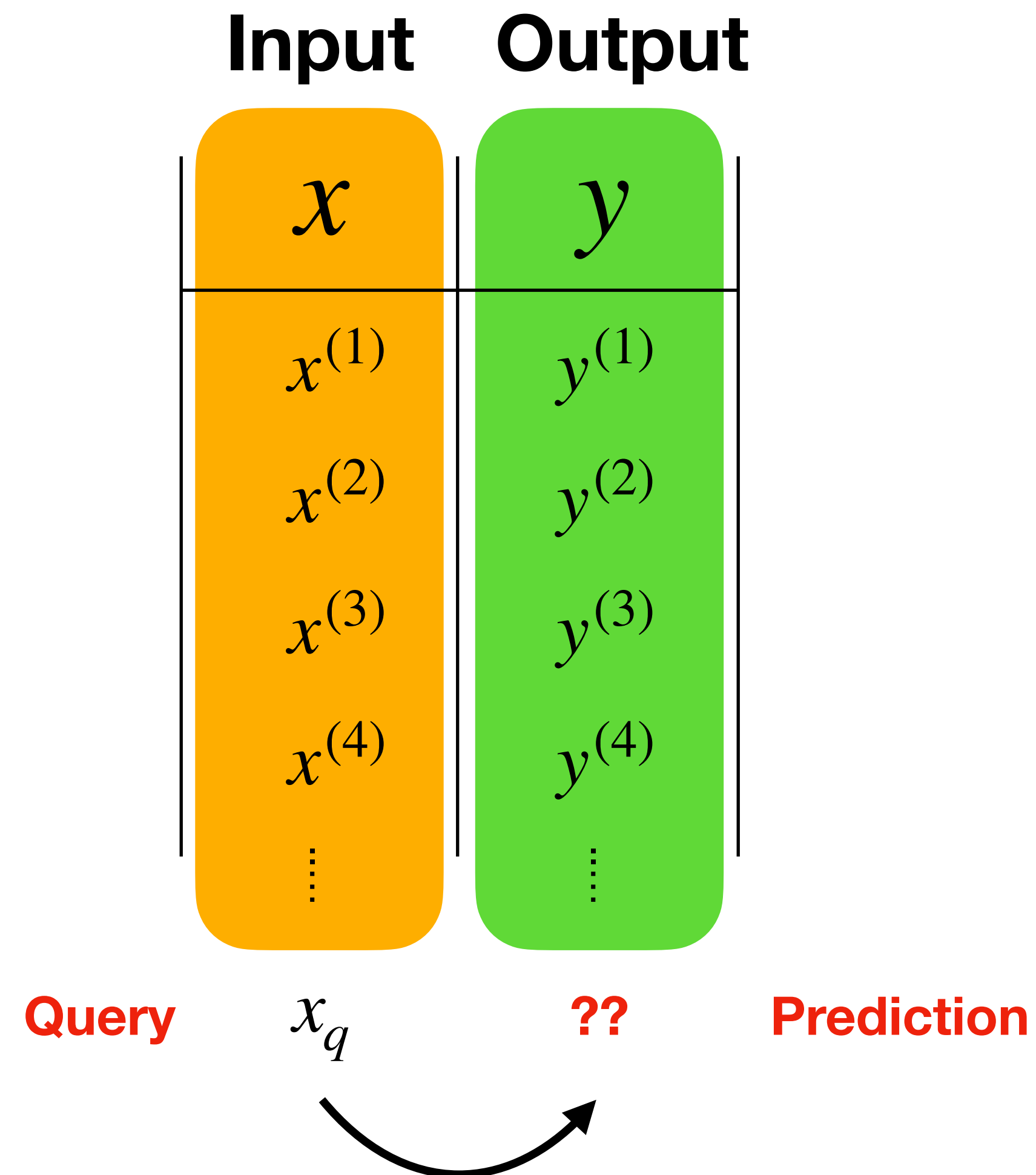


Inputs		Output
x_1	x_2	y
$x_1^{(1)}$	$x_2^{(1)}$	$y^{(1)}$
$x_1^{(2)}$	$x_2^{(2)}$	$y^{(2)}$
$x_1^{(3)}$	$x_2^{(3)}$	$y^{(3)}$
$x_1^{(4)}$	$x_2^{(4)}$	$y^{(4)}$
\vdots	\vdots	\vdots

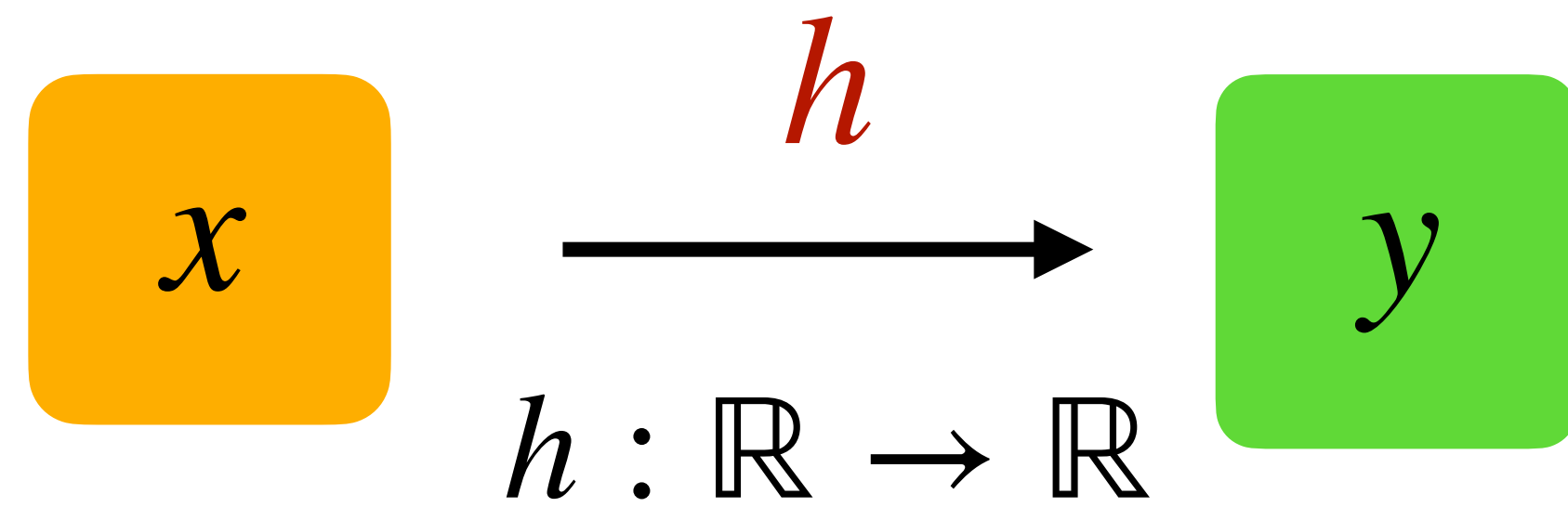
Given new input, what's the output?



Given new input, what's the output?

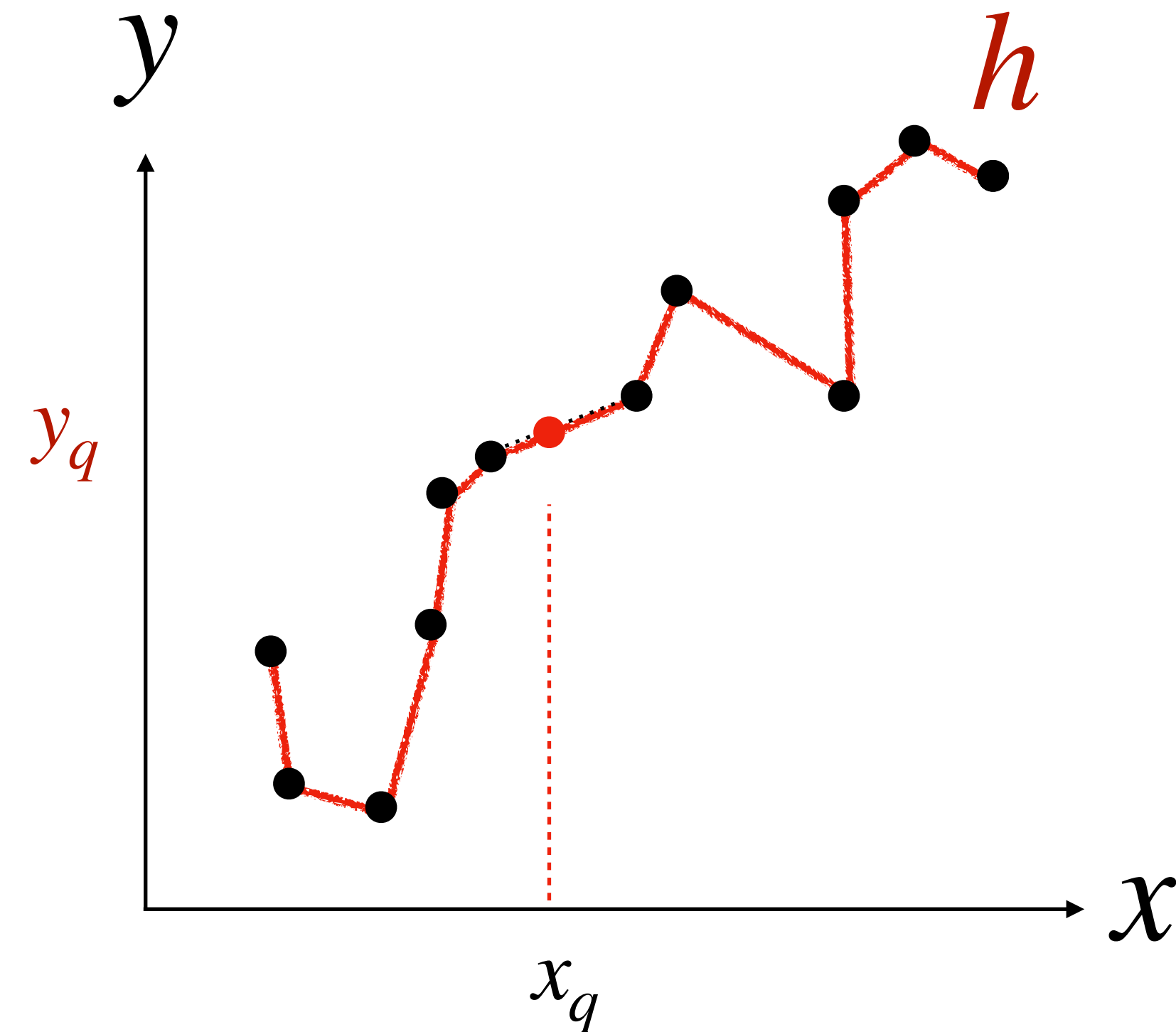


Given new input, what's the output?

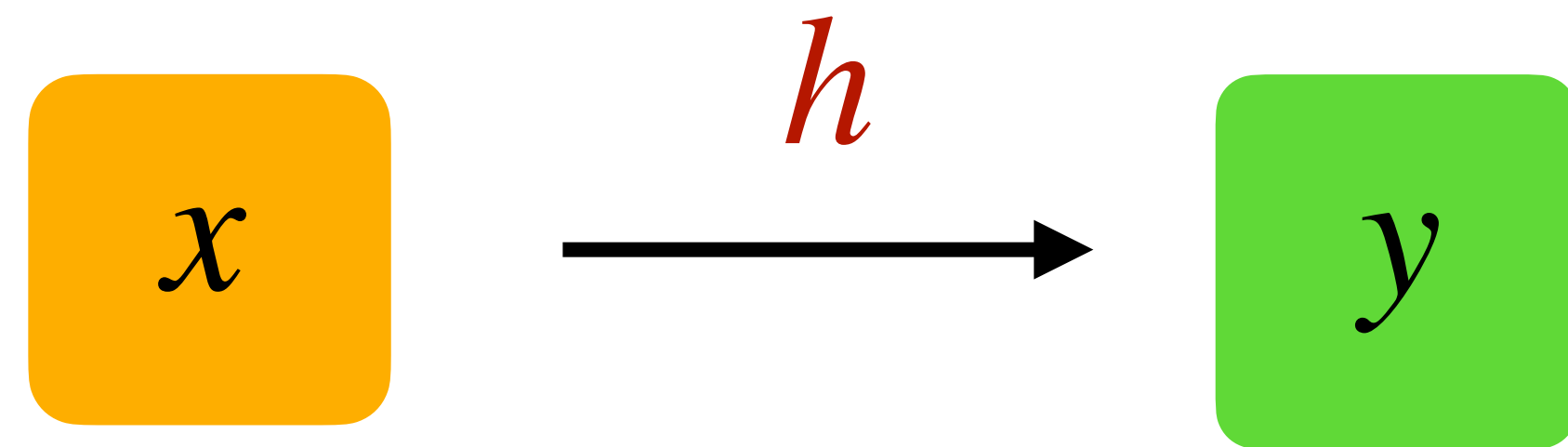


Given the data,
find a **function** h , a.k.a **hypothesis**,
that predicts an output, given an input

Linear Interpolation

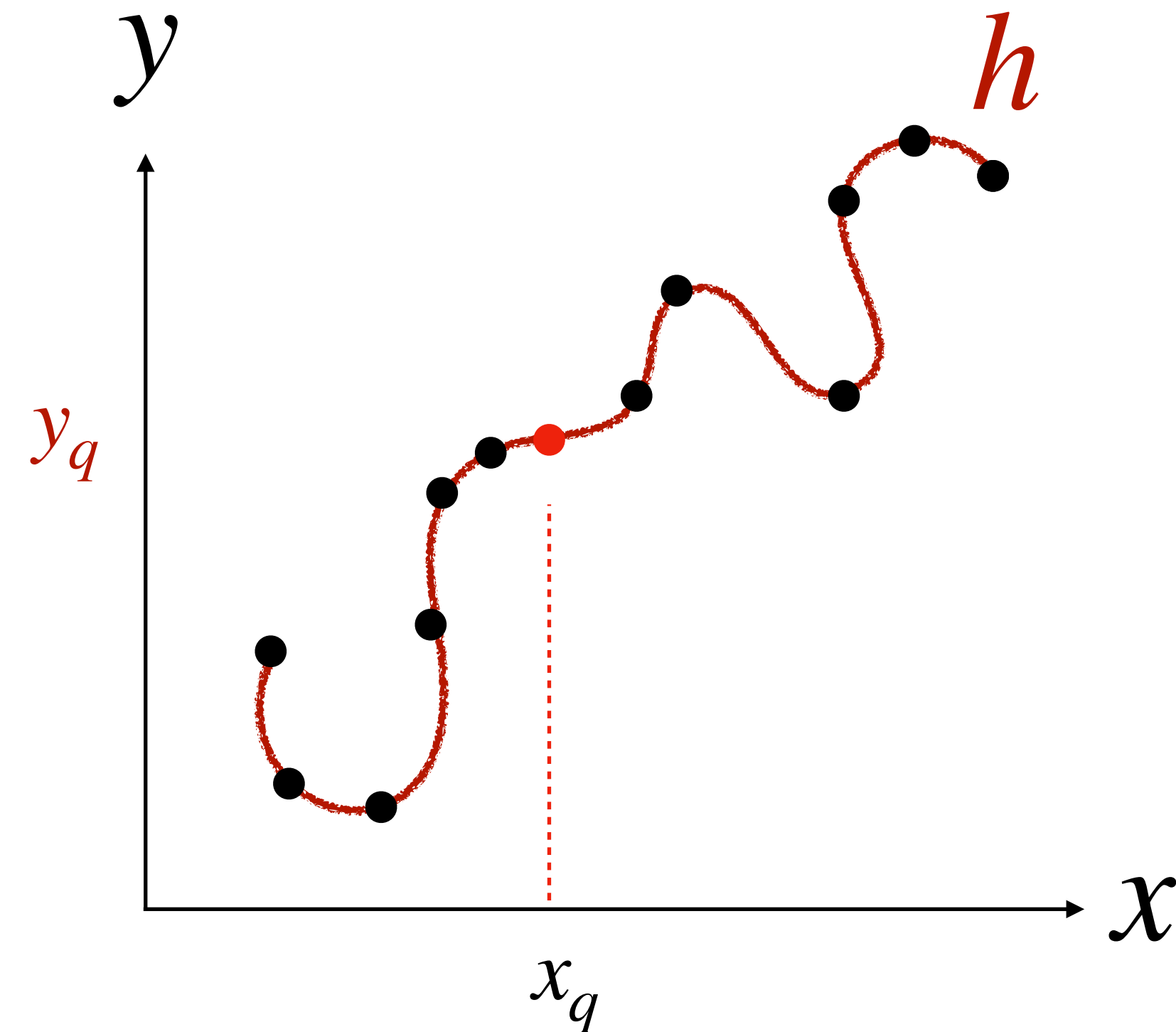


Given new input, what's the output?

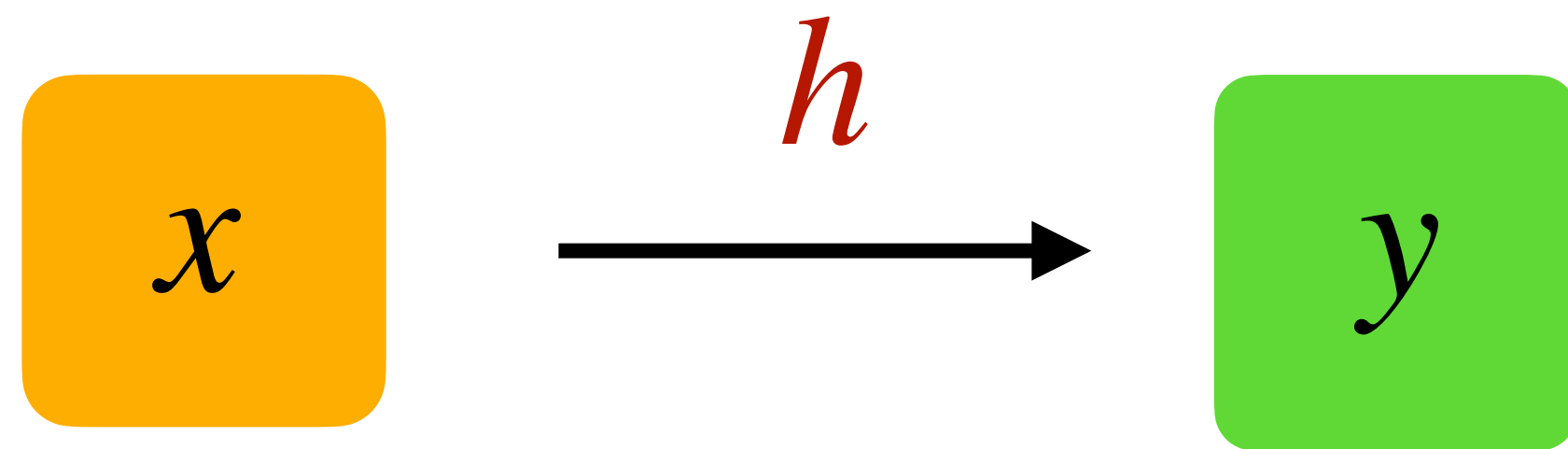


Given the data,
find a **function** h , a.k.a **hypothesis**,
that predicts an output, given an input

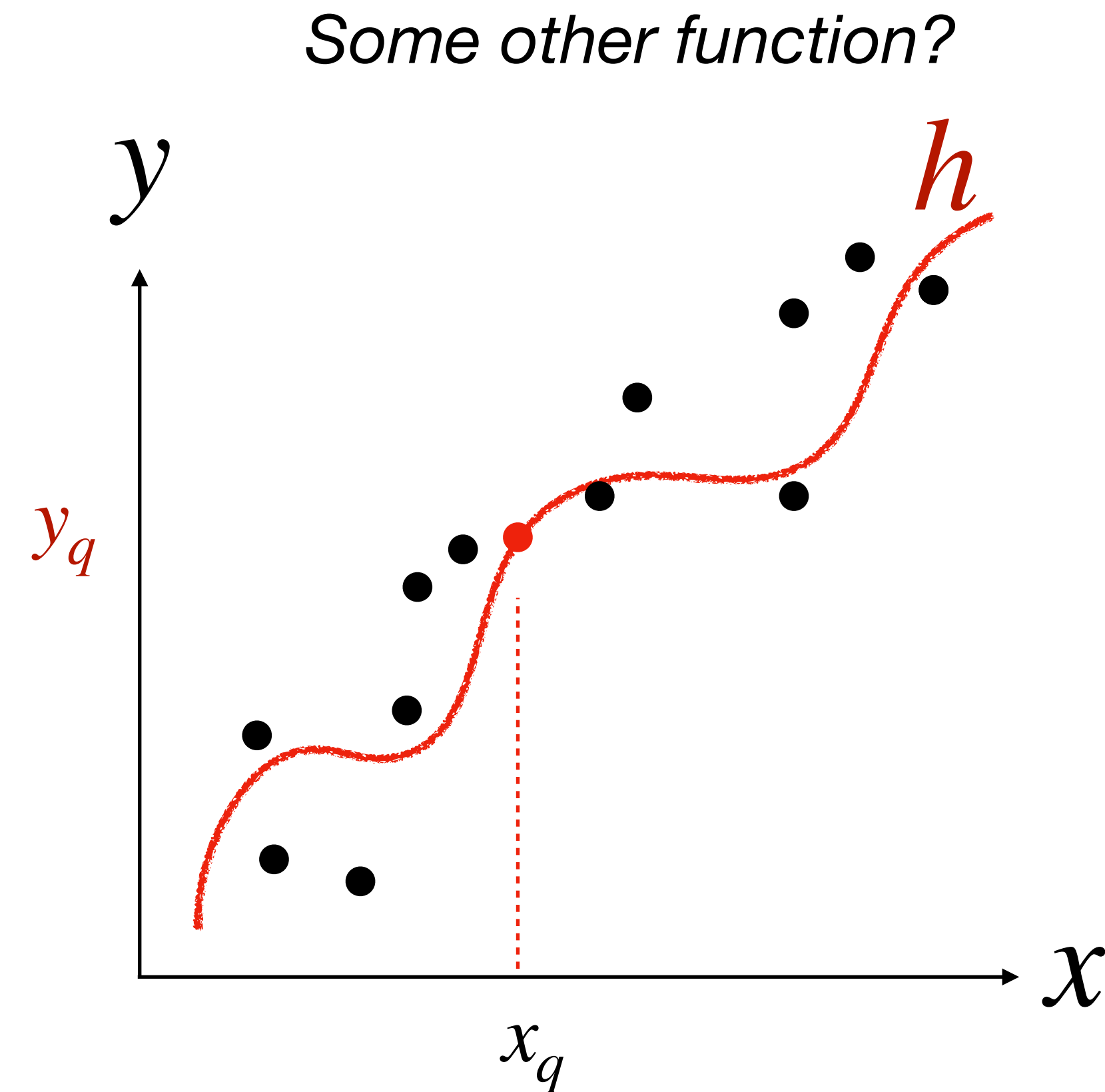
Polynomial Interpolation



Given new input, what's the output?

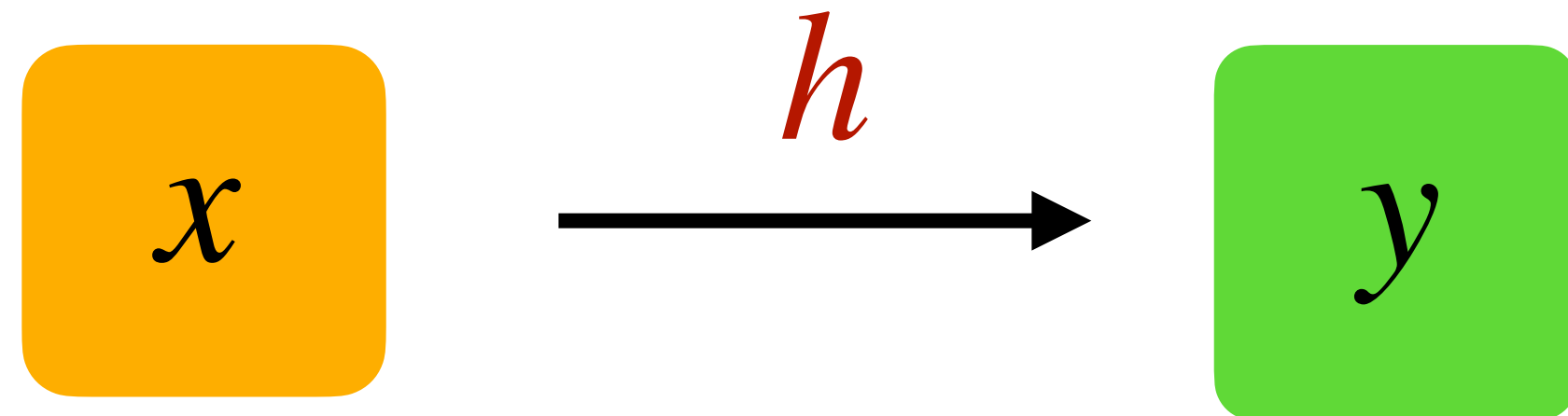


Given the data,
find a **function** h , a.k.a **hypothesis**,
that predicts an output, given an input



Given new input, what's the output?

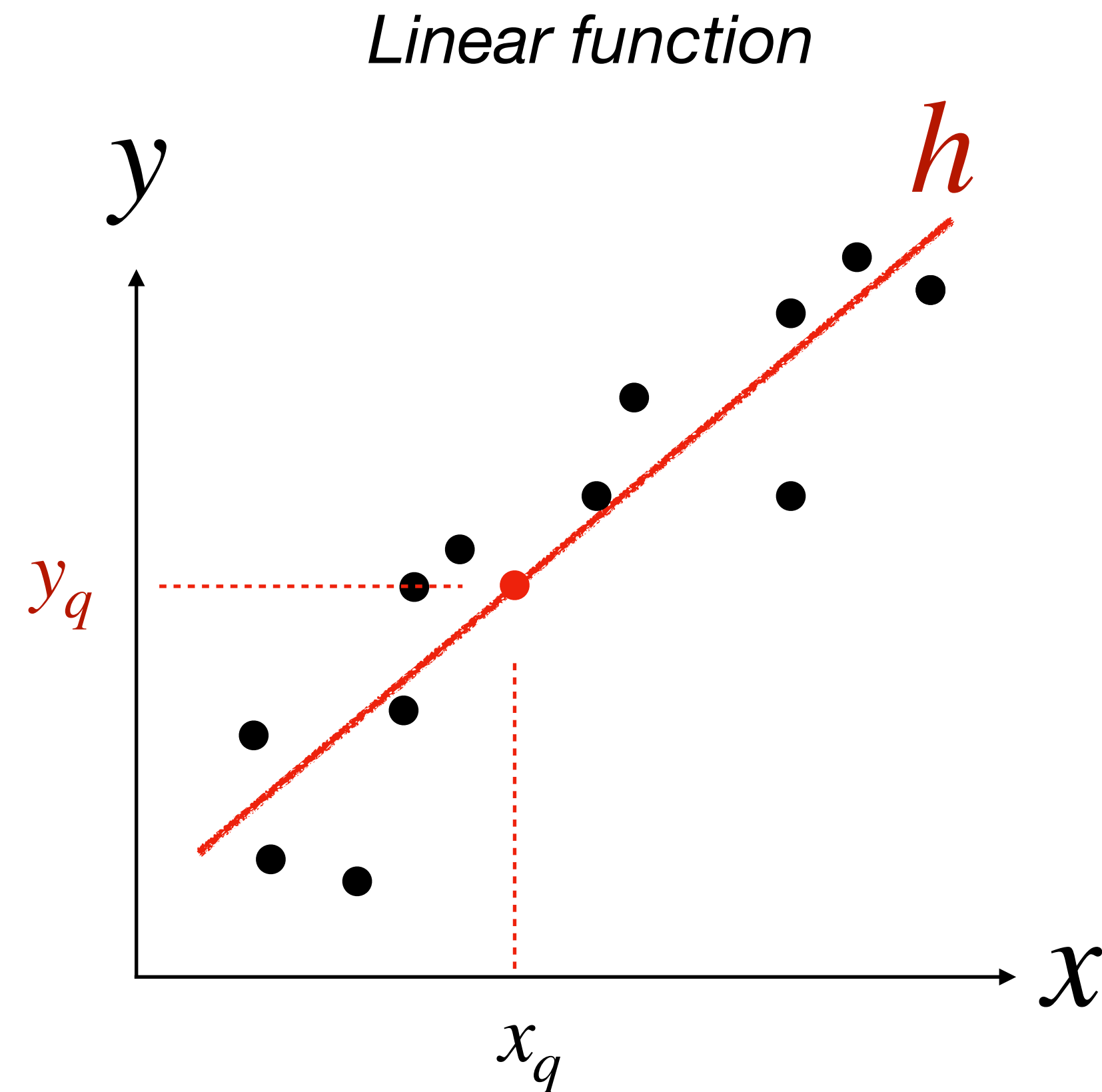
Assume a linear hypothesis



$$h(x) = ax + b$$

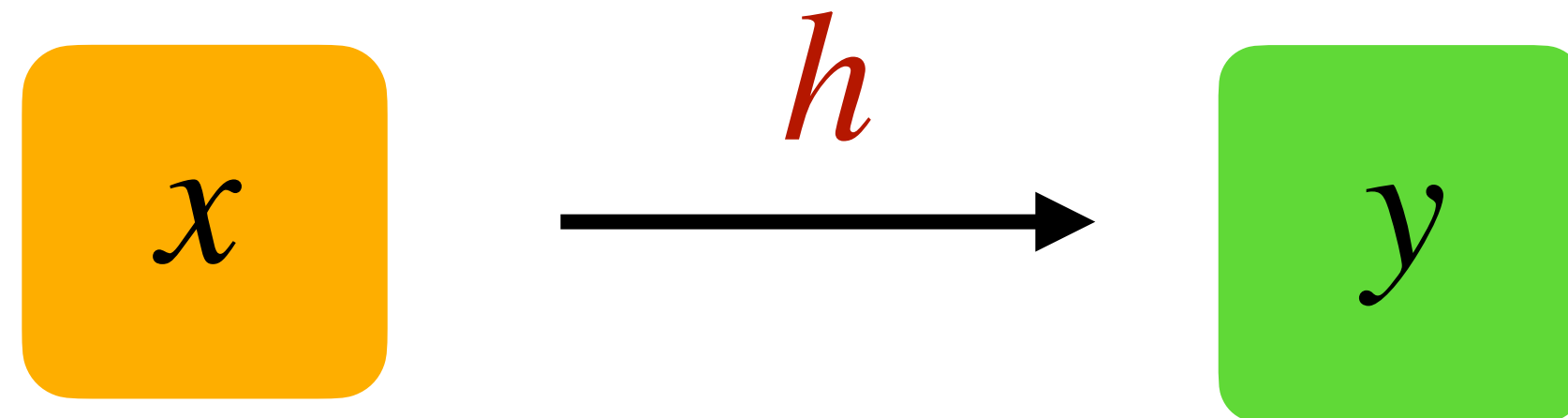
What are the **best** a and b
that **fit** the data?

a, b are **fitting** parameters



Assume a **linear** hypothesis

Assume a linear hypothesis



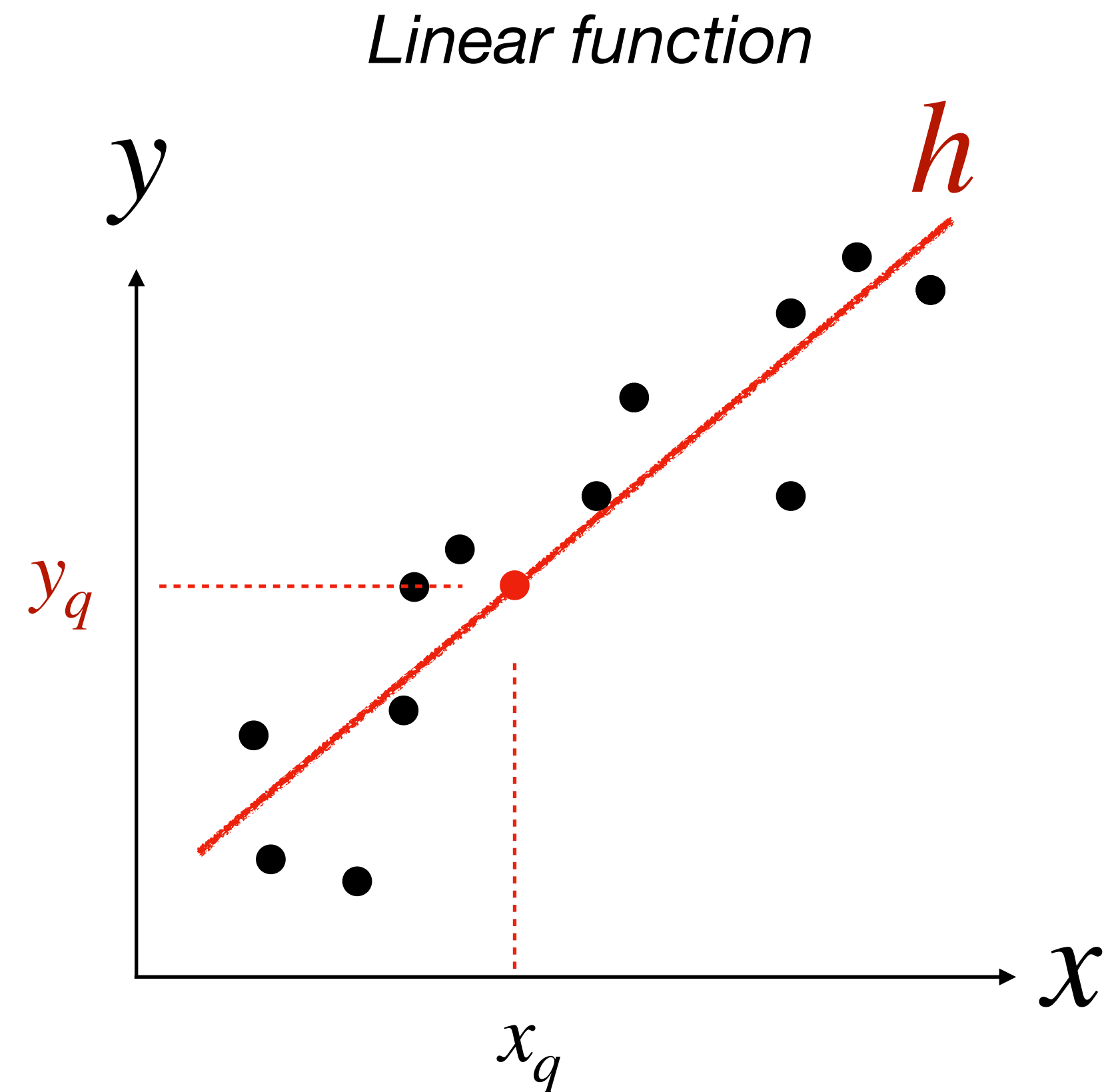
$$h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$h_{\theta}(x) = [\theta_0, \theta_1] \cdot [1, x]$$

Unknown
parameters

Input
features

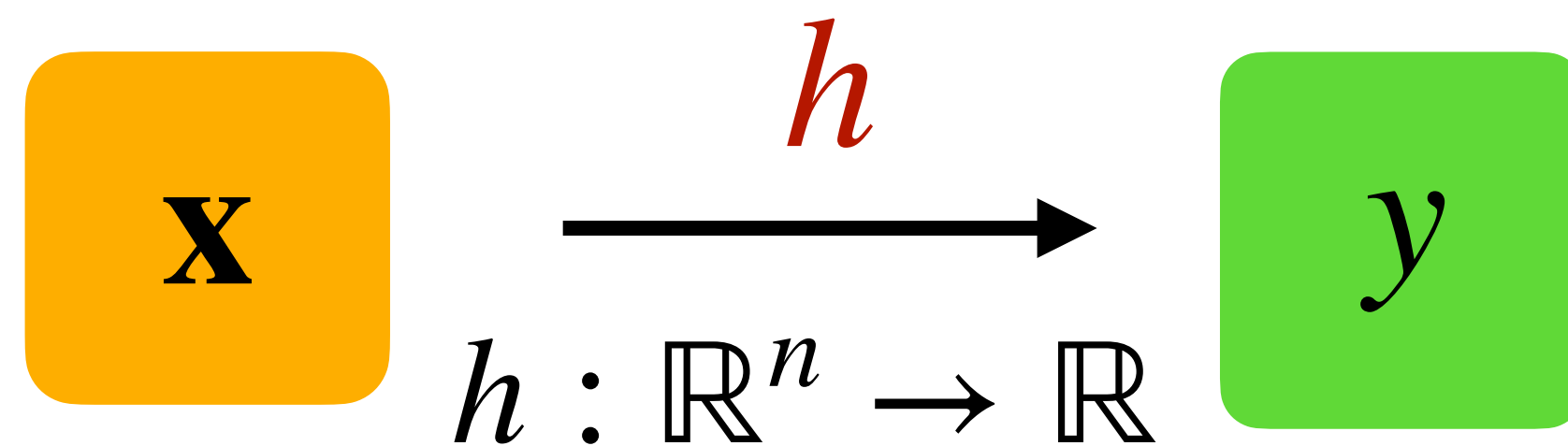
$$\theta \cdot \mathbf{x}$$



What's the **best** $\theta = [\theta_0, \theta_1]$, given the data ?

What happens if we have **more inputs**?

Assume a linear hypothesis



$$h_{\theta}(\mathbf{x}) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots$$

$$h_{\theta}(\mathbf{x}) = \underbrace{[\theta_0, \theta_1, \theta_2, \theta_3, \dots]}_{\theta} \cdot \underbrace{[1, x_1, x_2, x_3, \dots]}_{\mathbf{x}}$$

weights θ \mathbf{x} **inputs**

$$h_{\theta}(\mathbf{x}) = \theta \cdot \mathbf{x} = \theta^{\top} \mathbf{x}$$

Inputs

Output

x_1	x_2	y
$x_1^{(1)}$	$x_2^{(1)}$	$y^{(1)}$
$x_1^{(2)}$	$x_2^{(2)}$	$y^{(2)}$
$x_1^{(3)}$	$x_2^{(3)}$	$y^{(3)}$
$x_1^{(4)}$	$x_2^{(4)}$	$y^{(4)}$
\vdots	\vdots	\vdots

How do we pick the **best** parameters θ ?

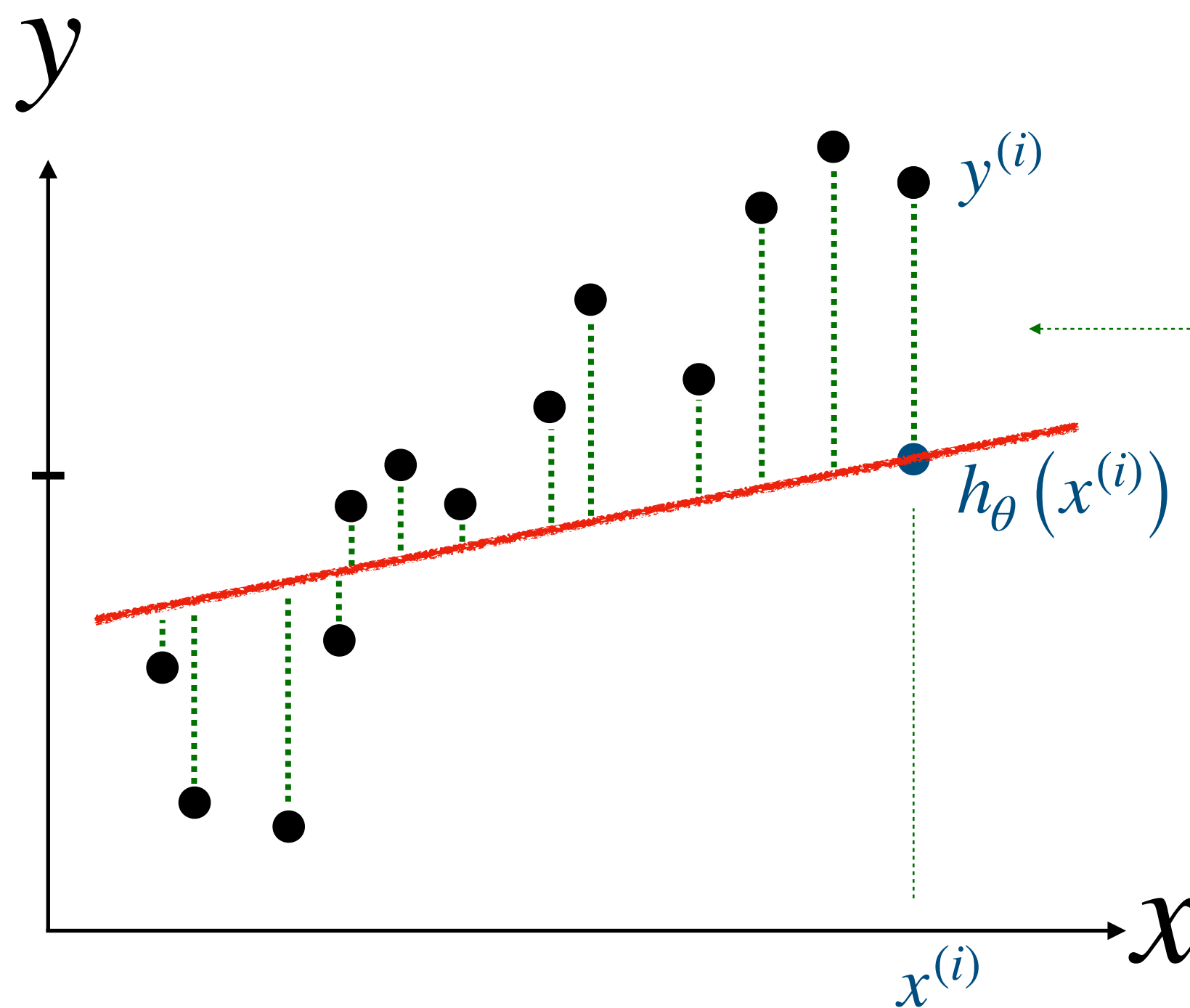
$$h_{\theta}(\mathbf{x}) = \theta^{\top} \mathbf{x} = \sum_{i=0}^d \theta_i x_i$$

Cost function

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

$$= \frac{1}{2} \sum_{i=1}^d \left(\theta^{\top} \mathbf{x}^{(i)} - y^{(i)} \right)^2$$

Ordinary least squares



distance $\left(h_{\theta}(x^{(i)}), y^{(i)} \right)$

$$h_{\theta}(x^{(i)}) - y^{(i)}$$

Residuals

$$\left| h_{\theta}(x^{(i)}) - y^{(i)} \right|$$

Absolute loss

$$\left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Square loss

Interactive Demo

https://colab.research.google.com/drive/1jEMvm_qlLneleOFDC5Andr7JstletVet?usp=sharing

Choose θ to **minimize** $J(\theta)$

Cost function

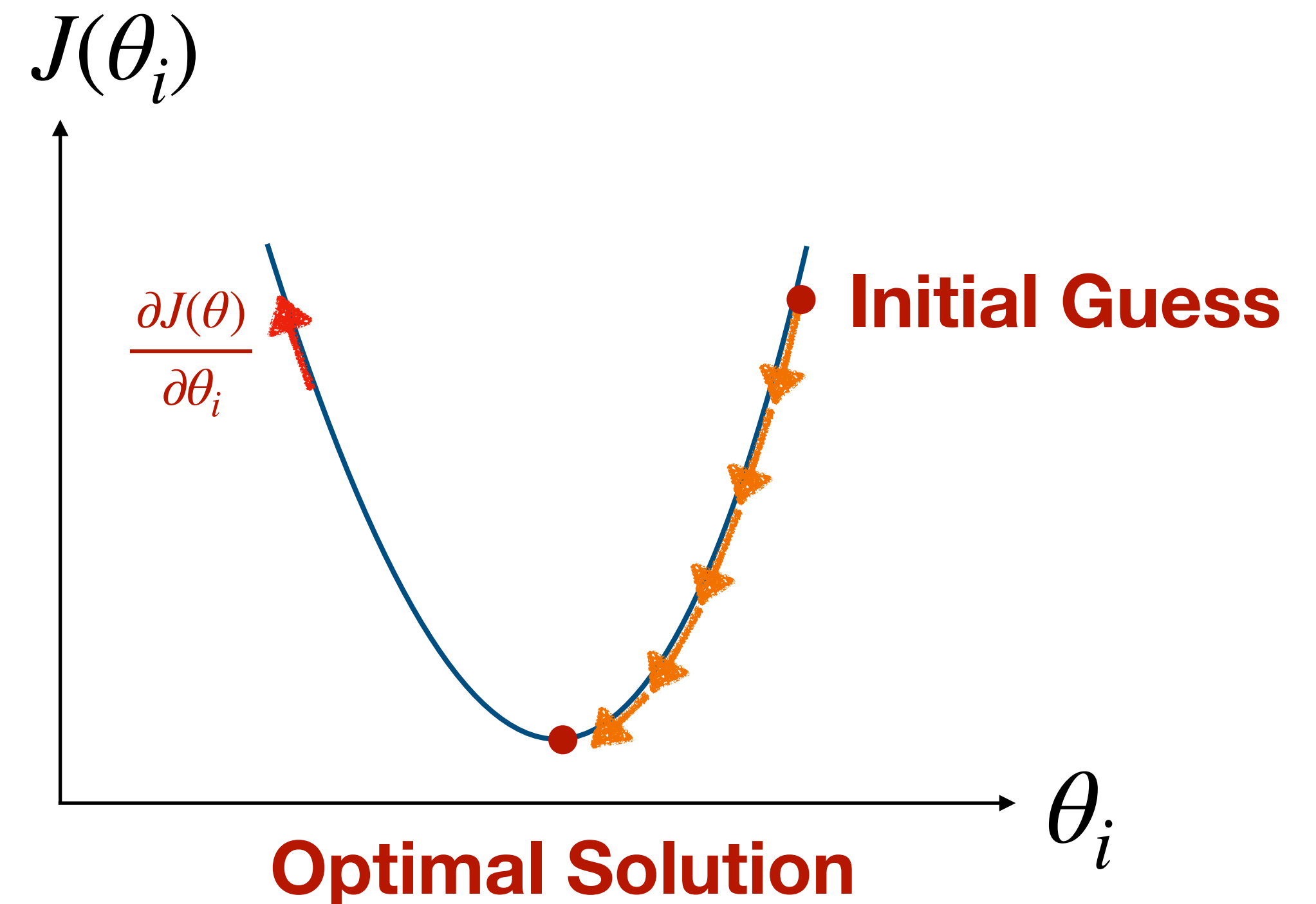
$$J(\theta) = \frac{1}{2} \sum_{i=1}^d \left(h_{\theta}(x^{(i)}) - y^{(i)} \right)^2$$

Gradient Descent Update

while not converged:

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

Learning Rate



Gradient can be computed explicitly

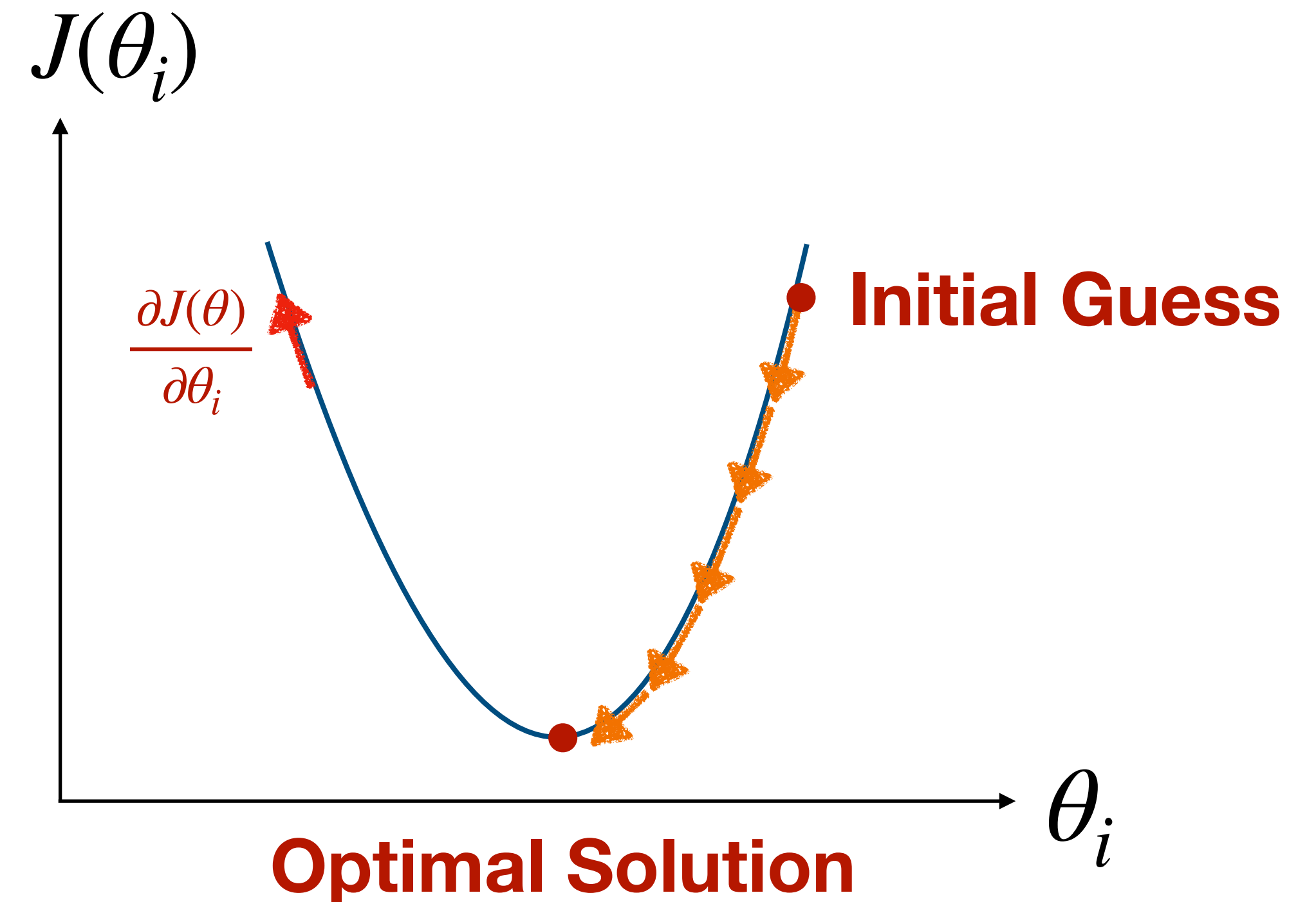
while not converged:

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

Learning Rate

Derive $\frac{\partial J(\theta)}{\partial \theta_i}$ explicitly, for one (x, y) pair

Assume $y = \theta_0 x + \theta_1$



Least Mean Squares (LMS)

A.K.A **Widrow-Hoff** learning rule

For a single training example $(x^{(i)}, y^{(i)})$:

$$\theta_j := \theta_j - \alpha \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) x_j^{(i)}$$

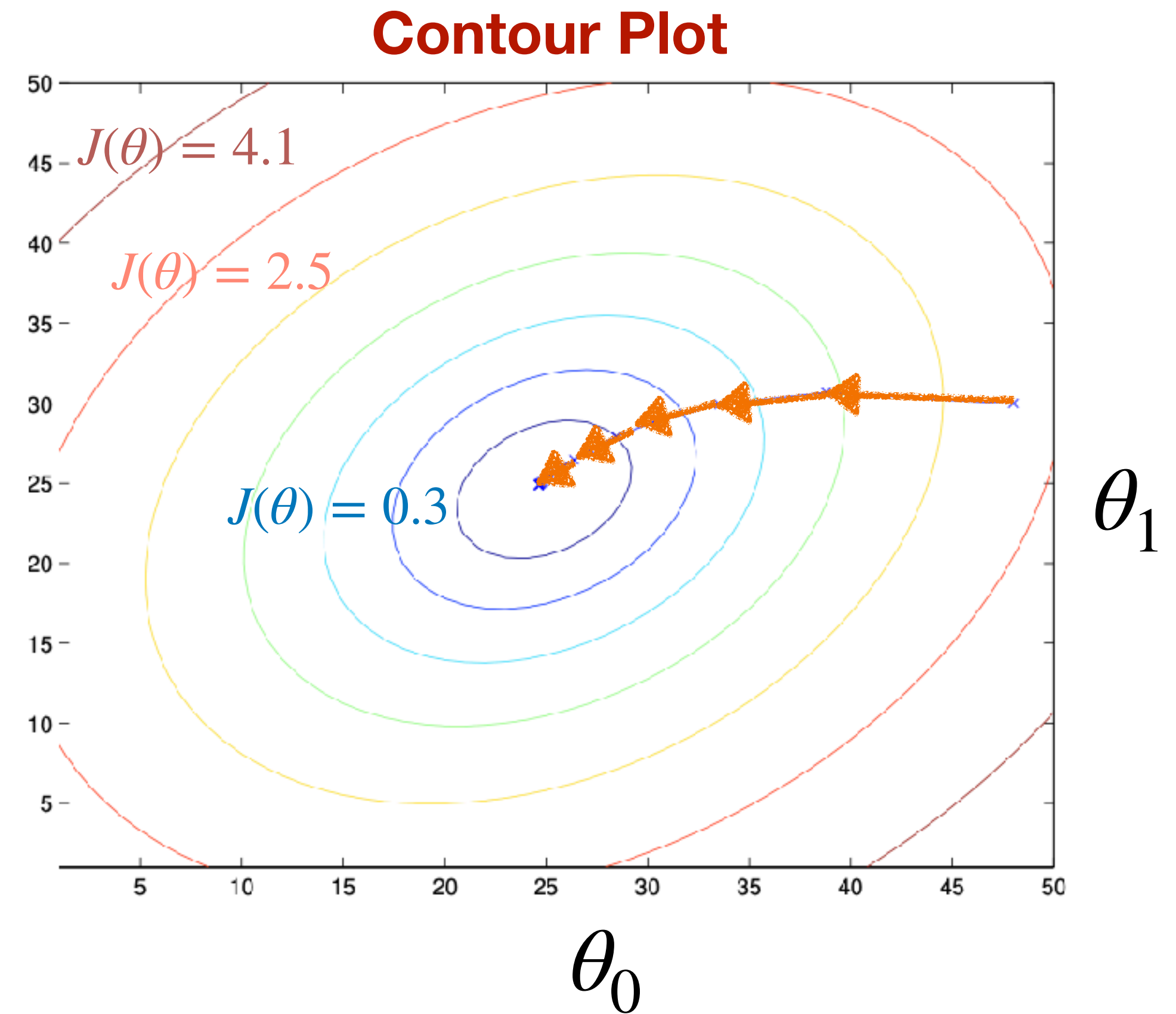
Batch Gradient Descent

for $t = 1 \dots T$: (Epochs)

for all parameters j :

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

Stack and vectorize



Least Mean Squares (LMS)

Batch Gradient Descent (vectorized)

for $t = 1 \dots T$:

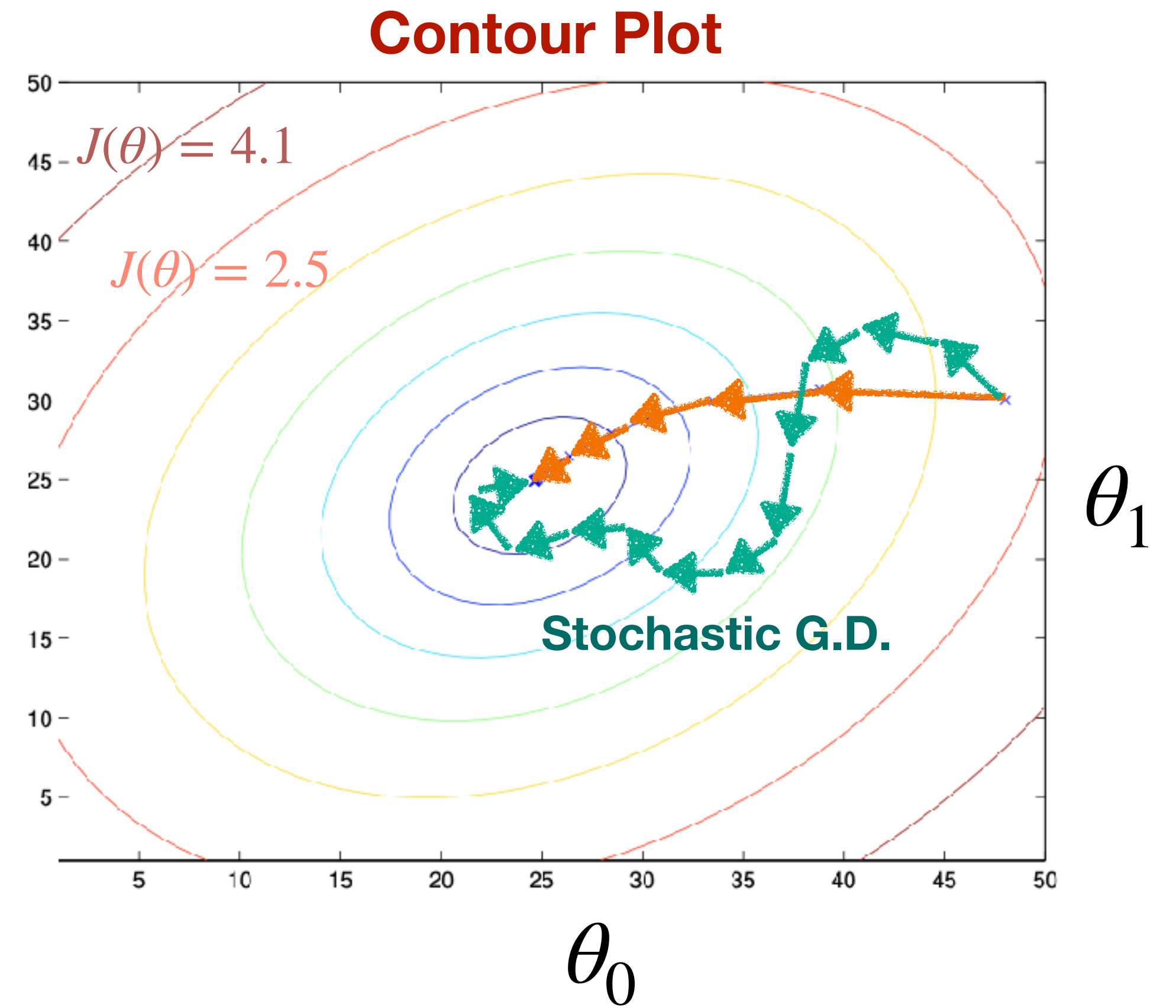
$$\theta := \theta - \alpha \sum_{i=1}^n \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) x^{(i)}$$

Stochastic Gradient Descent

for $t = 1 \dots T$:

for $i = 1 \dots n$:

$$\theta := \theta - \alpha \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) x^{(i)}$$



Summary

1. Assume a linear hypothesis

$$h_{\theta}(\mathbf{x}) = \theta^{\top} \mathbf{x} = \sum_{i=0}^d \theta_i x_i$$

2. Cost function

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

3. Minimize: Gradient Descent

$$\theta := \theta - \alpha \nabla_{\theta} J(\theta)$$

5. Predict unseen data

$$y_{pred} = h_{\hat{\theta}}(x_{new})$$

4. Optimal predictor

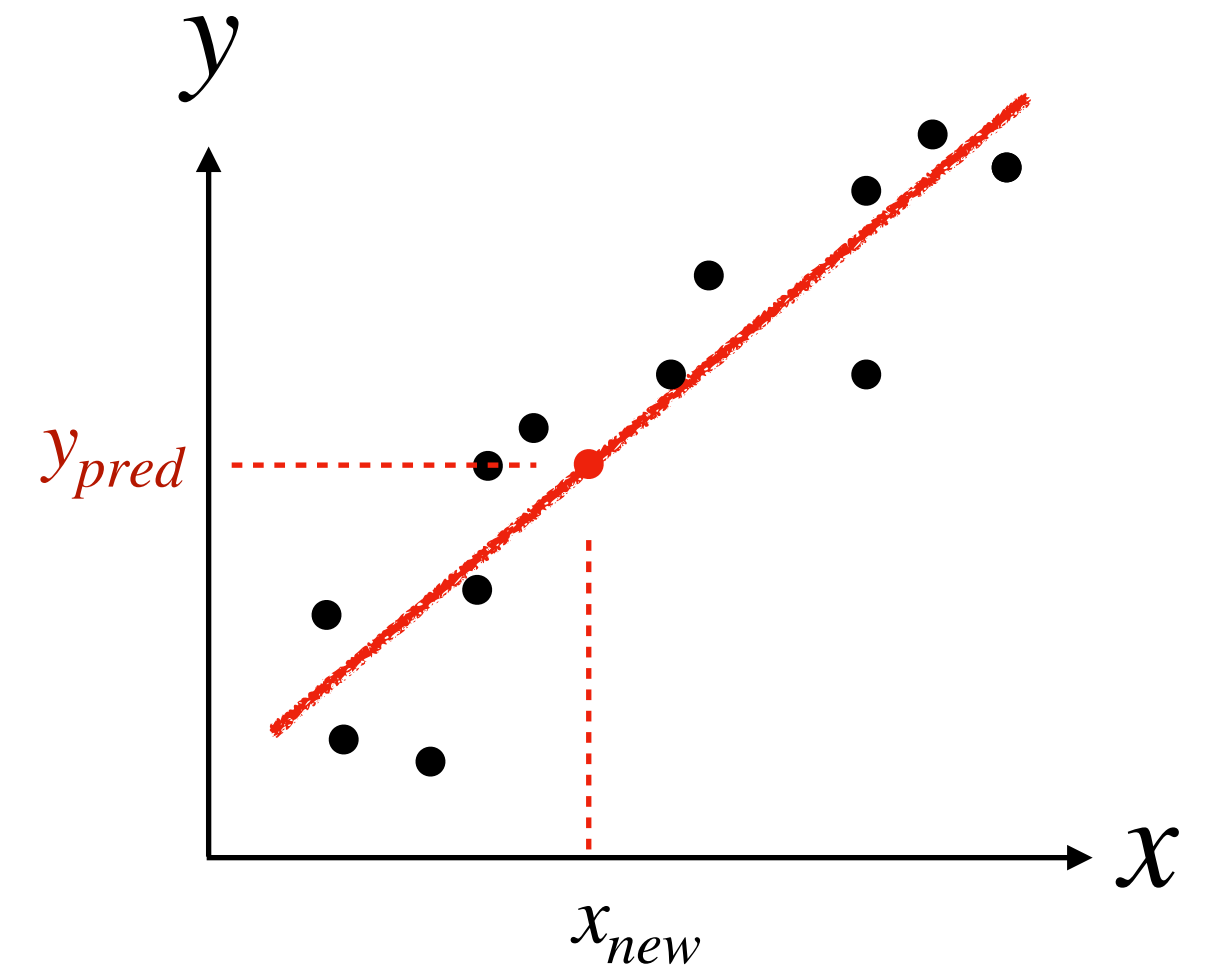
$$y = h_{\hat{\theta}}(x)$$

SGD

for $t = 1 \dots T$:

for $i = 1 \dots n$:

$$\theta := \theta - \alpha \left(h_{\theta}(x^{(i)}) - y^{(i)} \right) x^{(i)}$$



Can you find the minimum analytically?

Design matrix

$$X = \begin{bmatrix} \text{--} & x^{(1)} & \text{--} \\ \text{--} & x^{(2)} & \text{--} \\ \text{--} & \vdots & \text{--} \\ \text{--} & x^{(n)} & \text{--} \end{bmatrix}$$

Parameters

$$\vec{\theta} = \begin{bmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_m \end{bmatrix}$$

Output

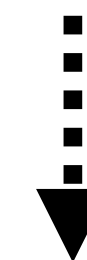
$$\vec{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(n)} \end{bmatrix}$$

Minimize

$$J(\theta) = \frac{1}{2} \left\| X\theta - \vec{y} \right\|_2^2$$



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



Normal Equation

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

Probabilistic Interpretation

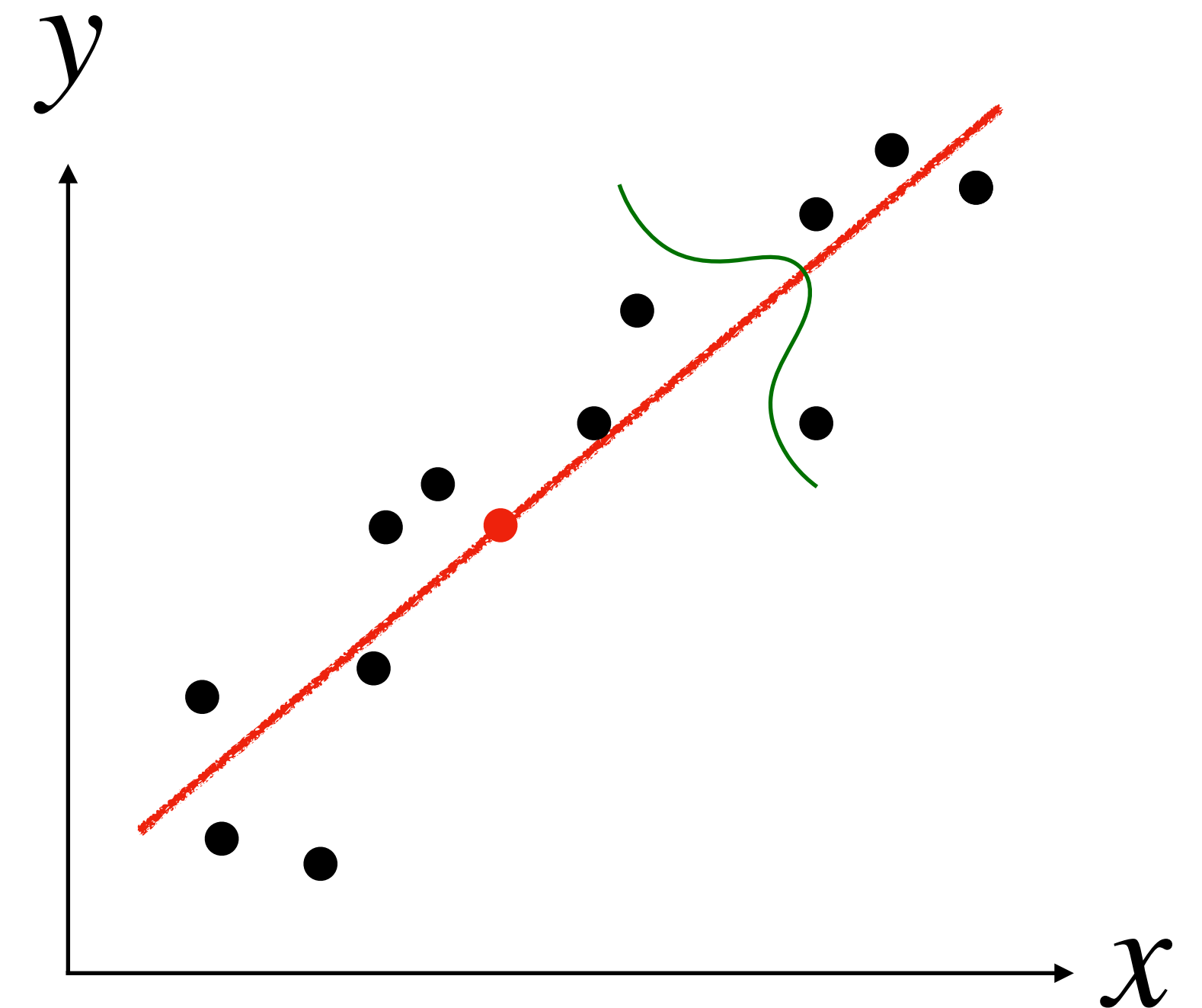
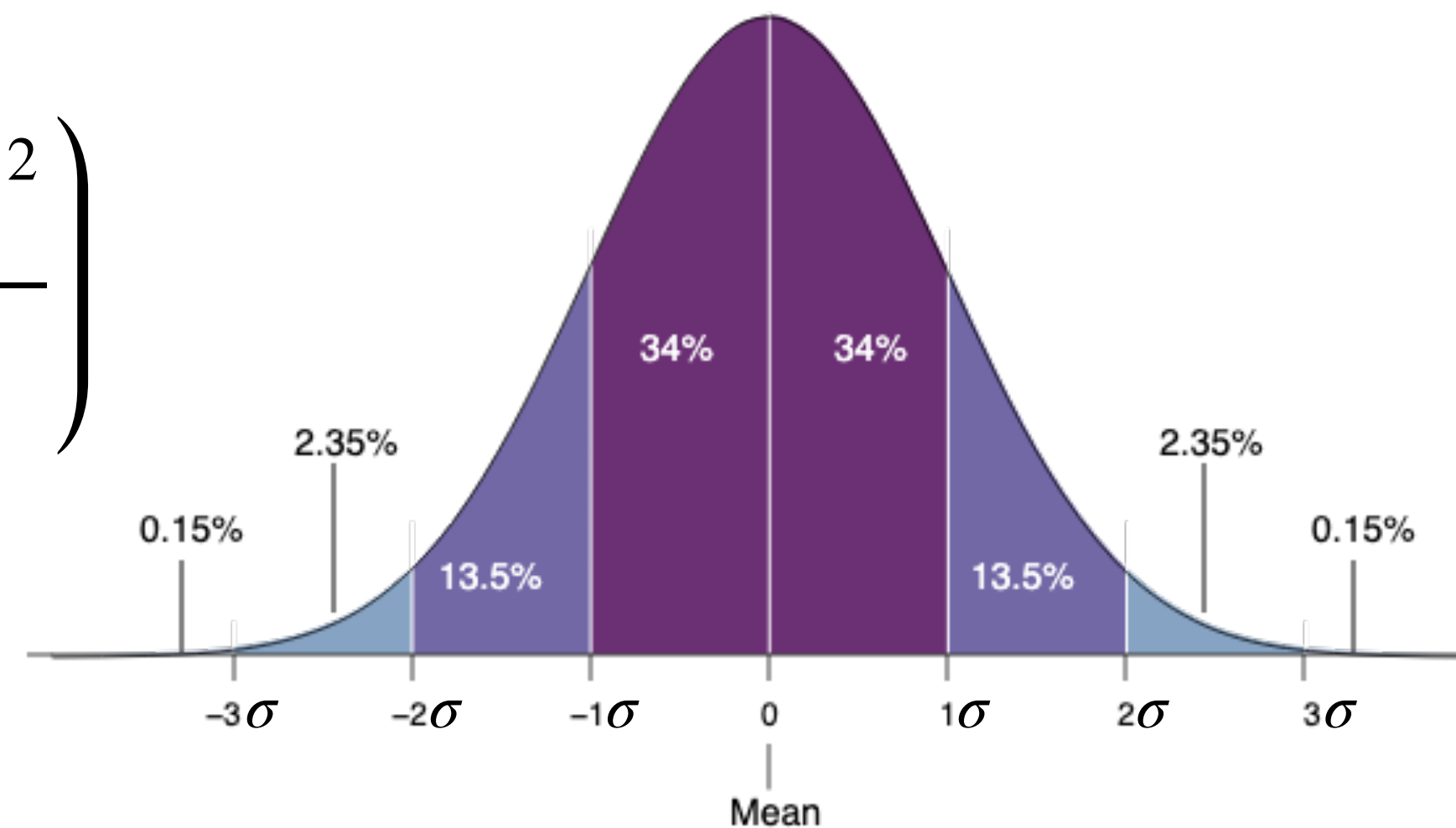
Assume noise is normally distributed around model

$$y^{(i)} = \theta^\top x^{(i)} + \varepsilon^{(i)}$$

Normally distributed

$$\mathcal{N}(0, \sigma^2)$$

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



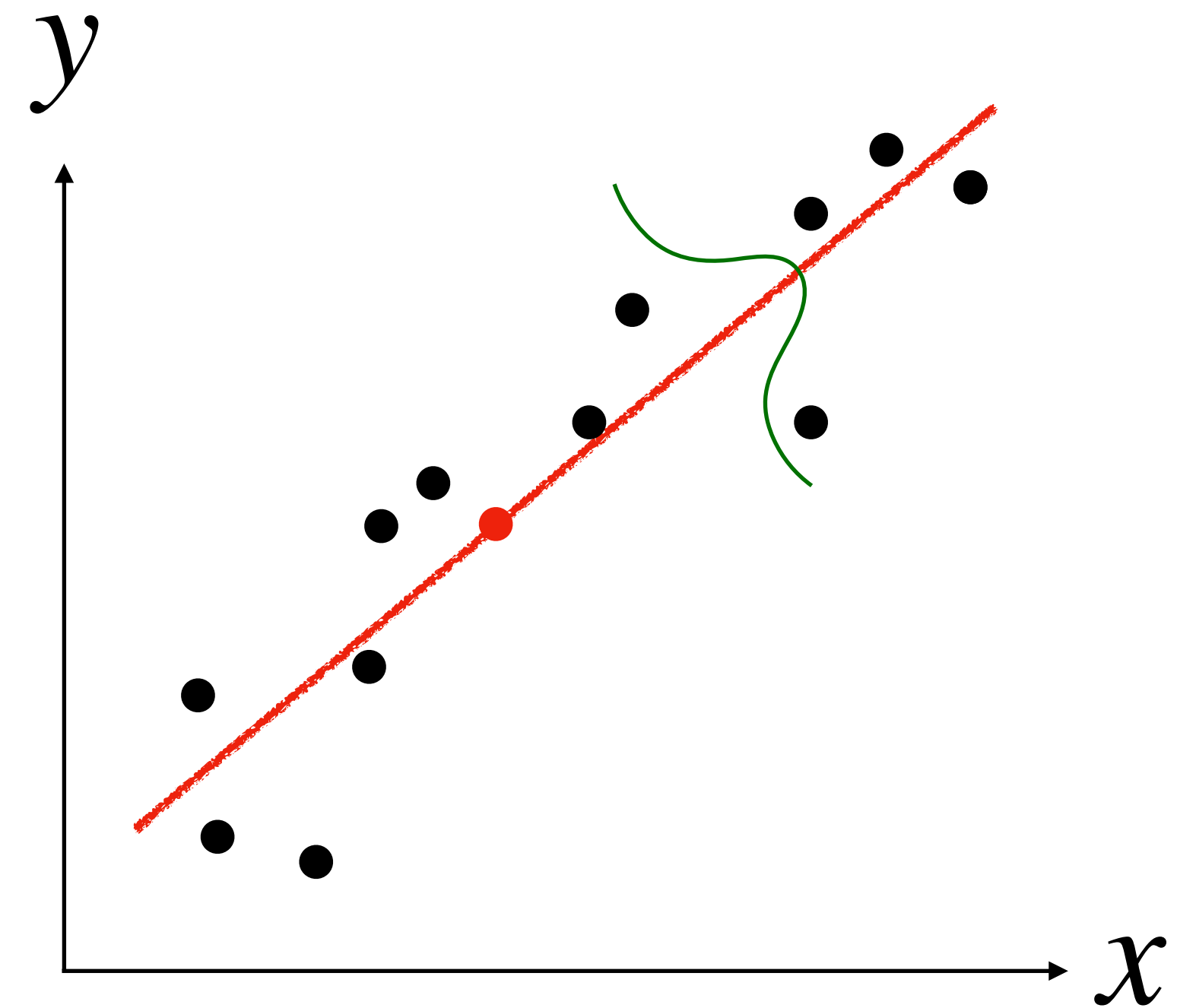
Probabilistic Interpretation

Assume noise is
normally distributed
around model

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

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

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

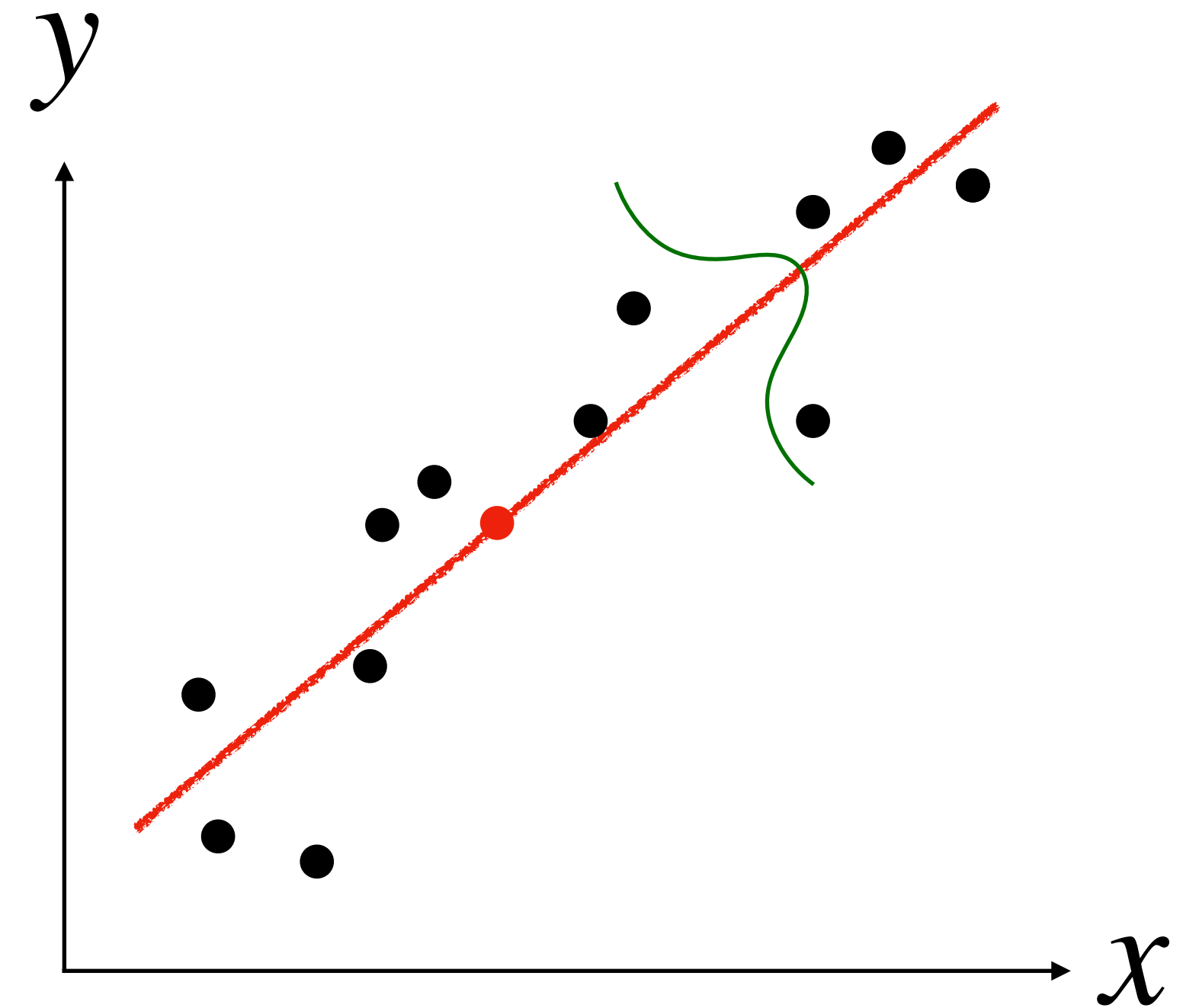


Likelihood of output given input

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

Log-likelihood

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



Maximize **Log-likelihood**

$$l(\theta) = \log L(\theta)$$

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

$$= n \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} \sum_{i=1}^n (y^{(i)} - \theta^\top x^{(i)})^2$$

$$\text{Maximize } l(\theta) \longrightarrow \text{Minimize } \frac{1}{2} \sum_{i=1}^n (y^{(i)} - \theta^\top x^{(i)})^2$$

