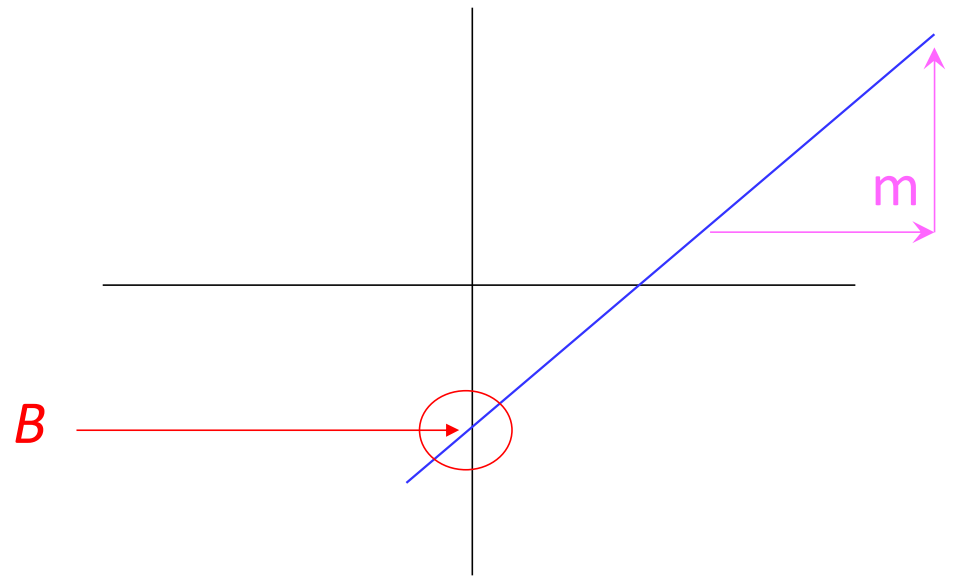# Linear Regression

# What is linear ?

- Remember this:
  - *Y=mX+B?*

A slope of 2 means that every 1-unit change in X yields a 2-unit change in Y.
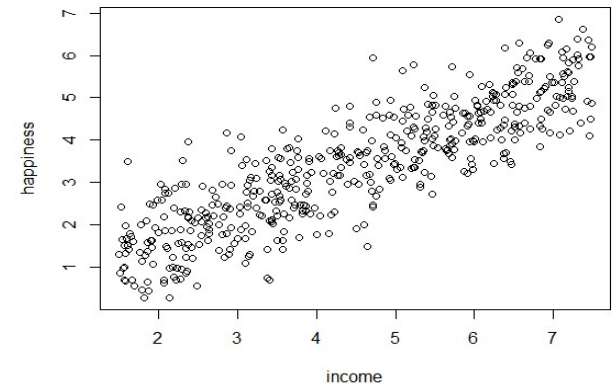
# Linear regression example



- Housing prices prediction

$$h_\theta(x) = \theta_0 + \theta_1 \times frontage + \theta_2 \times depth$$

- Estimating body weight by person height

- Estimate the happiness score of a person using its income

# Formulation

- We generally formulate the linear regression model in matrix form:
  - Y=Xw+$\varepsilon$
- the target value yi can be evaluated by
  - $y_i = \theta_0 + \theta_1 x_{i1} + \cdots + \theta_n x_{in} + \varepsilon_i$
  - Y represents a vector of length n containing the observed values
    Y=$(y_1, \ldots, y_m)^T$
  - $\varepsilon$ is a vector for errors $\varepsilon = (\varepsilon_1, \ldots, \varepsilon_m)^T$
  - X is a matrix of the features in which the column of ones incorporate the intercept
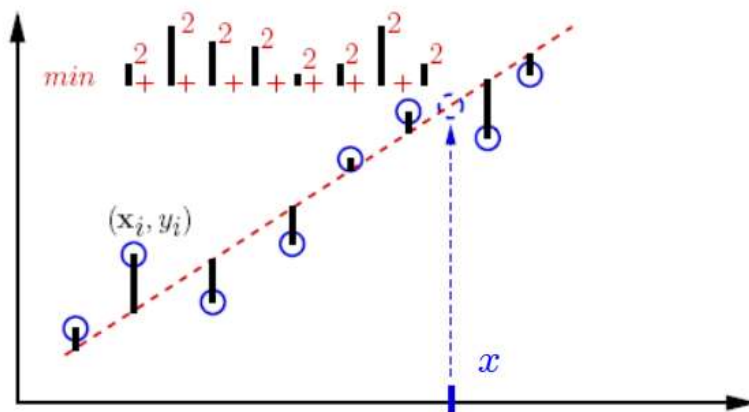
$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1n} \\ 1 & x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{m1} & x_{m2} & \cdots & x_{mn} \end{pmatrix}$$

- Hypothesis:

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \ldots + \theta_d x_d = \sum_{j=0}^{d} \theta_j x_j$$

Assume $x_0 = 1$

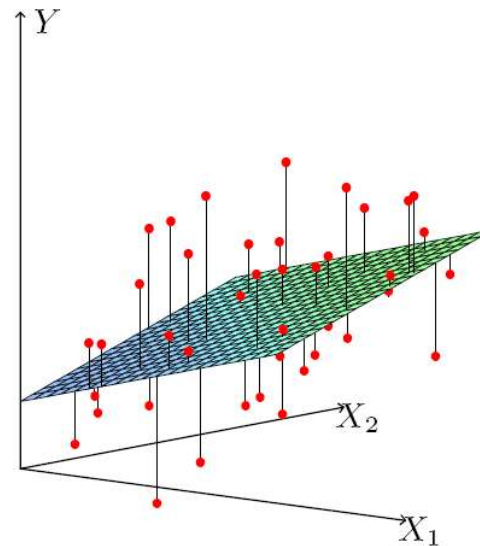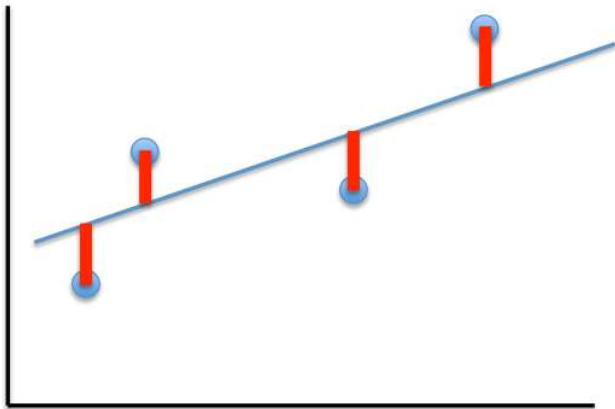- Fit model by minimizing sum of squared errors



east squares (LSQ)
The fitted line is used as a predictor

# Least square linear regression

- Cost Function

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

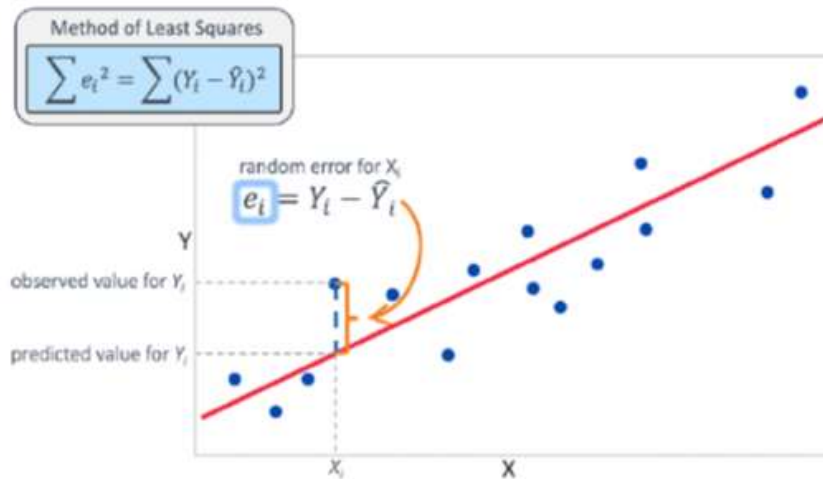- Fit by solving $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

# Optimization Problem

One suitable estimator of $\beta$ should be the one minimizing the sum of the squared errors $\|\epsilon\|_2^2 = \sum_{i=1}^{m} \epsilon_i^2 = \epsilon^T \epsilon$.

- $\sum_{i=1}^{m} \varepsilon_i^2 = \varepsilon^T \varepsilon \quad = (Y - X\theta)^T (Y - X\theta)$

$$= Y^T Y - 2\theta^T X^T Y + \theta^T X^T X \theta$$

### Method of Least Squares

$$\sum e_i^2 = \sum (Y_i - \hat{Y}_i)^2$$

random error for $X_i$
$$e_i = Y_i - \hat{Y}_i$$

Y

observed value for $Y_i$

predicted value for $Y_i$

$X_i$

X

### Vector derivative

| $f(\mathbf{x})$ | $\rightarrow$ | $\frac{df}{d\mathbf{x}}$ |
|---|---|---|
| $\mathbf{x}^T \mathbf{B}$ | $\rightarrow$ | $\mathbf{B}$ |
| $\mathbf{x}^T \mathbf{b}$ | $\rightarrow$ | $\mathbf{b}$ |
| $\mathbf{x}^T \mathbf{x}$ | $\rightarrow$ | $2\mathbf{x}$ |
| $\mathbf{B}\mathbf{x}$ | $\rightarrow$ | $\mathbf{B}^T$ |

- Differentiating this term and setting it to zero, we find that the estimate for $\theta$, which minimizes the squared error, satisfies the equation :
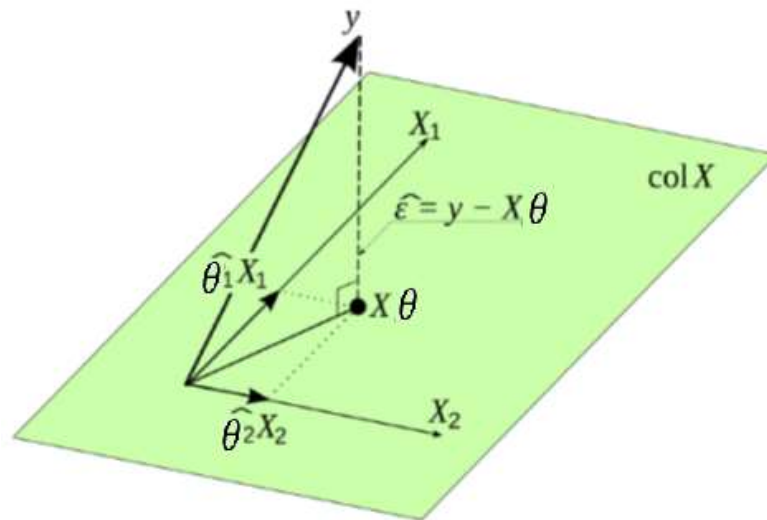
$$X^T X \theta = X^T Y$$

- Provided $X^T X$ is invertible :
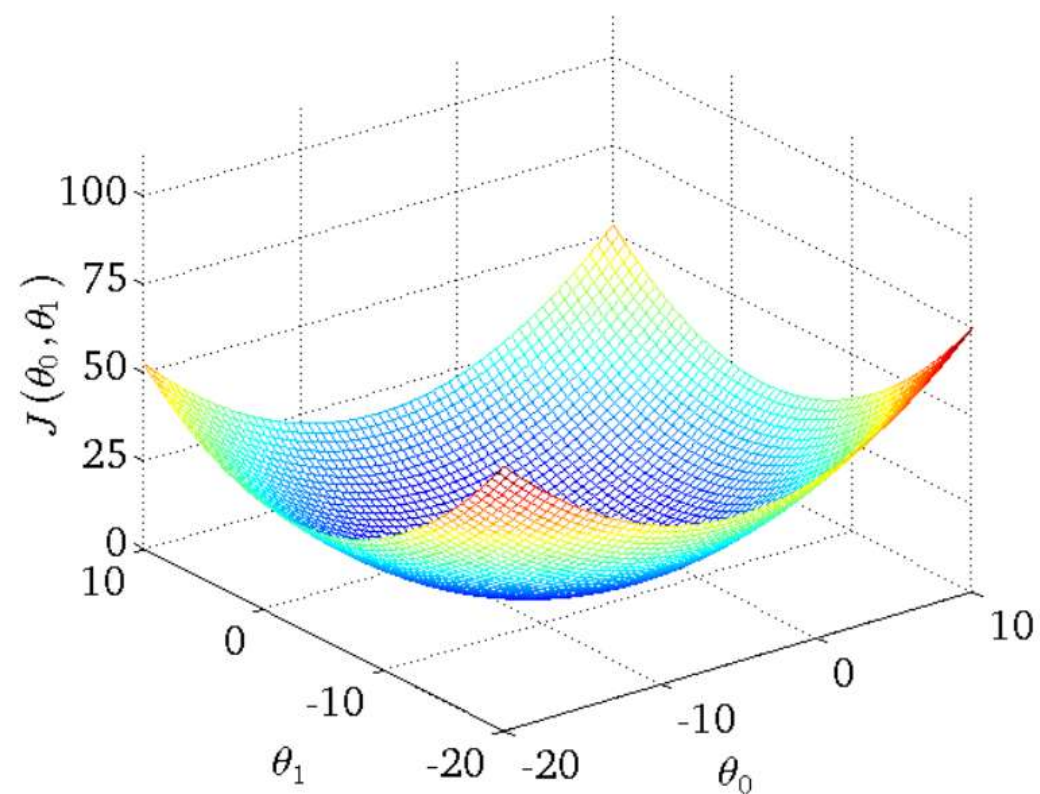
$$\theta = (X^T X)^{-1} X^T Y$$

# Geometric Approach

Another way to looking at this problem is to say we want a solution that lies in the space spanned by X become as close as possible to Y.

In this way, the systematic component $X\theta$ is projection of Y onto space spanned by X and residuals are Y- $X\theta$

# Intuition Behind Cost Function

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

# Intuition Behind Cost Function

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

# $J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# $h_\theta(x)$

## (for fixed $\theta_0, \theta_1$, this is a function of x)



# $J(\theta_0, \theta_1)$

## (function of the parameters $\theta_0, \theta_1$)

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

# Basic search Procedure

- Choose initial value for $\theta$
- Until we reach a minimum:
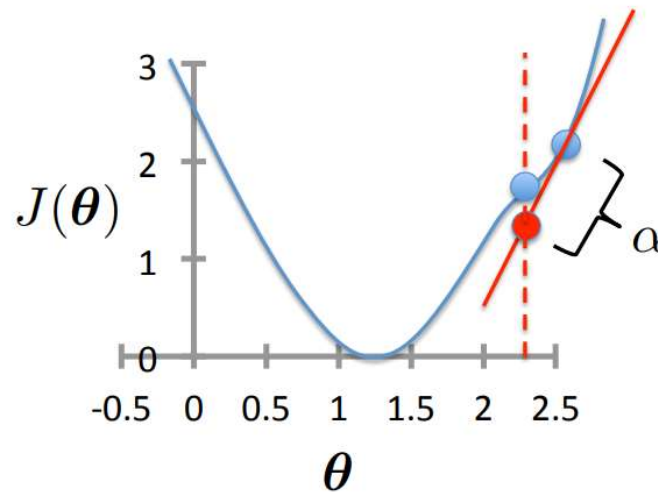  - Choose a new value for $\theta$ to reduce $J(\theta)$

# Gradient Descent

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence

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

simultaneous update
for $j = 0 \dots d$

learning rate (small)
e.g., $\alpha = 0.05$

- Initialize $\boldsymbol{\theta}$

- Repeat until convergence

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

simultaneous update
for $j = 0 \ldots d$

For Linear Regression:
$$\frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) = \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right)^2$$

$$= \frac{\partial}{\partial \theta_j} \frac{1}{2n} \sum_{i=1}^{n} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right)^2$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right) \times \frac{\partial}{\partial \theta_j} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right)$$

$$= \frac{1}{n} \sum_{i=1}^{n} \left( \sum_{k=0}^{d} \theta_k x_k^{(i)} - y^{(i)} \right) x_j^{(i)}$$

- Initialize $\boldsymbol{\theta}$

- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^{n} \left( h_{\boldsymbol{\theta}} \left( \boldsymbol{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)} \quad \text{simultaneous update for } j = 0 \ldots d$$

Assume convergence when $\|\boldsymbol{\theta}_{new} - \boldsymbol{\theta}_{old}\|_2 < \epsilon$
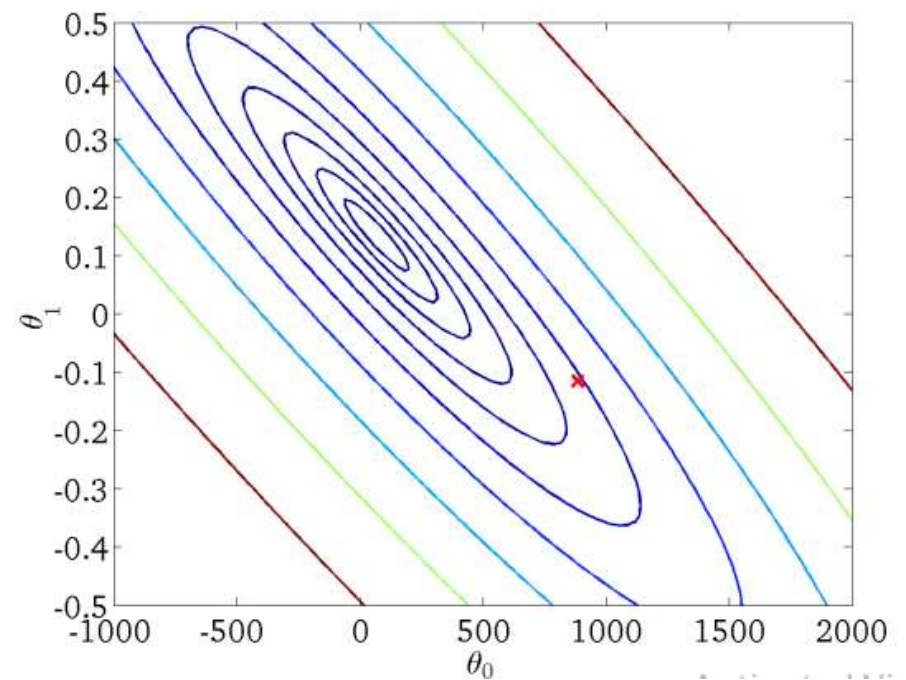
# Gradient Descent

$$h_\theta(x)$$

(for fixed $\theta_0$, $\theta_1$, this is a function of x)

$$J(\theta_0, \theta_1)$$

(function of the parameters $\theta_0$, $\theta_1$)



h(x) = -900 − 0.1 x

× Training data
— Current hypothesis

# $h_\theta(x)$

(for fixed $\theta_0$, $\theta_1$, this is a function of x)

# $J(\theta_0, \theta_1)$

(function of the parameters $\theta_0$, $\theta_1$)

# $h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)

$h_\theta(x)$

(for fixed $\theta_0, \theta_1$, this is a function of x)

$J(\theta_0, \theta_1)$

(function of the parameters $\theta_0, \theta_1$)
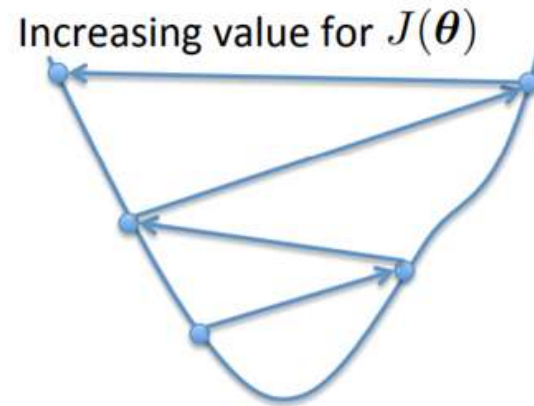
# Choosing α

### α too small

slow convergence

### α too large

Increasing value for $J(\boldsymbol{\theta})$

- May overshoot the minimum
- May fail to converge
- May even diverge

To see if gradient descent is working, print out $J(\boldsymbol{\theta})$ each iteration
- The value should decrease at each iteration
- If it doesn't, adjust α