

Linear regression with multiple variable

Multiple features

Multiple features (variables).

single feature	Size (feet ²)	Price (\$1000)	goal
	x	y	
	2104	460	
	1416	232	
	1534	315	
	852	178	
	

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

hypothesis

Multiple features (variables).

Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_1	x_2	x_3	x_4	y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

$m=47$

Notation:

n = number of features

$x^{(i)}$ = input (features) of i^{th} training example.

$x_j^{(i)}$ = value of feature j in i^{th} training example.

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 40 \end{bmatrix}$$

$$x_3^{(2)} = 2$$

Size (feet) ²	Number of bedroom s	Number of floors	Age of home (years)	Price (\$1000)
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

In the training set above, what is $x_1^{(4)}$?

- a) The size (in feet²) of the 1st home in the training set
- b) The age (in years) of the 1st home in the training set
- c) The size (in feet²) of the 4th home in the training set
- d) The age (in years) of the 4th home in the training set

Hypothesis:

Previously: $h_{\theta}(x) = \theta_0 + \theta_1 x$ single variable

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \theta_4 x_4$$

$$E.g. \quad h_{\theta}(x) = 80 + 0.3x_1 + 1.3x_2 + 4.2x_3 - 0.5x_4$$

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

For convenience of notation, define $x_0 = 1$. $x_0^{(i)} = 1$

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in R^{n+1} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in R^{n+1} \qquad \underbrace{[\theta_0 \quad \theta_1 \quad \theta_2 \quad \cdots \quad \theta_n]}_{\theta^T}$$

$$\begin{aligned} h_{\theta}(x) &= \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n \\ &= \theta^T x \end{aligned}$$

Multivariate linear regression.

Linear regression with
multiple variable

Gradient descent
for multiple variables

Hypothesis: $h_{\theta}(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$

Parameters: $\theta_0, \theta_1, \dots, \theta_n$ θ

Cost function:

$$\underbrace{J(\theta_0, \theta_1, \dots, \theta_n)}_{J(\theta)} = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

Gradient descent:

Repeat {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \dots, \theta_n) \quad J(\theta) - \text{derivative term}$$

}

(simultaneously update for every $j = 0, \dots, n$)

When there are n features, we define the cost function as

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

For linear regression, which of the following are also equivalent and correct definitions of $J(\theta)$?

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

(b) $J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$ (inner sum starts at 0)

(c) $J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=1}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$ (inner sum starts at 1)

(d) $J(\theta) = \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - \left(\sum_{j=0}^n y_j^{(i)} \right) \right)^2$

Gradient Descent

Previously (n=1):

Repeat {

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})$$

$$\frac{\partial}{\partial \theta_0} J(\theta)$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

(simultaneously update θ_0, θ_1)

}

New algorithm ($n \geq 1$):

Repeat {

$$\frac{\partial}{\partial \theta_j} J(\theta)$$

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

(simultaneously update θ_j for
 $j = 0, \dots, n$)

}

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...

Linear regression with
multiple variable

Gradient descent in practice I: Feature Scaling

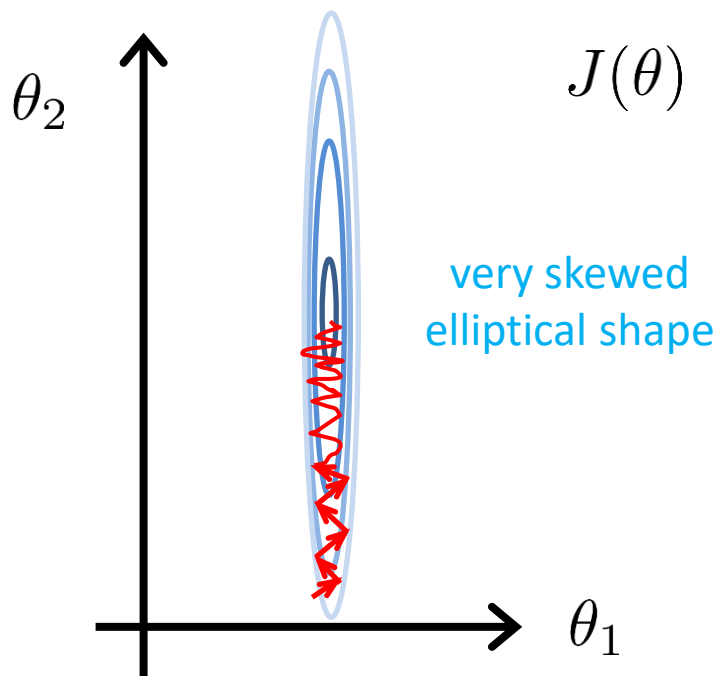
Fundamentals of Machine Learning

Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. $x_1 = \text{size (0-2000 feet}^2\text{)}$

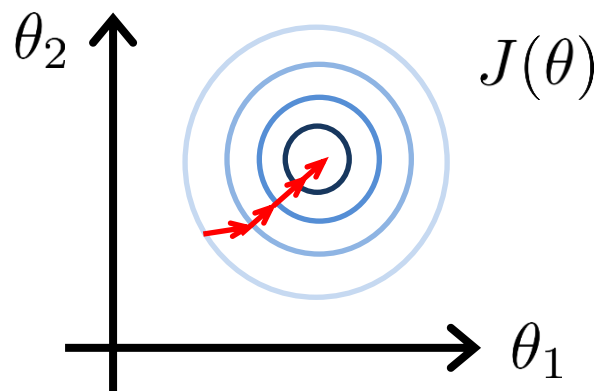
$x_2 = \text{number of bedrooms (1-5)}$



$$x_1 = \frac{\text{size (feet}^2\text{)}}{2000}$$

$$x_2 = \frac{\text{number of bedrooms}}{5}$$

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$



Feature Scaling

Get every feature into approximately a $-1 \leq x_i \leq 1$ range.

$$x_0 = 1$$

$$0 \leq x_1 \leq 2 \quad \checkmark$$

$$-3 \text{ to } 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.7 \quad \checkmark$$

$$-1/3 \text{ to } 1/3 \quad \checkmark$$

$$-100 \leq x_3 \leq 100 \quad \times$$

$$0.003 \leq x_4 \leq 0.01 \quad \times$$

Mean normalization

Replace x_i with $x_i - \mu_i$ to make features have approximately zero mean
(Do not apply to $x_0 = 1$).

E.g. $x_1 = \frac{\text{size} - 1000}{2000}$

average size = 1000

$$x_2 = \frac{\#bedrooms - 2}{5}$$

average bedrooms # = 2

$$-0.5 \leq x_1 \leq 0.5, -0.5 \leq x_2 \leq 0.5$$

$$x_i \leftarrow \frac{x_i - \mu_i}{S_i}$$

← average value of x_i in the training set

← Range of values of that feature

“maximum – minimum” or standard deviation

Suppose you are using a learning algorithm to estimate the price of houses in a city. You want one of your features x_i to capture the age of the house. In your training set, all of your houses have an age between 30 and 50 years, with an average age of 38 years. Which of the following would you use as features, assuming you use feature scaling and mean normalization?

(a) $x_i = \text{age of house}$

(b) $x_i = \frac{\text{age of house}}{50}$

(c) $x_i = \frac{\text{age of house} - 38}{50}$

(d) $x_i = \frac{\text{age of house} - 38}{20}$

Linear regression with
multiple variable

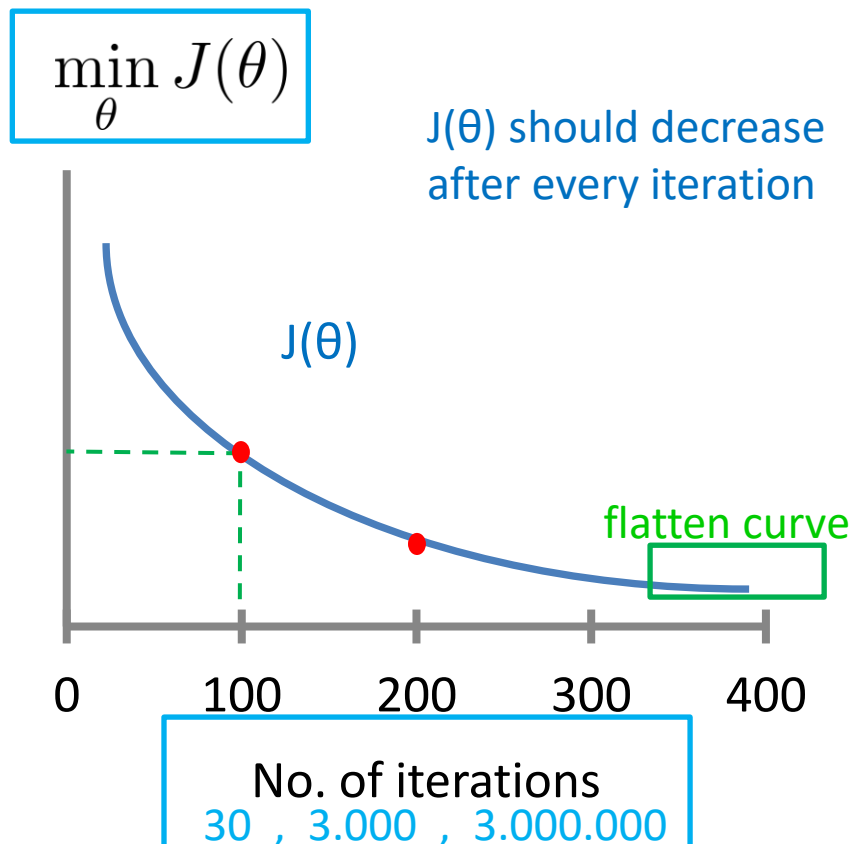
Gradient descent in
practice II: Learning rate

Gradient descent

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

- “Debugging”: How to make sure gradient descent is working correctly.
- How to choose learning rate α .

Making sure gradient descent is working correctly.

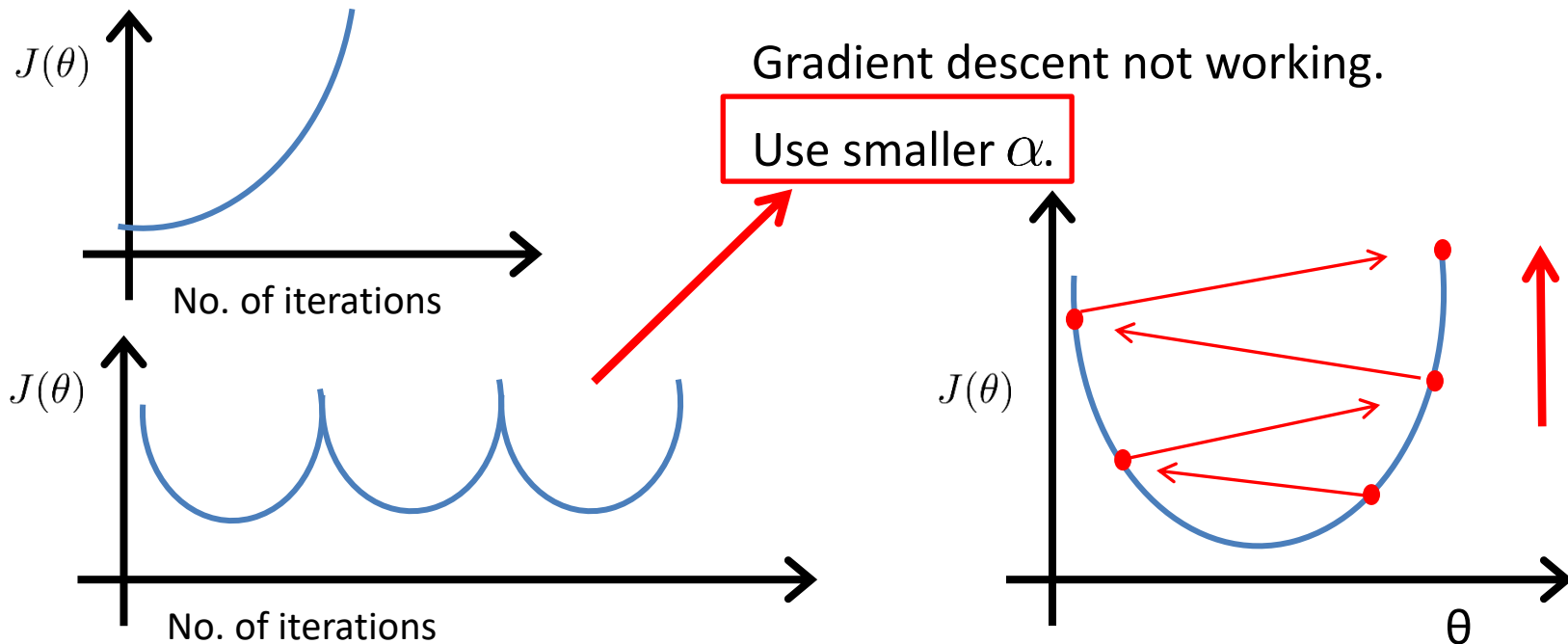


Example automatic convergence test:

Declare convergence if $J(\theta)$ decreases by less than 10^{-3} in one iteration.

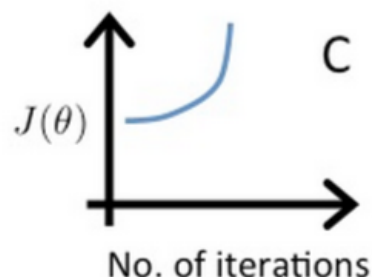
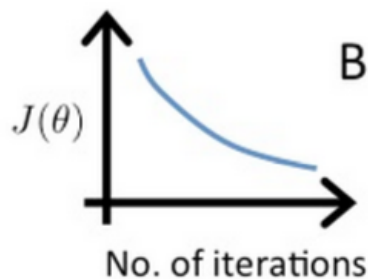
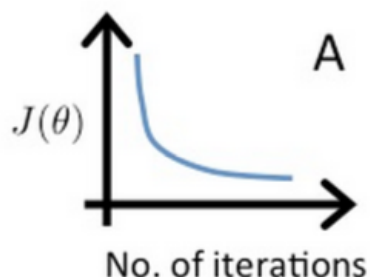
ϵ

Making sure gradient descent is working correctly.



- For sufficiently small α , $J(\theta)$ should decrease on every iteration.
- But if α is too small, gradient descent can be slow to converge.

Suppose a friend ran gradient descent three times, with $\alpha = 0.01$, $\alpha = 0.1$, and $\alpha = 1$, and got the following three plots (labeled A, B, and C):



Which plots corresponds to which values of α ?

- ☐ A is $\alpha = 0.01$, B is $\alpha = 0.1$, C is $\alpha = 1$.
- ☐ A is $\alpha = 0.1$, B is $\alpha = 0.01$, C is $\alpha = 1$.
- ☐ A is $\alpha = 1$, B is $\alpha = 0.01$, C is $\alpha = 0.1$.
- ☐ A is $\alpha = 1$, B is $\alpha = 0.1$, C is $\alpha = 0.01$.

Summary:

- If α is too small: slow convergence.
- If α is too large: $J(\theta)$ may not decrease on every iteration; may not converge.

Slow converge is
also possible

To choose α , try

$\dots, 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, \dots$

Linear regression with
multiple variable

Features and polynomial
regression

Fundamentals of Machine Learning

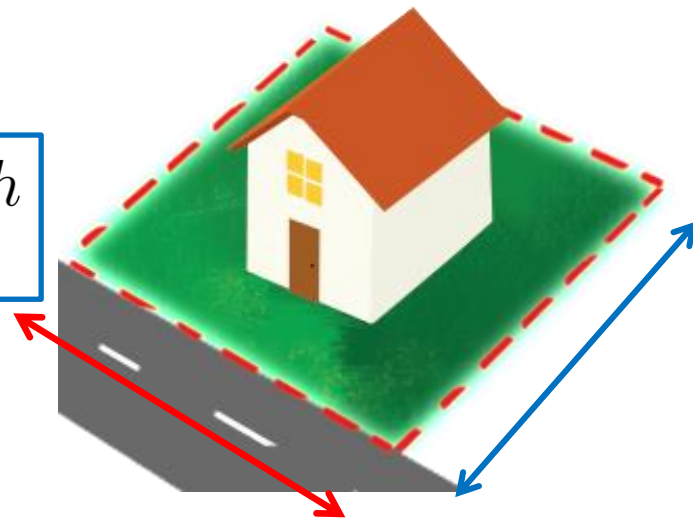
Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

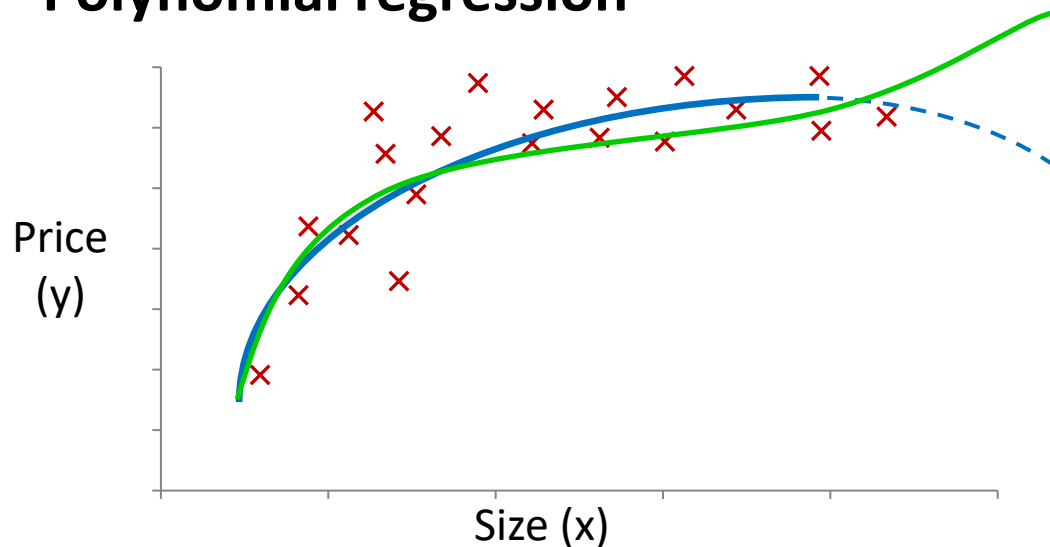
Area (land)

$x = \text{frontage} * \text{depth}$

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



Polynomial regression



$$\theta_0 + \theta_1 x + \theta_2 x^2$$

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3$$

$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1 (\text{size}) + \theta_2 (\text{size})^2 + \theta_3 (\text{size})^3 \end{aligned}$$

$$x_1 = (\text{size})$$

$$x_2 = (\text{size})^2$$

$$x_3 = (\text{size})^3$$

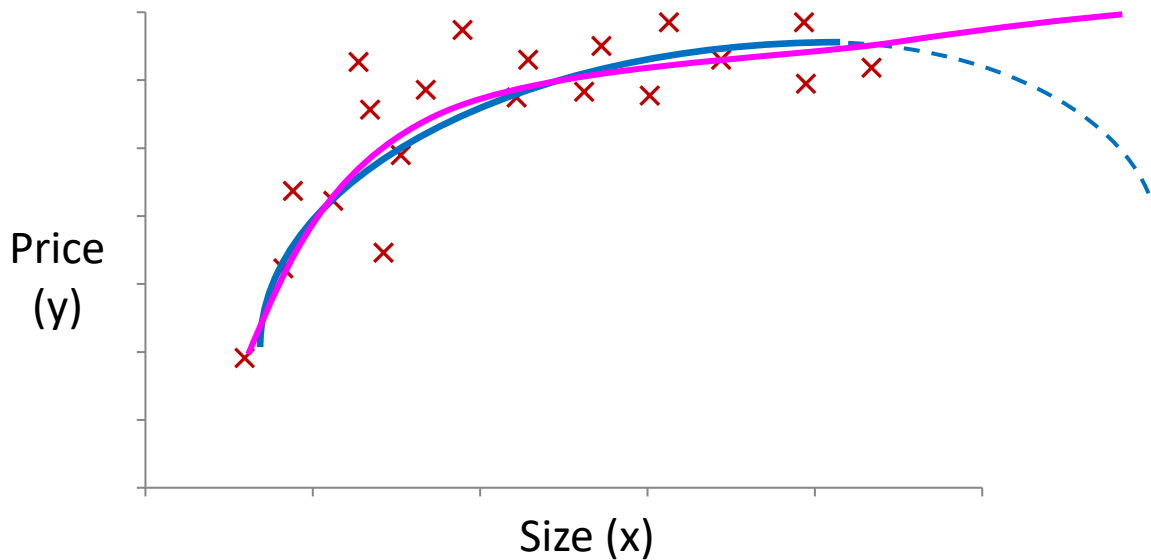
Size : 1 – 1000

Size²: 1 – 1.000.000

Size³: 1 – 10⁹

Feature scaling

Choice of features



$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})}$$

Suppose you want to predict a house's price as a function of its size. Your model is

$$h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2\sqrt{(\text{size})}.$$

Suppose size ranges from 1 to 1000 (feet²). You will implement this by fitting a model

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2.$$

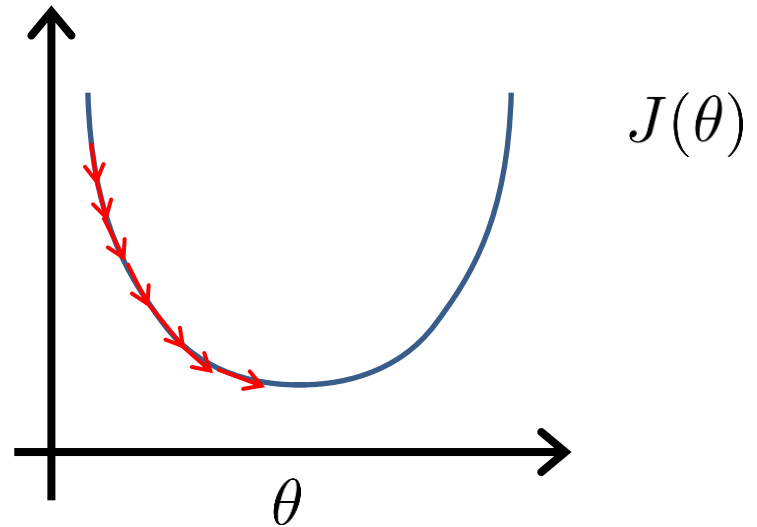
Finally, suppose you want to use feature scaling (without mean normalization). Which of the following choices for x_1 and x_2 should you use? (Note: $\sqrt{1000} \approx 32$.)

- ☐ $x_1 = \text{size}, x_2 = 32\sqrt{(\text{size})}$
- ☐ $x_1 = 32(\text{size}), x_2 = \sqrt{(\text{size})}$
- ☐ $x_1 = \frac{\text{size}}{1000}, x_2 = \frac{\sqrt{(\text{size})}}{32}$
- ☐ $x_1 = \frac{\text{size}}{32}, x_2 = \sqrt{(\text{size})}.$

Linear regression with multiple variable

Normal equation

Gradient Descent



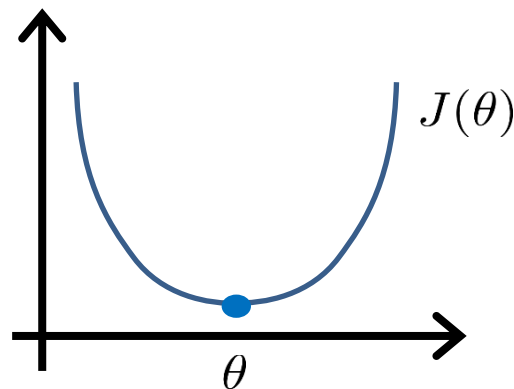
Normal equation: Method to solve for θ analytically.

Intuition: If 1D ($\theta \in \mathbb{R}$)

$$J(\theta) = a\theta^2 + b\theta + c$$

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

Solve for θ



$$\theta \in \mathbb{R}^{n+1} \quad J(\theta_0, \theta_1, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \quad (\text{for every } j)$$

Solve for $\theta_0, \theta_1, \dots, \theta_n$

Examples: $m = 4$.

extra x_0	Size (feet ²) x_1	Number of bedrooms x_2	Number of floors x_3	Age of home (years) x_4	Price (\$1000) y
1	2104	5	1	45	460
1	1416	3	2	40	232
1	1534	3	2	30	315
1	852	2	1	36	178

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix}$$

$m \times (n+1)$

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

m – dimensional vector

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

m **examples** $(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})$; n **features**.

$$x^{(i)} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1}$$

$X =$
(design matrix)

$$\begin{bmatrix} \text{----- } (x^{(1)})^T \text{-----} \\ \text{----- } (x^{(2)})^T \text{-----} \\ \vdots \\ \text{----- } (x^{(m)})^T \text{-----} \end{bmatrix}$$

$m \times (n+1)$ dimensional matrix

E.g. If $x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(n)} \end{bmatrix} \quad m \times 2$$

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

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

Suppose you have the training in the table below:

age (x_1)	height in cm (x_2)	weight in kg (y)
4	89	16
9	124	28
5	103	20

You would like to predict a child's weight as a function of his age and height with the model

$$\text{weight} = \theta_0 + \theta_1 \text{age} + \theta_2 \text{height}.$$

☐ $X = \begin{bmatrix} 1 & 4 & 89 \\ 1 & 9 & 124 \\ 1 & 5 & 103 \end{bmatrix}, y = \begin{bmatrix} 16 \\ 28 \\ 20 \end{bmatrix}$

☐ $X = \begin{bmatrix} 4 & 89 & 1 \\ 9 & 124 & 1 \\ 5 & 103 & 1 \end{bmatrix}, y = \begin{bmatrix} 16 \\ 28 \\ 20 \end{bmatrix}$

☐ $X = \begin{bmatrix} 1 & 4 & 89 \\ 1 & 9 & 124 \\ 1 & 5 & 103 \end{bmatrix}, y = \begin{bmatrix} 1 & 16 \\ 1 & 28 \\ 1 & 20 \end{bmatrix}$

☐ $X = \begin{bmatrix} 1 & 4 & 89 \\ 1 & 9 & 124 \\ 1 & 5 & 103 \end{bmatrix}, y = \begin{bmatrix} 16 \\ 28 \\ 20 \end{bmatrix}$

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

$(X^T X)^{-1}$ is inverse of matrix $X^T X$.

Matlab: `pinv(X' * X) * X' * y` X' – transpose of X (X^T)

Feature scaling is not necessary

$$0 \leq x_1 \leq 1$$

$$0 \leq x_2 \leq 1000$$

$$0 \leq x_1 \leq 10^{-6}$$

....



m training examples, n features.

Gradient Descent

- Need to choose α .
- Needs many iterations.
- Works well even when n is large.

$$n = 10^6$$

Normal Equation

- No need to choose α .
- Don't need to iterate.
- Need to compute $(X^T X)^{-1}$ $n \times n$ $O(n^3)$
- Slow if n is very large.

$$n = 100$$

$$n = 1.000$$

$$n = 10.000$$

Linear regression with
multiple variable

Normal equation
and non-invertibility

Normal equation

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

- What if $X^T X$ is non-invertible? (singular/degenerate)
- Octave: `pinv(X' * X) * X' * y`

pinv – pseudo-inverse

inv – inverse

What if $X^T X$ is non-invertible?

- Redundant features (linearly dependent).

E.g. x_1 = size in feet²

x_2 = size in m²

1m = 3.28 feet

$x_1 = (3.28) x_2$

delete one of the feature

- Too many features (e.g. $m \leq n$).

- Delete some features, or use regularization.

$m=10$

$n=100$