# Linear Regression with multiple variables

## Multiple features

**Multiple features (variables).**

| Size (feet²) $x_1$ | Number of bedrooms $x_2$ | Number of floors $x_3$ | Age of home (years) $x_4$ | Price ($1000) $y$ |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| ... | ... | ... | ... | ... |

$m = 47$

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

Notation:

$n$ = number of features    $n = 4$

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

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

$x_3^{(2)} = 2$

Multivariate linear regression:

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

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \qquad \theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\begin{bmatrix} \theta_0 & \theta_1 & \cdots & \theta_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$\theta^T$ , $(n+1) \times 1$ matrix , $\theta^T x$

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n$$

$$= \boxed{\theta^T x.}$$

## Gradient descent for multiple variables

Hypothesis: $h_\theta(x) = \theta^T x = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$    $x_0 = 1$

Parameters: $\theta_0, \theta_1, \ldots, \theta_n$    $\theta$    $n+1$ - dimensional vector

Cost function:
$$J(\theta_0, \theta_1, \ldots, \theta_n) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)})^2$$
$J(\theta)$

Gradient descent:

Repeat {
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_n) \; J(\theta)$$
}    (simultaneously update for every $j = 0, \ldots, n$)

**Gradient Descent**

Previously (n=1):

Repeat {

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

$$\underbrace{\qquad}_{\frac{\partial}{\partial \theta_0} J(\theta)}$$

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

(simultaneously update $\theta_0, \theta_1$)

}

New algorithm $(n \geq 1)$:

Repeat {

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

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

(simultaneously update $\theta_j$ for $j = 0, \ldots, n$)

$$x_0^{(i)} = 1$$

}

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

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

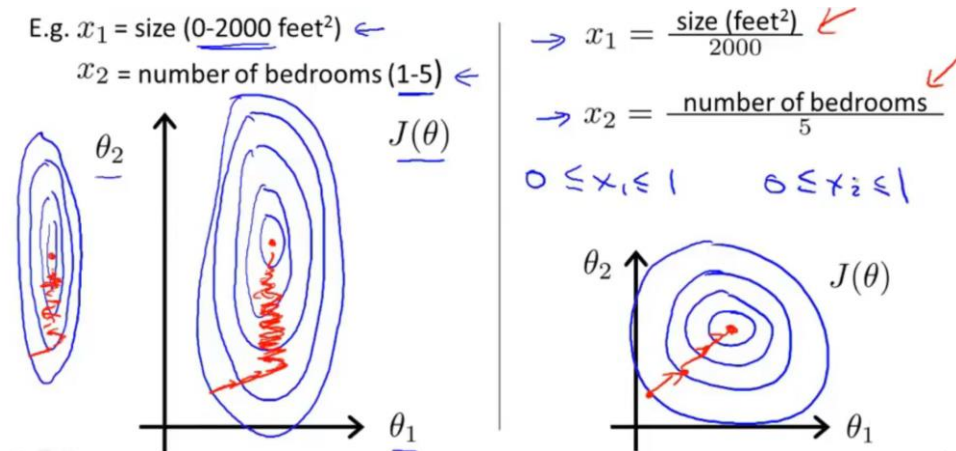$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^{m} (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

...

Andrew Ng

# Gradient descent in practice I: Feature Scaling

Feature Scaling

Idea: Make sure features are on a similar scale.(Then gradient descents can converge more quickly)

E.g. $x_1$ = size (0-2000 feet²)

$x_2$ = number of bedrooms (1-5)



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

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

$$0 \leq x_1 \leq 1 \qquad 0 \leq x_2 \leq 1$$



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

$$x_0 = 1$$

$$0 \leq x_1 \leq 3 \quad \checkmark$$

$$-2 \leq x_2 \leq 0.5 \quad \checkmark$$

$$-100 \leq x_3 \leq \boxed{100} \quad ✗$$

$$-0.0001 \leq x_4 \leq \boxed{0.0001} \quad ✗$$

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

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

## 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. $\rightarrow$ $x_1 = \dfrac{size - 1000}{2000}$      Average $size = 100$

$x_2 = \dfrac{\#bedrooms - 2}{5}$      $1\text{-}5$ bedrooms
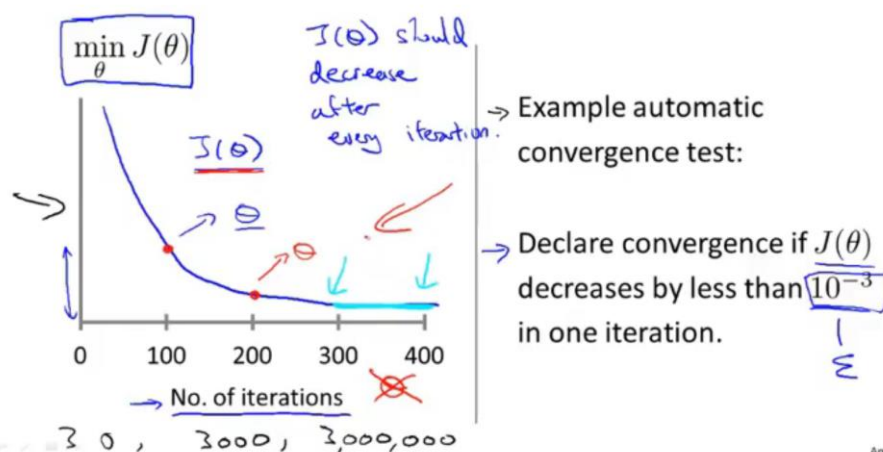
$-0.5 \le x_1 \le 0.5$   $-0.5 \le x_2 \le 0.5$

$x_1 \leftarrow$    $\dfrac{x_1 - \mu_1}{S_1} \leftarrow$ avg value of $x_1$ in training set

$\qquad$ range (max - min)
$\qquad$ (or standard deviation)

$x_2 \leftarrow \dfrac{x_2 - \mu_1}{S_2}$

# Gradient descent in practice II: Learning rate

## Making sure gradient descent is working correctly.

$\boxed{\min_\theta J(\theta)}$

$J(\theta)$ should decrease after every iteration.

$J(\theta)$

$\theta$
$\theta$

0   100   200   300   400
$\rightarrow$ No. of iterations
$3\,0\,,\quad 3000\,,\quad 3,000,000$

$\rightarrow$ Example automatic convergence test:

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

$\varepsilon$

## Making sure gradient descent is working correctly.

$J(\theta)$

No. of iterations

Gradient descent not working.

$\boxed{\text{Use smaller } \alpha.} \leftarrow$

$J(\theta)$

No. of iterations

$J(\theta)$

No. of iterations

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

## Summary:

$J(\theta)$ graph with arrow pointing down, #iters axis

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

To choose $\alpha$, try

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

$3\times$  $\sim 3\times$  $3\times$  $\sim 3\times$

Andrew Ng

# Feature and polynomial regression

define new features: x = frontage * depth.

## Housing prices prediction

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

$x_1$  $x_2$

Area

$x = \underline{frontage * depth}$

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

$\curvearrowright$ land area

## Polynomial regression

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

$$\to \boxed{\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3}$$

Price (y) vs Size (x)

$$h_\theta(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$
$$= \theta_0 + \theta_1 (size) + \theta_2 (size)^2 + \theta_3 (size)^3$$

$\to x_1 = (size)$
$\to x_2 = (size)^2$
$\to x_3 = (size)^3$

Size: $1 - 1000$
Size$^2$: $1 - 1000,000$
Size$^3$: $1 - 10^9$

Andrew Ng

**Choice of features**



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

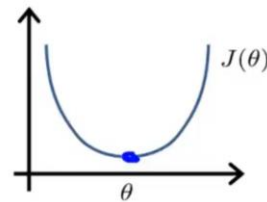$$\boxed{h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2\sqrt{(size)}}$$

# Normal equation

Normal equation: Method to solve for theta analytically.

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

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

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

Solve for $\theta$



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

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

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

**Examples:** $m = 4$.

| $x_0$ | Size (feet²) $x_1$ | 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 |

$$X = \begin{bmatrix} 1 & 2104 & 5 & 1 & 45 \\ 1 & 1416 & 3 & 2 & 40 \\ 1 & 1534 & 3 & 2 & 30 \\ 1 & 852 & 2 & 1 & 36 \end{bmatrix} \qquad y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 178 \end{bmatrix}$$

$m \times (n+1)$ $\qquad\qquad\qquad$ m-dimensional vector

$$\boxed{\theta = (X^TX)^{-1}X^Ty}$$

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

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

$m \times (n+1)$

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

$$X = \begin{bmatrix} 1 & x_1^{(i)} \\ 1 & x_2^{(i)} \\ \vdots & \vdots \\ 1 & x_m^{(i)} \end{bmatrix}$$

$m \times 2$

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

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

Andrew Ng

Feature scaling isn't actually necessary in normal equation.

$$\theta = (X^TX)^{-1}X^Ty \quad \leftarrow$$

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

Set $A = X^TX$

$(X^TX)^{-1} = A^{-1}$

Octave: **pinv(X' *X) *X' *y**

$pinv(X^T * X) * X^T * y$

$\theta = (X^TX)^{-1}X^Ty$    $\min_\theta J(\theta)$

$X'$    $X^T$

~~Feature Scaling~~

$0 \le x_1 \le 1$
$0 \le x_2 \le 1000$
$0 \le x_3 \le 10^{-5}$ ✓

$m$ **training examples,** $n$ **features.**

|  Gradient Descent  |  Normal Equation  |
| --- | --- |
| • Need to choose $\alpha$. | • No need to choose $\alpha$. |
| • Needs many iterations. | • Don't need to iterate. |
| • Works well even when $n$ is large. | • Need to compute $(X^TX)^{-1}$  $n \times n$  $O(n^3)$ |
|  | • Slow if $n$ is very large. |

$n = 10^6$

$n = 100$
$n = 1000$
$n = 10000$

## Normal equation and non-invertibility
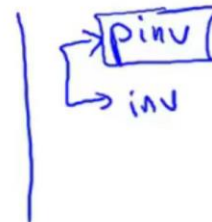
Normal equation

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

$X^T X$

- What if $X^T X$ is non-invertible? (singular/ degenerate)

- Octave: `pinv(X'*X)*X'*y`

pinv → inv

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

- Redundant features (linearly dependent).
  E.g. $x_1$ = size in feet$^2$
  ~~$x_2$ = size in m$^2$~~

  $x_1 = (3.28)^2 x_2$

  $1m = 3.28$ feet

  → $m = 10$ ←
  → $n = 100$ ←
  $\theta \in \mathbb{R}^{101}$

- Too many features (e.g. $m \le n$).
  - Delete some features, or use regularization.

  ↓ later