# Linear Regression with multiple variables

## Multiple features

Machine Learning

# Multiple features (variables).

| Size (feet²) | Price ($1000) |
|:---:|:---:|
| $x$ | $y$ |
| 2104 | 460 |
| 1416 | 232 |
| 1534 | 315 |
| 852 | 178 |
| ... | ... |

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

# Multiple features (variables).

| Size (feet²) | 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 |
| ... | ... | ... | ... | ... |

$$m = 47$$

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

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

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.

Hypothesis:

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

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

e.g. $h_\theta(x) = 80 + 0.1 x_1 + 0.01 x_2 + 3 x_3 - 2 x_4$

age

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

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

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

$$\underbrace{[\theta_0 \ \theta_1 \cdots \theta_n]}_{\theta^T} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

$(n+1) \times 1$ matrix

$\theta^T x$

$x$

$$h_\theta(x) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n \qquad \overset{=1}{\downarrow}$$

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

Multivariate linear regression. $\leftarrow$

# Linear Regression with multiple variables

# Gradient descent for multiple variables

Machine Learning

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

$\longrightarrow x_0 = 1$

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 {

$\longrightarrow \quad \theta_j := \theta_j - \alpha \dfrac{\partial}{\partial \theta_j} J(\theta_0, \ldots, \theta_n) \quad J(\theta)$

}

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

# Gradient Descent

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 $\theta_j$ for $j = 0, \ldots, n$)

}

Previously (n=1):

Repeat {

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

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

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

(simultaneously update $\theta_0, \theta_1$)

}

$x_0^{(i)} = 1$

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

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

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

$\ldots$

Andrew Ng

# Linear Regression with multiple variables

## Gradient descent in practice I: Feature Scaling

If you have a problem where you have multiple features, if you make sure that the features are on a similar scale, by which I mean make sure that the different features take on similar ranges of values, then gradient descents can converge more quickly.
我: 應該是因為alpha對所有dimenstion都是一樣的值.

Machine Learning

# Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. $x_1$ = size (0-2000 feet²)

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

$J(\theta)$



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

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

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

$J(\theta)$

# Feature Scaling

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

$$x_0 = 1$$

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

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

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

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

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

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

don't worry if your features are not exactly on the same scale or exactly in the same range of values. But so long as they're all close enough to this gradient descent it should work okay.

# Mean normalization

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

for those of you that are being super careful technically if we're taking the range as max minus min this 5 here will actually become a 4.

Avenge       size = 100

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

$x_2 = \frac{\#bedrooms - 2}{5 \quad 4}$

1-5  bedrooms

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

$x_1 \leftarrow \frac{x_1 - \mu_1}{S_1}$

$\leftarrow$ avg value of $x_1$ in traning set

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

range (max - min)
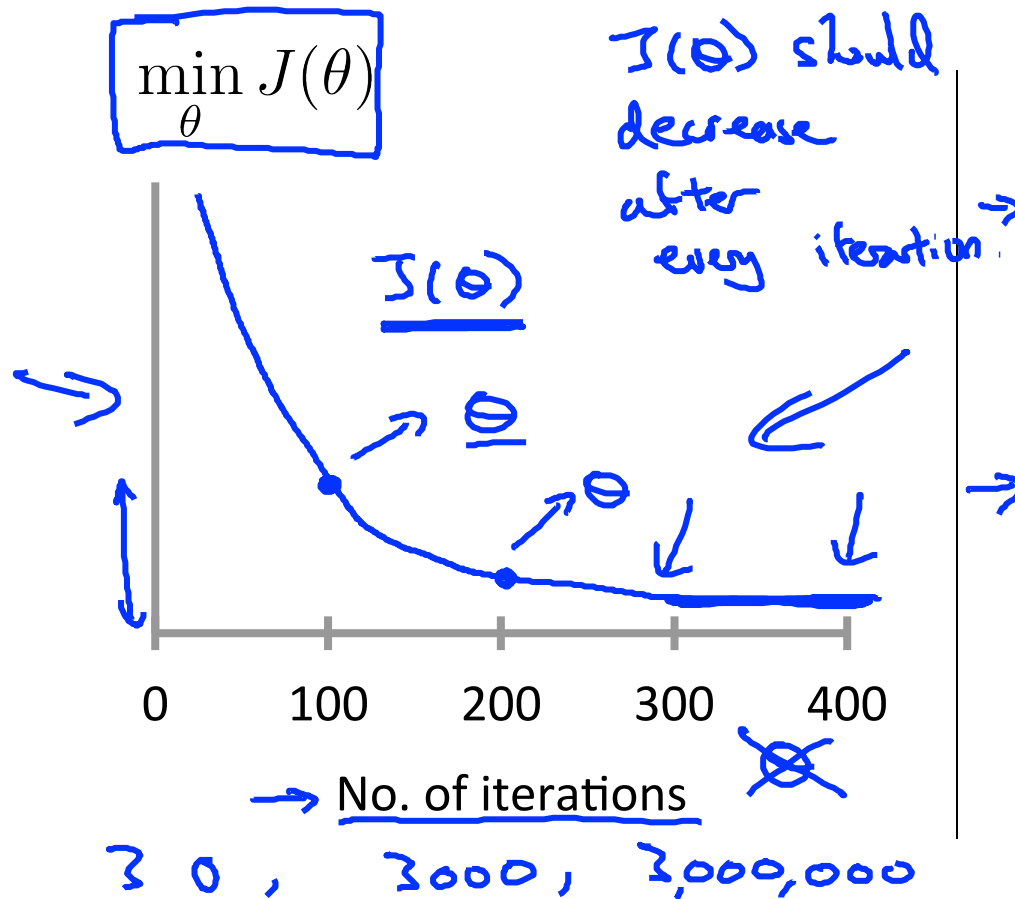(or standard deviation)

# Linear Regression with multiple variables

## Gradient descent in practice II: Learning rate

Machine Learning

**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 $\alpha.$

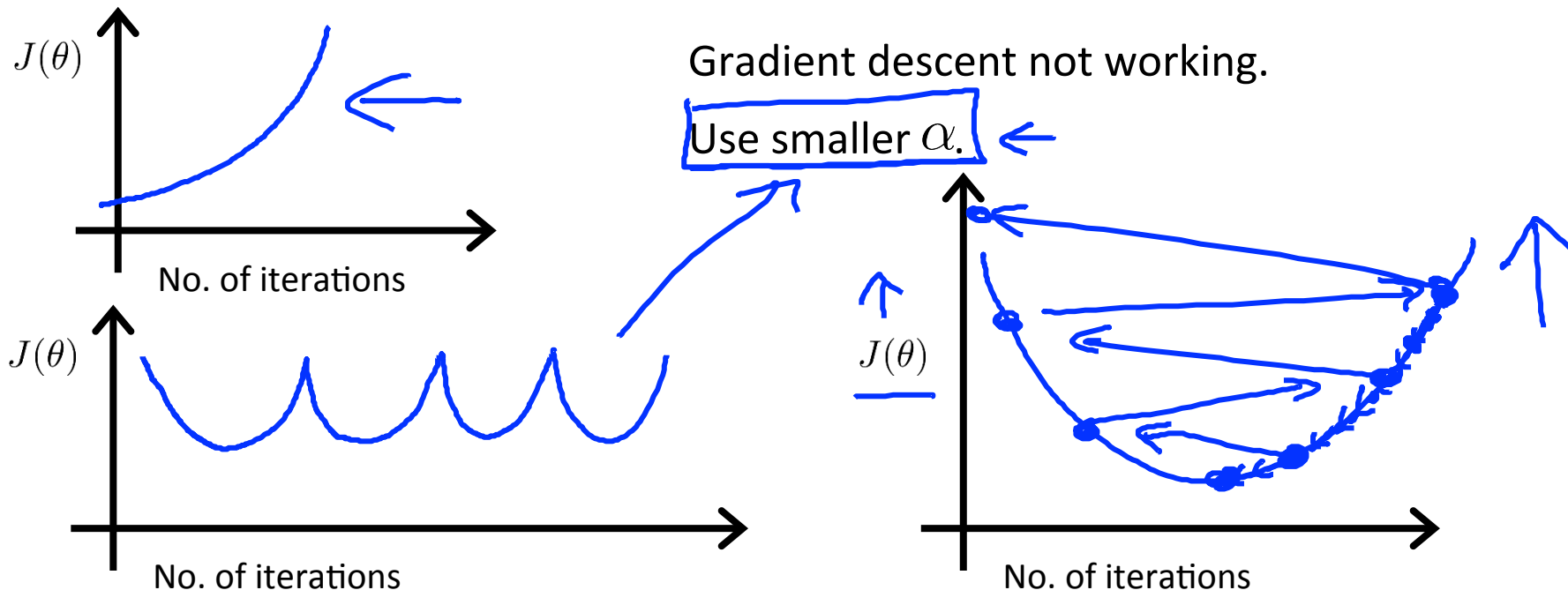**Making sure gradient descent is working correctly.**



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

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

$J(\theta)$

No. of iterations

30, 3000, 3,000,000

→ Example automatic convergence test:

→ 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.

Use smaller $\alpha$.

$J(\theta)$

No. of iterations

$J(\theta)$

No. of iterations

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

Andrew Ng

## Summary:

- 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)

$J(\theta)$

#iters

To choose $\alpha$, try  *these are factor of ten differences*

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

3x  ≈3x  3x  ≈3x

Linear Regression with multiple variables

Features and polynomial regression

Machine Learning

# Housing prices prediction

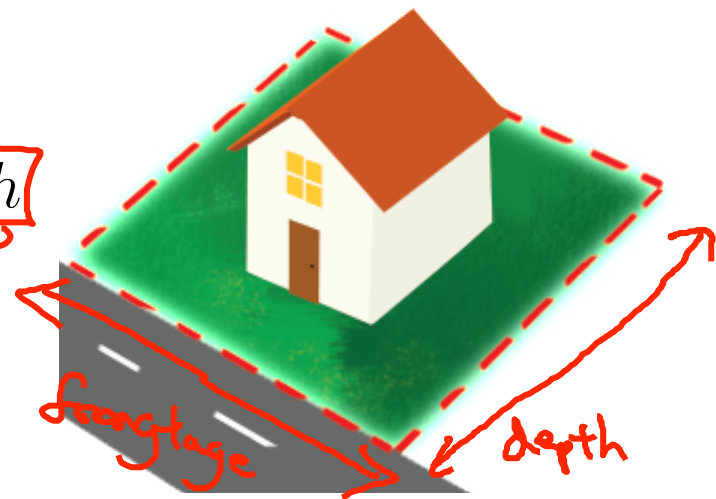$$h_\theta(x) = \theta_0 + \theta_1 \times \boxed{frontage} + \theta_2 \times \boxed{depth}$$

$x_1$

$x_2$

What you can do is actually create new features by yourself:

Area

$x = frontage * depth$

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

← land area

frontage
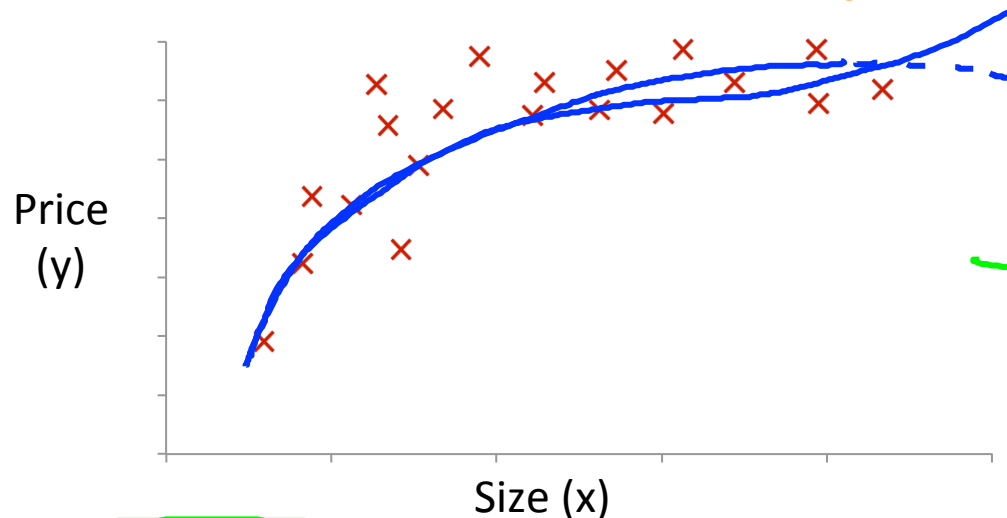
depth

Closely related to the idea of choosing your features is this idea called polynomial regression.
the natural way to do that is to set the first feature x one to be the size of the house,
and set the second feature x two to be the square of the size of the house...

# Polynomial regression

Price
(y)

Size (x)

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

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

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

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

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

Andrew Ng

# Choice of features



Price (y)

Size (x)

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

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

Linear Regression with multiple variables

Normal equation

Machine Learning

# Gradient Descent



$J(\theta)$

$\theta$

Normal equation: Method to solve for $\theta$ analytically.

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

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

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

Solve for $\theta$



$J(\theta)$

$\theta$

---

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

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

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

Examples: $m = 4$.

| $x_0$ | 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 |

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

M × (n+1)

m - dimensional vector

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

由本頁知, 後面講的還是線性回歸, 而不是高階的.

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

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

$X \\ \text{(design matrix)}$

$$= \begin{bmatrix} \rule{1cm}{0.4pt} (x^{(1)})^T \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} (x^{(2)})^T \rule{1cm}{0.4pt} \\ \vdots \\ \rule{1cm}{0.4pt} (x^{(m)})^T \rule{1cm}{0.4pt} \end{bmatrix}$$

$m \times (n+1)$

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

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

$m \times 2$

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

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

Andrew Ng

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

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

Set $\underset{\underbrace{\quad}}{A} = \underset{\underbrace{\quad}}{X^T}\underset{\underbrace{\quad}}{X}$

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

Octave: **pinv $\boxed{\textbf{(X' *X)}}$ *X' *y**

$$\underbrace{pinv(X^T * X) * X^T * y}$$

$$\theta = (X^TX)^{-1}X^Ty \qquad \underset{\theta}{min} \ J(\theta)$$

$X'$ $\qquad$ $X^T$

~~Feature Scaling~~

$0 \le x_1 \le 1$

$0 \le x_2 \le 1000$

$0 \le x_3 \le 10^{-5}$ $\checkmark$

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

## Gradient Descent

- Need to choose $\alpha$.
- Needs many iterations.
- Works well even when $n$ is large.

n: number of features

$n = 10^6$

## Normal Equation

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

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

Andrew Ng

Linear Regression with multiple variables

Normal equation and non-invertibility (optional)

Machine Learning

# 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`

$$pinv$$

$$inv$$
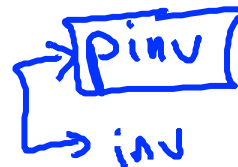
Octave hast two functions for inverting matrices.

One is called pinv, and the other is called inv.

One's called the pseudo-inverse, one's called the inverse.

But you can show mathematically that so long as you use the pinv function then this will actually compute the value of data that you want even if X transpose X is non-invertible.

Andrew Ng

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

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

  $\rightarrow m = 10 \leftarrow$

  $\rightarrow n = 100 \leftarrow$

  $\Theta \in \mathbb{R}^{101}$

  - Delete some features, or use regularization.

    $\downarrow$ later

Andrew Ng