Machine Learning

# Linear Regression with multiple variables

## Multiple features

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

$\qquad\qquad\qquad\qquad\qquad \rightarrow \qquad\quad \rightarrow \quad \rightarrow$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad age$

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

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

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

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

$$= \underbrace{\begin{bmatrix} \theta_0 & \theta_1 & \cdots & \theta_n \end{bmatrix}}_{\theta^T \; 1 \times (n+1) \atop matrix} \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \; = \; \theta^T x$$

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

Multivariate linear regression.

Machine Learning

Linear Regression with multiple variables

Gradient descent for multiple variables

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

$\hookrightarrow x_0 = 1$

**Parameters:** $\theta_0, \theta_1, \ldots, \theta_n$  $\theta$

$n+1$ - dimensional vector

**Cost function:**

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

$J(\theta)$

**Gradient descent:**

Repeat {

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

}

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

# Gradient Descent

**Previously ($n=1$):**

Repeat $\{$

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

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

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

$x_1$

(simultaneously update $\theta_0, \theta_1$)

$\}$

---

**New algorithm ($n \geq 1$):**

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

Repeat $\{$

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

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

$\}$

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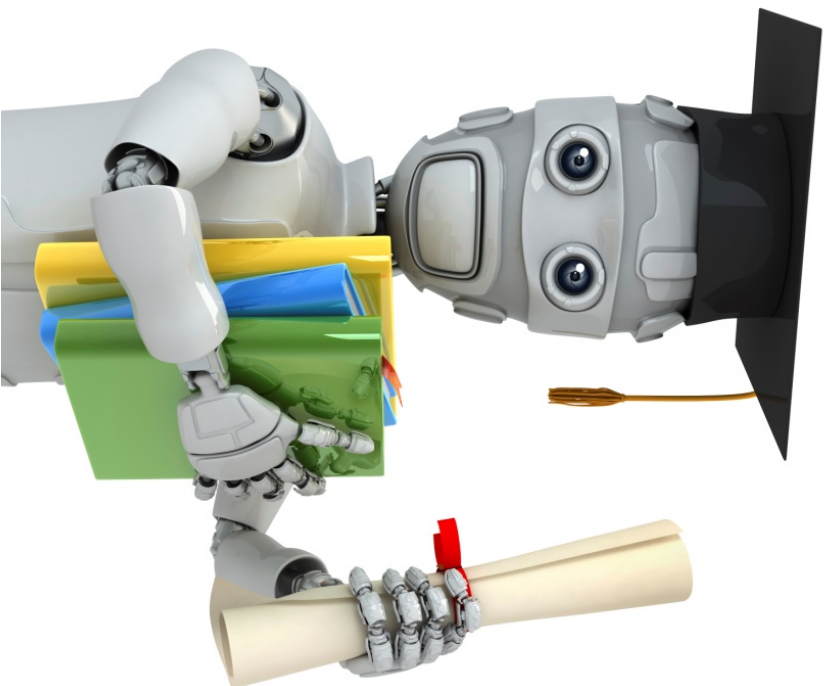
$x_0^{(i)} = 1$

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

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

$\ldots$

Machine Learning

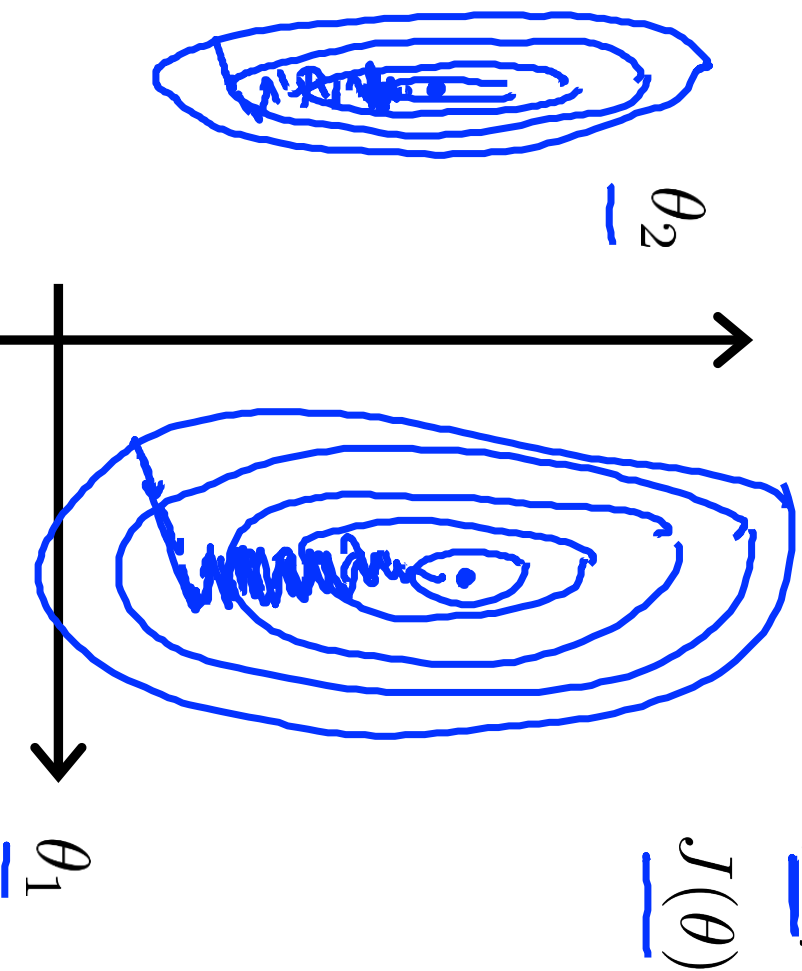# Linear Regression with multiple variables

## Gradient descent in practice I: Feature Scaling

# Feature Scaling

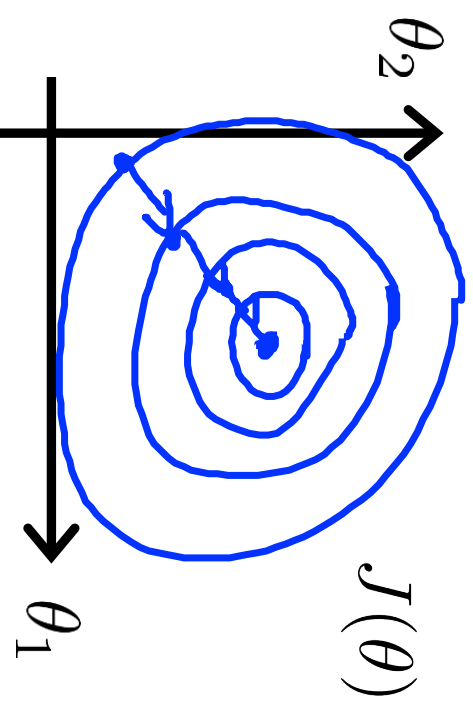Idea: Make sure features are on a similar scale.

E.g. $x_1$ = size (0-2000 feet²)

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

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

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

$$0 \leq x_1 \leq 1$$

$$0 \leq x_2 \leq 1$$

# Feature Scaling

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

$x_0 = 1$

$0 \le x_1 \le 3$

$-2 \le x_2 \le 0.5$

$\boxed{-100 \le x_3 \le 100}$ ✗

$\boxed{-0.0001 \le x_4 \le 0.0001}$ ✗

$-3 \text{ to } 3$

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

# 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.

$$x_1 = \frac{size - 1000}{2000}$$

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

$$-0.5 \le x_1 \le 0.5 \qquad -0.5 \le x_2 \le 0.5$$

Average size $\approx 1000$

$1\text{-}5$ bedrooms

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

$\mu_1$ — avg value of $x_1$ in training set

$S_1$ — range (max - min)
(or standard deviation)

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

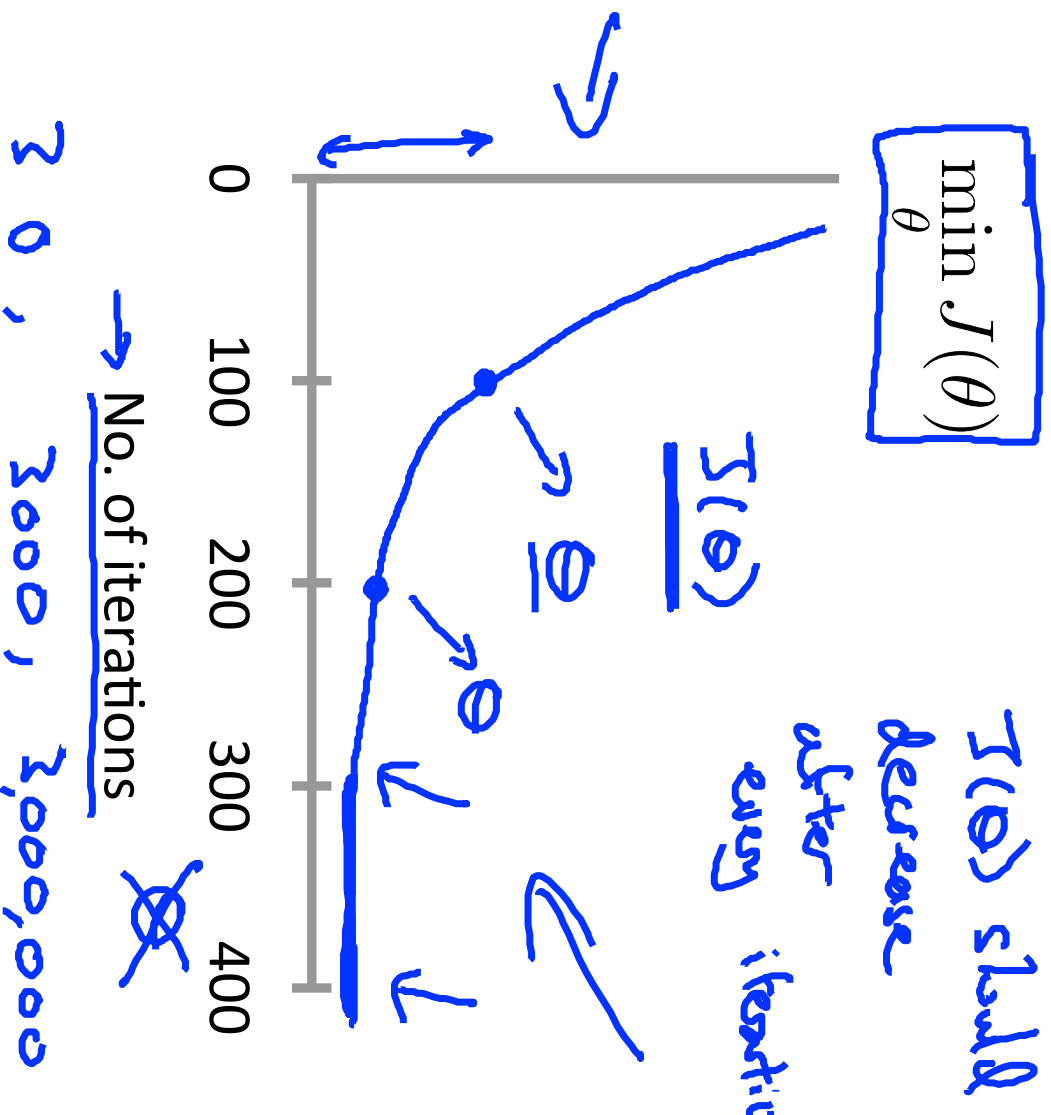# Linear Regression with multiple variables

## 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 $\alpha$.

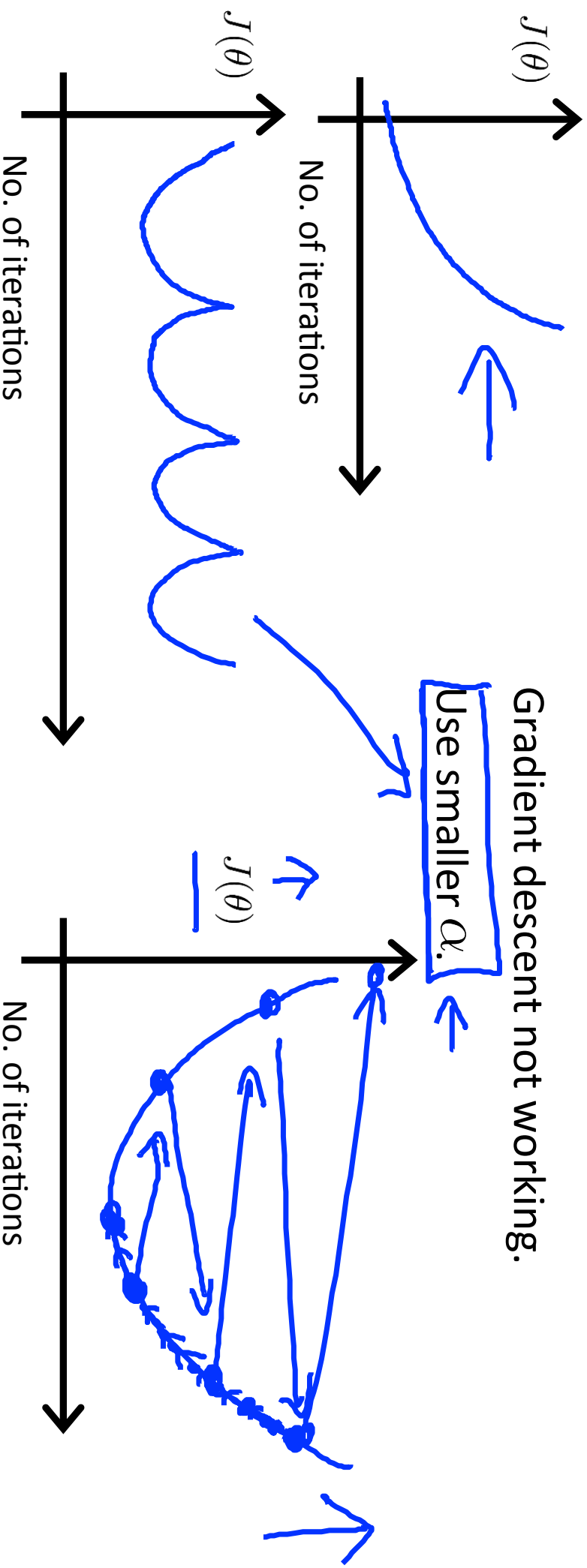# Making sure gradient descent is working correctly.

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

$J(\theta)$

$\theta$

$\theta$

0     100     200     300     400

No. of iterations

3 0 , 3000 , 3,000,000

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

Example automatic
convergence test:

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

$\varepsilon$

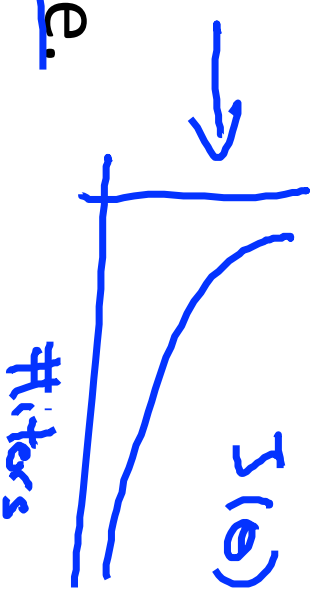# Making sure gradient descent is working correctly.

$J(\theta)$

No. of iterations

Gradient descent not working.

Use smaller $\alpha$.

$J(\theta)$

No. of iterations

$J(\theta)$

No. of iterations

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

# 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)
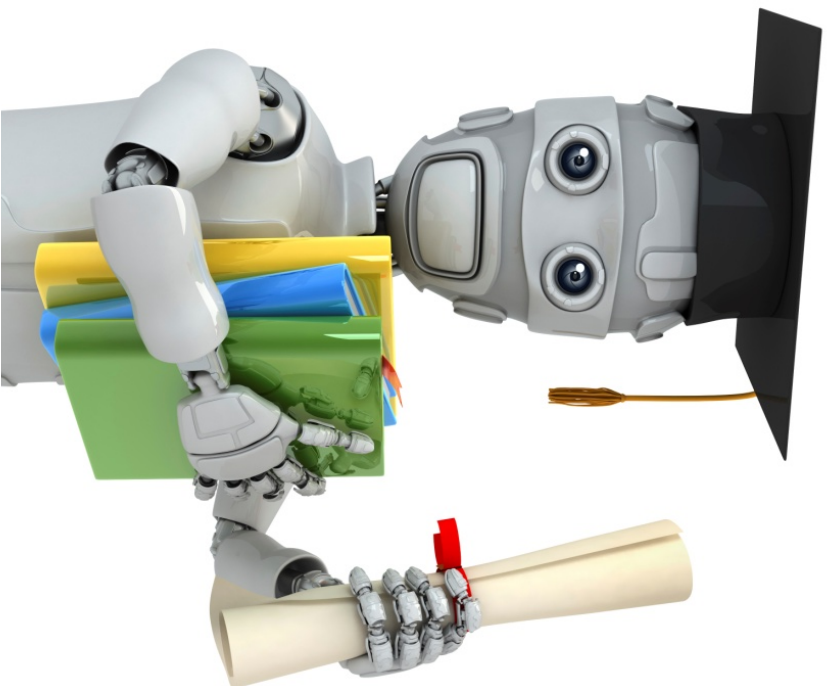
$J(\theta)$

#iters

To choose $\alpha$, try

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

$\times 3 \quad \times 3 \quad \times 3 \quad \times 3 \quad \times 3 \quad \times 3$

Machine Learning

# Linear Regression with multiple variables

## Features and polynomial regression

# Housing prices prediction

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

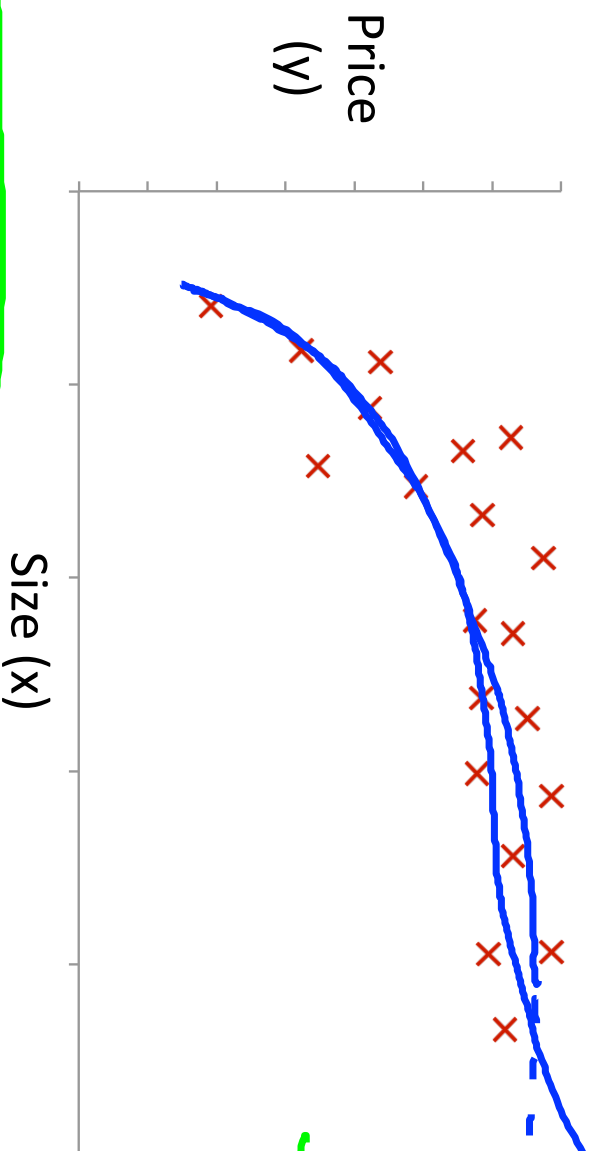$$\underset{x_1}{\phantom{-}} \qquad \underset{x_2}{\phantom{-}}$$

$$\underline{Area}$$

$$x = frontage * depth$$

$$h_\theta(x) = \Theta_0 + \Theta_1 x \quad \curvearrowleft \text{land area}$$



frontage

depth

# Polynomial regression

Price (y)

Size (x)



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

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

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

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

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

# Choice of features

Price
(y)

Size (x)

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

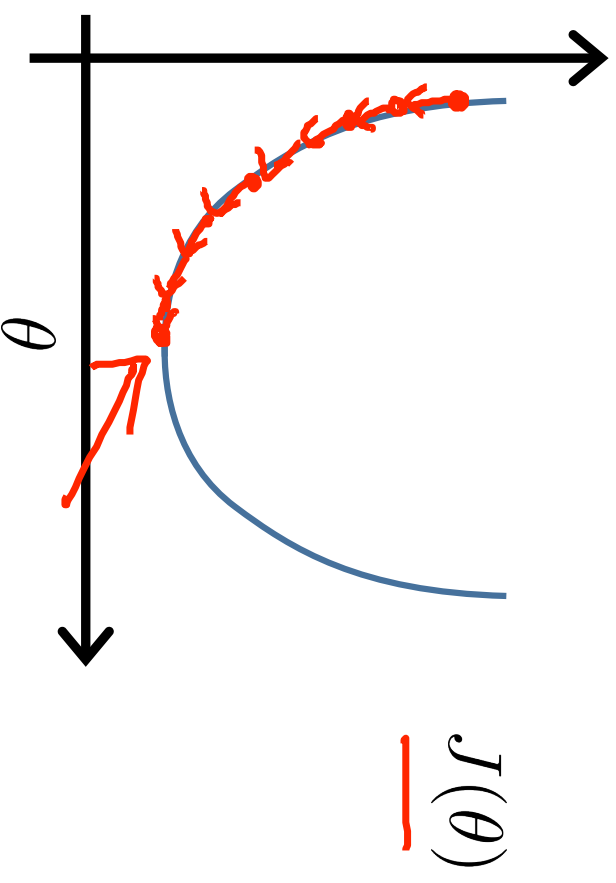$h_\theta(x) = \theta_0 + \theta_1(size) + \theta_2\sqrt{(size)}$

# Linear Regression with multiple variables

## Normal equation

# Gradient Descent



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

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

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

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

Solve for $\theta$



---

$\theta \in \mathbb{R}^{n+1}$

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

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

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 |

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

$M \times (n+1)$

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

m - dimensional vector

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

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

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

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

$m \times (n+1)$

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

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

$$X = \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(3)} \end{bmatrix}$$

$m \times 2$

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

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

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

Set $\underbrace{A: X^T X}$

$\underbrace{(X^T X)^{-1}} = A^{-1}$

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

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

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

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

$X'$  $X^T$  ~~feature scaling~~

$0 \le x_1 \le 3$
$0 \le x_2 \le 100$
$-5 \le x_3 \le 10$

$\underline{m}$ **training examples,** $\underline{n}$ **features.**

## Gradient Descent

- Need to choose $\alpha$.

- Needs many iterations.

- Works well even when $\underline{n}$ is large.

$n = 10^6$

$\longrightarrow$

## Normal Equation

- No need to choose $\alpha$.

- Don't need to iterate.

- Need to compute $\underbrace{(X^T X)^{-1}}_{n \times n}$

- Slow if $\underline{n}$ is very large. $O(n^3)$

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

# Linear Regression with multiple variables

## Normal equation and non-invertibility (optional)

Machine Learning

# Normal equation

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

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

- Octave: **pinv(x'*x)*x'*y**

$X^T X$

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$^2$

  $1m = 3.28$ feet

  $x_1 = (3.28)^2 x_2$

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

  $\rightarrow n = 10$
  $\rightarrow n = 100$
  $\rightarrow m = 10$
  $\rightarrow m = 100$

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

  - Delete some features, or use regularization.

  $\downarrow$
  later