# Linear Regression with multiple variables

## Multiple features

Machine Learning

# Multiple features (variables).

| Size (feet²) $x$ | Price ($1000) $y$ |
|---|---|
| 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) |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |
| … | … | … | … | … |

Notation:

$n$ = number of features

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

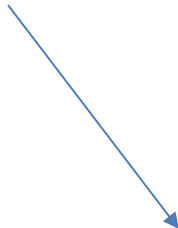$x_j^{(i)}$ = value of feature $j$ in $i^{th}$ training example.

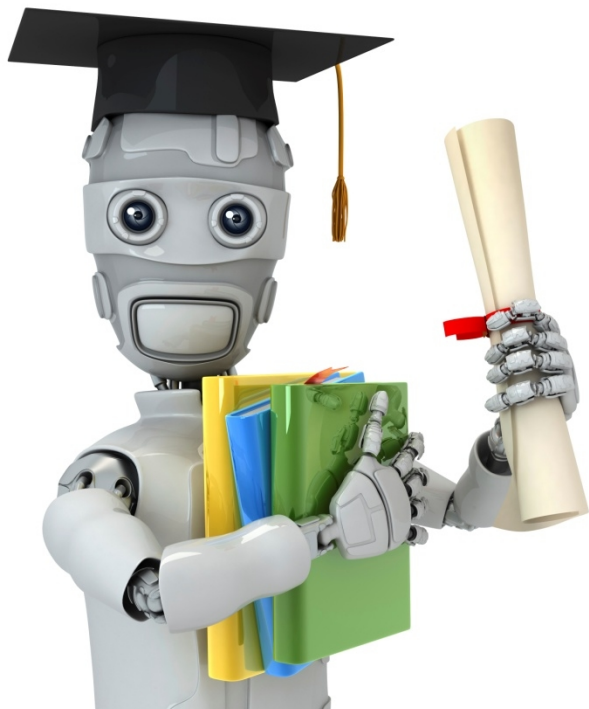Hypothesis:

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

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

For convenience of notation, define $x_0 = 1$.

Multivariate linear regression.

1 * (n+1)

Machine Learning

# Linear Regression with multiple variables

## Gradient descent for multiple variables

Hypothesis:

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

Parameters:

$$\theta_0, \theta_1, \ldots, \theta_n$$

Cost function:

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

Gradient descent:

Repeat {

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

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

# **Gradient Descent**

Previously (n=1):

Repeat {

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

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

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

(simultaneously update $\theta_0, \theta_1$)

}

New algorithm $(n \geq 1)$ :

Repeat {

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

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

}

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

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

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

...

# Linear Regression with multiple variables

## Gradient descent in practice I: Feature Scaling

Machine Learning
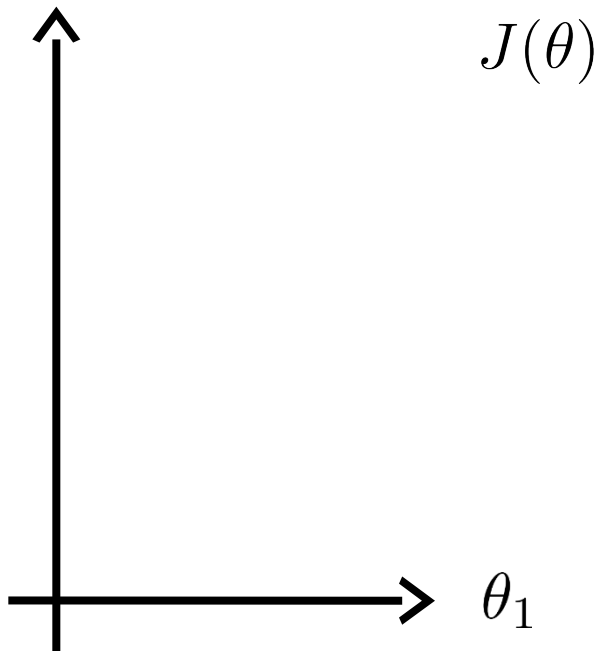
# Feature Scaling

Idea: Make sure features are on a similar scale.

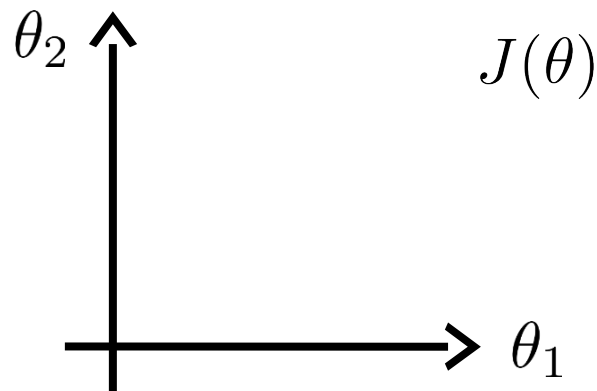E.g. $x_1$ = size (0-2000 feet²)

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

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

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

$\theta_2$

$J(\theta)$

$\theta_1$

$\theta_2$

$J(\theta)$

$\theta_1$

# Feature Scaling

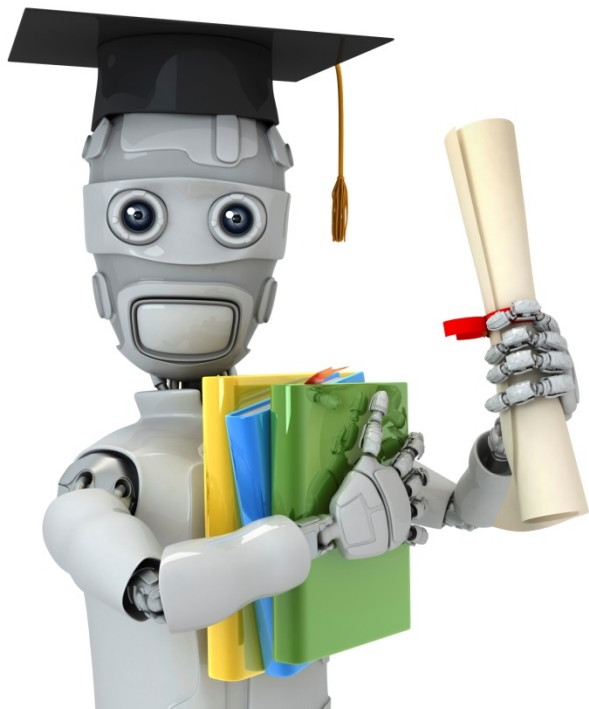Get every feature into approximately $-1 \leq x_i \leq 1$ range.

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

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

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

Machine Learning

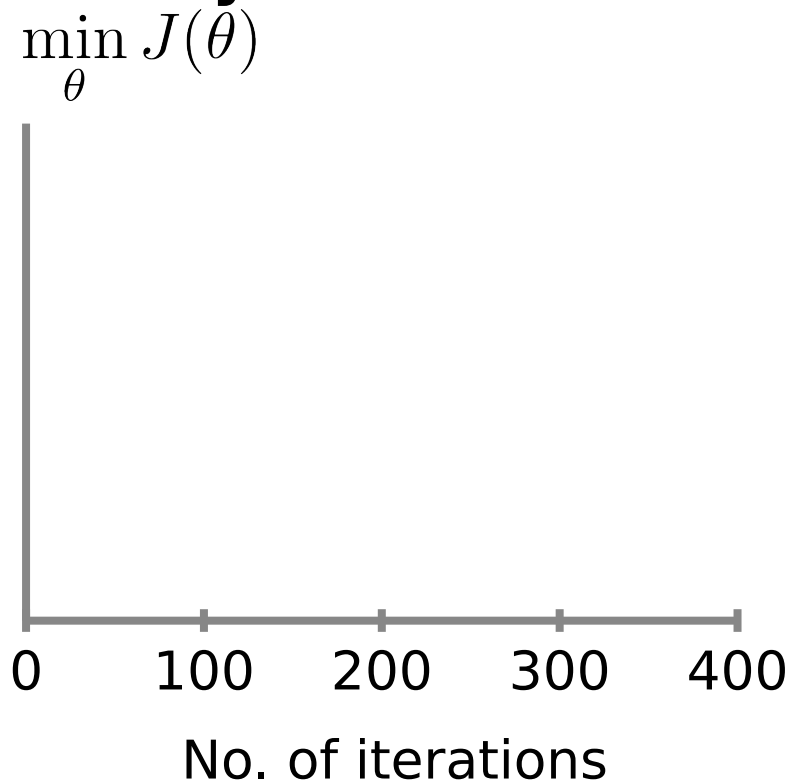# Linear Regression with multiple variables

## 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.
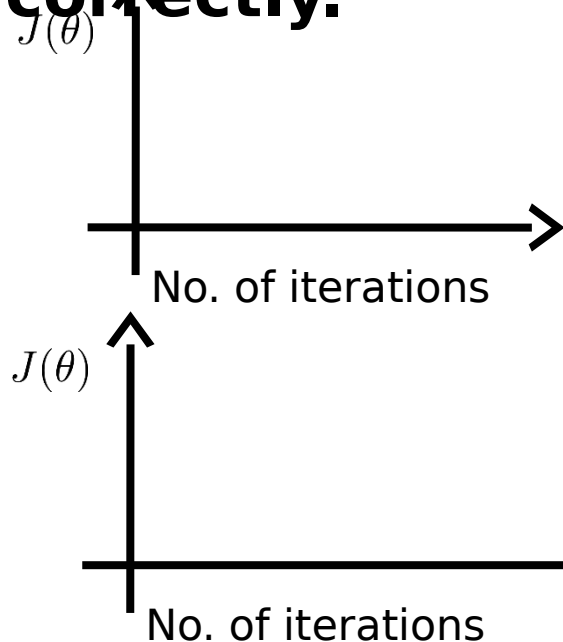
- How to choose learning rate $\alpha$.

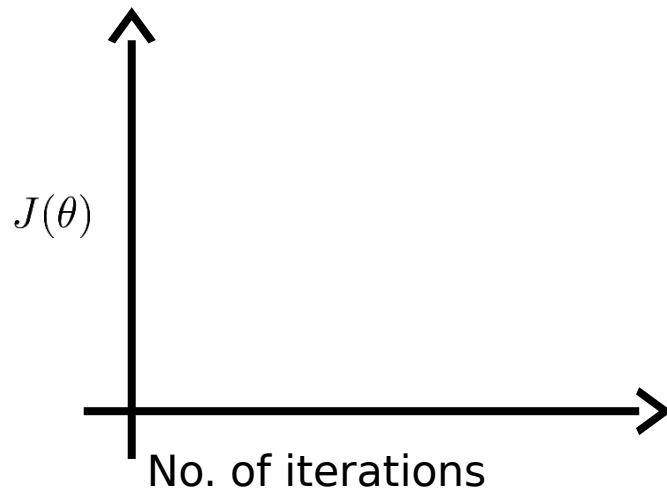# Making sure gradient descent is working correctly.

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

No. of iterations

Example automatic convergence test:

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

# 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.
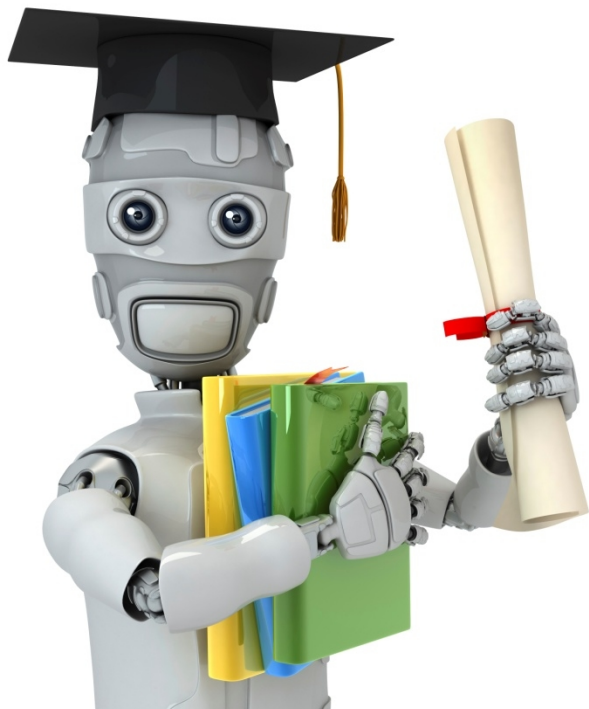- But if $\alpha$ is too small, gradient descent can be slow to

**Summary:**

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

To choose $\alpha$, try

$$\ldots, 0.001, \qquad , 0.01, \qquad , 0.1, \qquad , 1, \ldots$$

Linear Regression with multiple variables

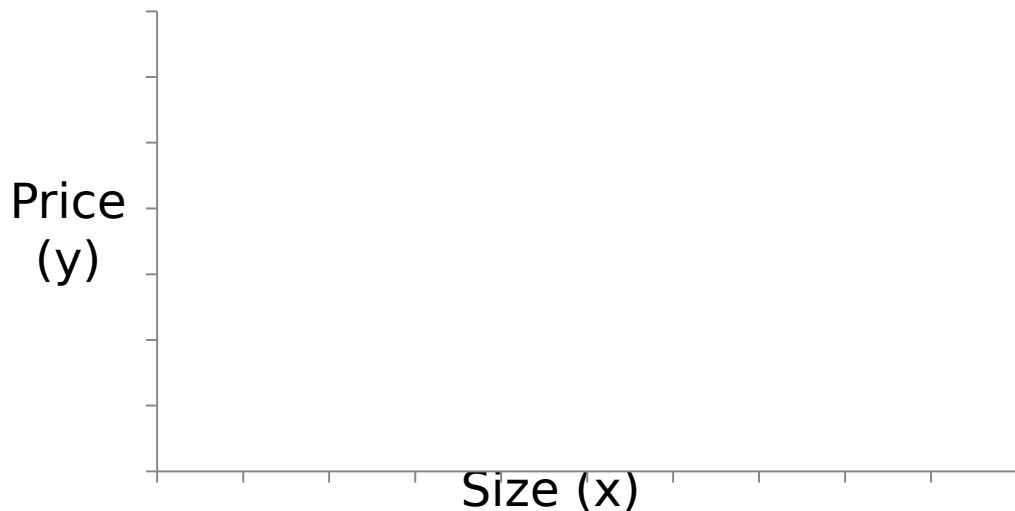Features and polynomial regression

Machine Learning

# Housing prices prediction

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

# Polynomial regression

Price (y) ↑      → Size (x)

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

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

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

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

# Choice of features


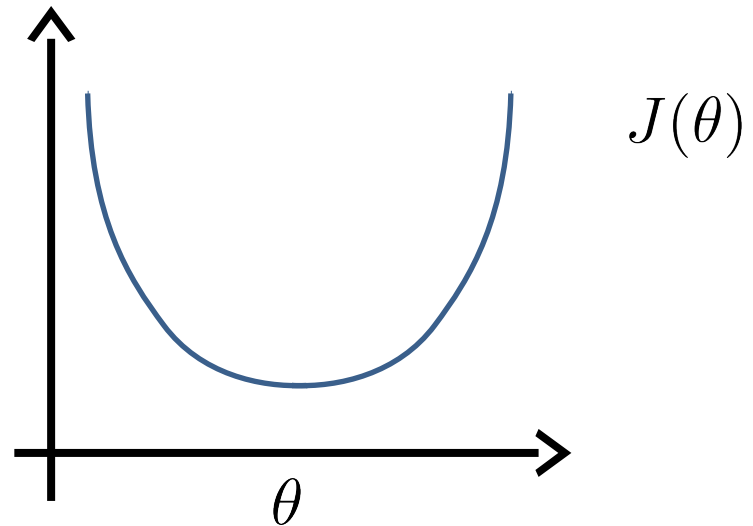
Price (y) / Size (x)

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

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

# Linear Regression with multiple variables

## Normal equation

Machine Learning

# Gradient Descent

$J(\theta)$

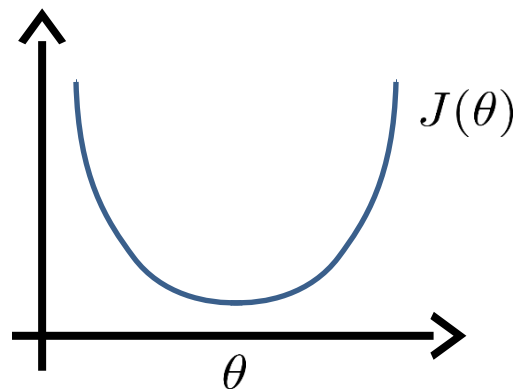$\theta$

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

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

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



$J(\theta)$

$\theta$

---

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

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

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

# Examples: $n = 4$.

| $x_0$ | 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 |

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

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

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

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

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

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

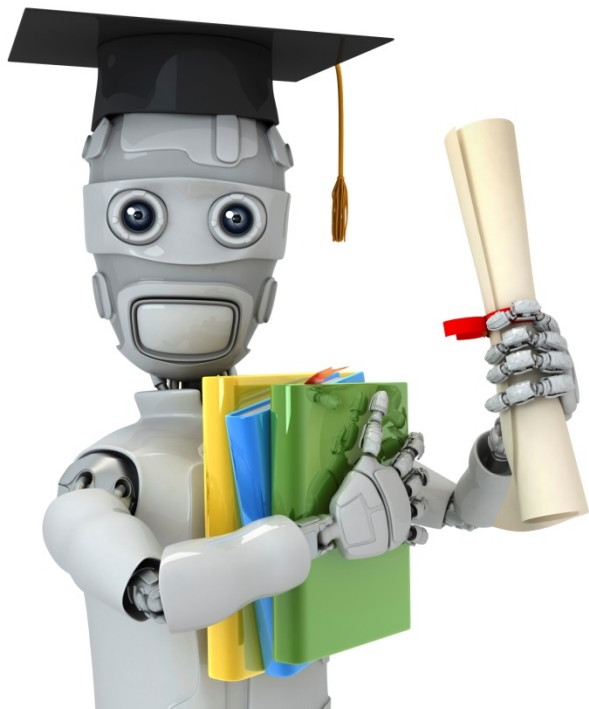$(X^T X)^{-1}$ is inverse of matrix $X^T X$

.

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

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

| Gradient Descent | Normal Equation |
|---|---|

Gradient Descent

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

Normal Equation

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

Machine Learning

Linear Regression with multiple variables

Normal equation and non-invertibility (optional)

# 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**

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

- Redundant features (linearly dependent).
  E.g. $x_1 =$ size in feet²
  $x_2 =$ size in m²

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