# Normal Equation

Computer Parameters Analytically

*Linear Regression with Multiple Variables*

# Gradient Descent



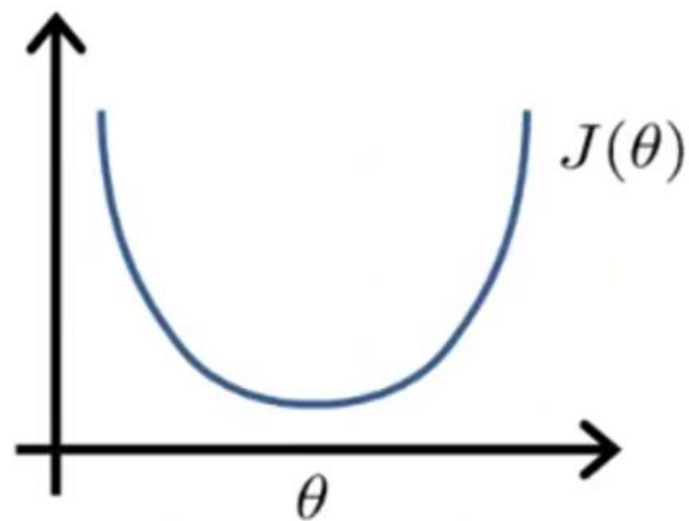$$J(\theta)$$

# Gradient Descent

$$J(\theta)$$

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

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

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

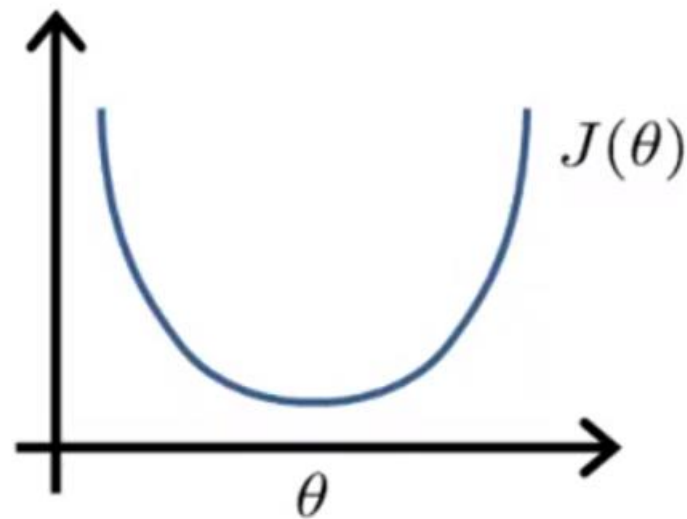How to minimize a function?

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

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

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

$S$



$J(\theta)$

$\theta$

Andrew Ng

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

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

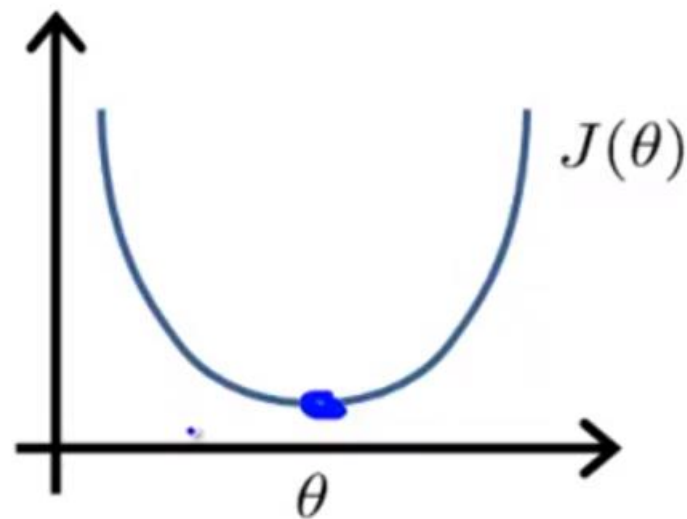$\dfrac{d}{d\theta} J(\theta) = \dots \overset{set}{=} 0$

Solve for $\theta$

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

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

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

Solve for $\theta$


$J(\theta)$

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

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

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

# Examples: $m = 4$.

| 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 |

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 |

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 |

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

# 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 |

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

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 |

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

# 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

Andrew Ng

# 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

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

Andrew Ng

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

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

$m$ **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)

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

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

$$X = \begin{bmatrix} \rule{0.5cm}{0.4pt} \; (x^{(i)})^\top \; \rule{0.5cm}{0.4pt} \\ \\ \\ \end{bmatrix}$$

(design matrix)

$m$ **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)

$$= \begin{bmatrix} \underline{\quad\quad} (x^{(1)})^{\top} \underline{\quad\quad} \\ \underline{\quad\quad} (x^{(2)})^{\top} \underline{\quad\quad} \\ \vdots \end{bmatrix}$$

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

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

$$X \quad \text{(design matrix)} \quad = \begin{bmatrix} \rule{1cm}{0.4pt} (x^{(1)})^\top \rule{1cm}{0.4pt} \\ \rule{1cm}{0.4pt} (x^{(2)})^\top \rule{1cm}{0.4pt} \\ \vdots \\ \rule{1cm}{0.4pt} (x^{(m)})^\top \rule{1cm}{0.4pt} \end{bmatrix}$$

$m$ **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)

$$= \begin{bmatrix} \underline{\quad} (x^{(1)})^\top \underline{\quad} \\ \underline{\quad} (x^{(2)})^\top \underline{\quad} \\ \vdots \\ \underline{\quad} (x^{(m)})^\top \underline{\quad} \end{bmatrix}$$

$m \times (n+1)$

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

$m$ **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)

$$= \begin{bmatrix} \underline{\quad} (x^{(1)})^\top \underline{\quad} \\ \underline{\quad} (x^{(2)})^\top \underline{\quad} \\ \vdots \\ \underline{\quad} (x^{(m)})^\top \underline{\quad} \end{bmatrix}$$

$m \times (n+1)$

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

$m$ **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) $= \begin{bmatrix} \underline{\quad} (x^{(1)})^\top \underline{\quad} \\ \underline{\quad} (x^{(2)})^\top \underline{\quad} \\ \vdots \\ \underline{\quad} (x^{(m)})^\top \underline{\quad} \end{bmatrix}$

$m \times (n+1)$

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

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

Andrew Ng

$m$ **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)

$$= \begin{bmatrix} \underline{\quad} (x^{(1)})^T \underline{\quad} \\ \underline{\quad} (x^{(2)})^T \underline{\quad} \\ \vdots \\ \underline{\quad} (x^{(m)})^T \underline{\quad} \end{bmatrix}$$

$m \times (n+1)$

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

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

$m \times 2$

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

Andrew Ng

$m$ **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)

$$= \begin{bmatrix} \underline{\quad} (x^{(1)})^\top \underline{\quad} \\ \underline{\quad} (x^{(2)})^\top \underline{\quad} \\ \vdots \\ \underline{\quad} (x^{(m)})^\top \underline{\quad} \end{bmatrix}$$

$m \times (n+1)$

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

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

$m \times 2$

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

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

Andrew Ng

# Exercise

- 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
- $weight = \theta_0 + \theta_1 x_1 + \theta_2 x_2$
- What are X and y?

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

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

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

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

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

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

Set $\underline{A} = \underline{X^T X}$

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

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

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

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

Set $\underline{A} = \underline{X^T X}$

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

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

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

$(X^T X)^{-1} X^T y$

$X' \qquad X^T$

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

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

When to choose gradient descent and when to choose normal equations???

Set $A = X^T X$

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

No need for feature scaling!

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

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

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

$$\min_\theta J(\theta)$$

$X'$ $\quad$ $X^T$

Feature Scaling

$0 \leq x_1 \leq 1$

$0 \leq x_2 \leq 1000$

$0 \leq x_3 \leq 10^{-5}$

Andrew Ng

$m$ **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. |

# $m$ 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. | |

$m$ **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 $(X^T X)^{-1}$ |
| | • Slow if $n$ is very large. |

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

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

$m$ **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 $(X^T X)^{-1}$   $n \times n$   $O(n^3)$ |
| | • Slow if $n$ is very large. |

$n = 100$

$n = 1000$

$m$ **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 $(X^T X)^{-1}$ $n \times n$  $O(n^3)$ |
|  | • Slow if $n$ is very large. |

$n = 100$

$n = 1000$

$n = 10000$

$m$ **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. |

$n = 10^6$

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

| Gradient Descent | Normal Equation |
| --- | --- |
| Need to choose alpha | No need to choose alpha |
| Needs many iterations | No need to iterate |
| $O(kn^2)$ | $O(n^3)$, need to calculate inverse of $X^TX$ |
| Works well when n is large | Slow if n is very large |