

Introduction to ML



George Igwegbe
Artificial Intelligence
Saturday (AI6), Lagos



WK3

1. Welcome
2. What is ML?
3. Supervised Learning
4. Unsupervised Learning
5. Model Representation
6. Cost function
7. Cost function 2
8. Gradient Descent
9. Gradient Descent 2
10. Gradient Descent for Linear Regression
11. What's Next
12. Matrices and Vectors
13. Addition and Scalar Multiplication
14. Matrix Vector Multiplication
15. Matrix Vector Multiplication 2
16. Matrix Vector Multiplication 3
17. Matrix Vector Properties
18. Inverse and Transpose
19. Multiple Features



Machine Learning

BUZZWORD

Multivariate Linear Regression

Singular

Degenerate

Normal Equation

Condition of adding matrices

Condition of multiplying matrices

Dimension of matrix

Transpose of matrix

1-indexed vs 0-indexed

Inverse of matrix

MULTIPLE FEATURES

Size (feet ²)	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

- Difficult to visualize->3D
- Notations becomes complicated

LINEAR ALGEBRA

Size (feet ²)	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

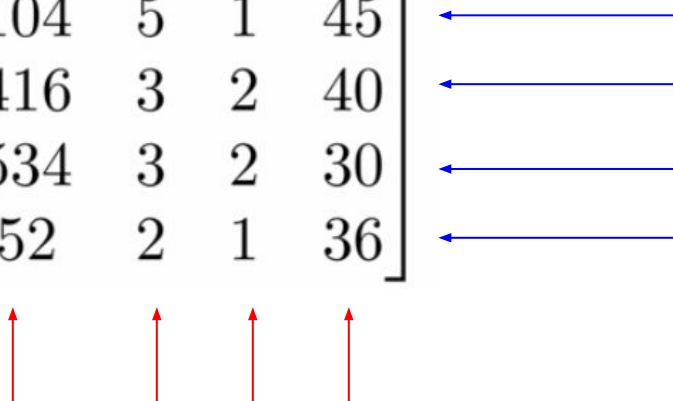
$$X = \begin{bmatrix} 2104 & 5 & 1 & 45 \\ 1416 & 3 & 2 & 40 \\ 1534 & 3 & 2 & 30 \\ 852 & 2 & 1 & 36 \end{bmatrix}$$

MATRIX

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 172 \end{bmatrix}$$

VECTOR

MATRIX: Rectangular array of numbers

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


Dimension of matrix: number of rows x number of columns

(? x ?) matrix

Matrix Elements (entries of matrix)

$$A = \begin{bmatrix} 1402 & 191 \\ 1371 & 821 \\ 949 & 1437 \\ 147 & 1448 \end{bmatrix}$$

A_{ij} = " i, j entry" in the i^{th} row, j^{th} column.

$$A_{11} = 1402$$

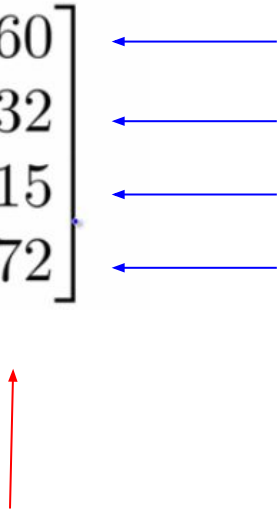
$$A_{12} = 191$$

$$A_{32} = 1437$$

$$A_{41} = 147$$

$$\cancel{A_{43}} = \text{Undefined (error)}$$

VECTOR: An $n \times 1$ matrix.

$$y = \begin{bmatrix} 460 \\ 232 \\ 315 \\ 172 \end{bmatrix}$$


(? x ?) matrix

Dimension of matrix: number of rows x number of columns

VECTOR: An $n \times 1$ matrix.

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

$\uparrow \uparrow$
 $n = 4$
 \leftarrow 4-dimensional vector.

$y_i = i^{\text{th}}$ element

$$y_1 = 460$$

$$y_2 = 232$$

$$y_3 = 315$$

1-indexed vs 0-indexed:

$$y[i] \quad y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \quad \text{1-indexed}$$
$$y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix} \quad \text{0-indexed}$$

$y[0]$

Matrix Addition

$$\begin{array}{l} \rightarrow \\ \rightarrow \\ \rightarrow \end{array} \begin{bmatrix} \textcircled{1} & 0 \\ \textcircled{2} & 5 \\ \textcircled{3} & 1 \end{bmatrix} + \begin{bmatrix} \textcircled{4} & 0.5 \\ \textcircled{2} & 5 \\ \textcircled{0} & 1 \end{bmatrix} = \begin{bmatrix} 5 & 0.5 \\ 4 & 10 \\ 3 & 2 \end{bmatrix}$$

3×2 matrix 3×2 3×2

$$\begin{array}{l} \rightarrow \\ \rightarrow \\ \rightarrow \end{array} \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} + \begin{bmatrix} 4 & 0.5 \\ 2 & 5 \end{bmatrix} = \text{error}$$

3×2 2×2

Scalar Multiplication

$$\begin{array}{c} 3 \times \begin{bmatrix} \textcircled{1} & 0 \\ \textcircled{2} & 5 \\ \textcircled{3} & 1 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 6 & 15 \\ 9 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 2 & 5 \\ 3 & 1 \end{bmatrix} \times 3 \\ \underline{3 \times 2} \qquad \underline{3 \times 2} \end{array}$$

$$\begin{array}{c} \boxed{\begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} / 4 =} \end{array} \quad \frac{1}{4} \begin{bmatrix} 4 & 0 \\ 6 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ \frac{3}{2} & \frac{3}{4} \end{bmatrix}$$

Condition of multiplication of matrix

$$\begin{array}{ccccc} A & \times & x & = & y \\ \left[\begin{array}{c} \\ \\ \end{array} \right] & \times & \left[\begin{array}{c} \\ \\ \end{array} \right] & = & \left[\begin{array}{c} \\ \\ \end{array} \right] \\ \text{m x n matrix} & & \text{n x 1 matrix} & & \text{m-dimensional} \\ \text{(m rows,} & & \text{(n-dimensional} & & \text{vector} \\ \text{n columns)} & & \text{vector)} & & \end{array}$$

To get y_i , multiply A 's i^{th} row with elements of vector x , and add them up.

Example of Matrix Multiplication

$$\begin{bmatrix} 1 & 2 & 1 & 5 \\ 0 & 3 & 0 & 4 \\ -1 & -2 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 3 \\ 2 \\ 1 \end{bmatrix} = \begin{bmatrix} 14 \\ 13 \\ -7 \end{bmatrix}$$

3×4 4×1 3×1

$$1 \times 1 + 2 \times 3 + 1 \times 2 + 5 \times 1 = 14$$

$$0 \times 1 + 3 \times 3 + 0 \times 2 + 4 \times 1 = 13$$

$$-1 \times 1 + (-2) \times 3 + 0 \times 2 + 0 \times 1 = -7$$

Application of Matrix Multiplication

House sizes:

- 2104
- 1416
- 1534
- 852

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix}$$

Handwritten notes: A red arrow points to the first row with "4x2". A green box highlights the first row. A pink box highlights the second row.

$$h_{\theta}(x) = -40 + 0.25x$$

$h_{\theta}(x)$

Handwritten notes: "2x1" and "Vector" with a red arrow pointing to the vector below.

$$\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$$

Handwritten note: A red "X" is written to the left of the vector.

Handwritten note: "4x1 matrix" with a red arrow pointing to the matrix below.

$$\begin{bmatrix} -40 \times 1 + 0.25 \times 2104 \\ -40 \times 1 + 0.25 \times 1416 \end{bmatrix}$$

Handwritten notes: "h_{\theta}(2104)" points to the first row. "h_{\theta}(1416)" points to the second row. A pink box highlights the second row.

Condition of multiplication of matrix

$$\begin{array}{ccccc} \underline{A} & \times & B & = & C \\ \left[\begin{array}{c} \\ \\ \end{array} \right] & \times & \left[\begin{array}{c} \\ \\ \end{array} \right] & = & \left[\begin{array}{c} \\ \\ \end{array} \right] \end{array}$$

$m \times n$ matrix
(m rows,
 n columns)

$n \times o$ matrix
(n rows,
 o columns)

$m \times o$
matrix

Example of Matrix Multiplication

$$\overset{2 \times 2}{\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix}} \overset{2 \times 2}{\begin{bmatrix} 0 & 1 \\ 3 & 2 \end{bmatrix}} = \overset{2 \times 2}{\begin{bmatrix} 9 & 4 \\ 15 & 12 \end{bmatrix}}$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 0 \\ 3 \end{bmatrix} = \begin{bmatrix} 1 \times 0 + 3 \times 3 \\ 2 \times 0 + 5 \times 3 \end{bmatrix} = \begin{bmatrix} 9 \\ 15 \end{bmatrix}$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 5 \end{bmatrix} \begin{bmatrix} 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 1 \times 1 + 3 \times 2 \\ 2 \times 1 + 5 \times 2 \end{bmatrix} = \begin{bmatrix} 7 \\ 12 \end{bmatrix}$$

House sizes:

$$\begin{cases} 2104 \\ 1416 \\ 1534 \\ 852 \end{cases}$$

Have 3 competing hypotheses:

$$1. h_{\theta}(x) = -40 + 0.25x$$

$$2. h_{\theta}(x) = 200 + 0.1x$$

$$3. h_{\theta}(x) = -150 + 0.4x$$

Matrix

Matrix

$$\begin{bmatrix} 1 & 2104 \\ 1 & 1416 \\ 1 & 1534 \\ 1 & 852 \end{bmatrix} \times$$

$$\begin{bmatrix} -40 \\ 0.25 \end{bmatrix}$$

$$\begin{bmatrix} 200 \\ 0.1 \end{bmatrix}$$

$$\begin{bmatrix} -150 \\ 0.4 \end{bmatrix} =$$

$$\begin{bmatrix} 486 \\ 314 \\ 344 \\ 173 \end{bmatrix} \begin{bmatrix} 410 \\ 342 \\ 353 \\ 285 \end{bmatrix} \begin{bmatrix} 692 \\ 416 \\ 464 \\ 191 \end{bmatrix}$$

Prediction
of first
 h_{θ}

Predictions
of 2nd
 h_{θ}

Commutativity: Properties of Matrices

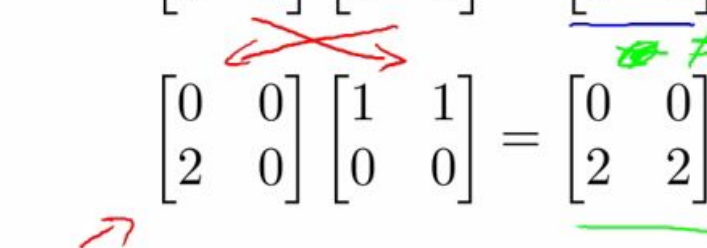
$$3 \times 5 = 5 \times 3$$


"Commutative"

Let A and B be matrices. Then in general,

$$\underline{A \times B} \neq \underline{B \times A.} \text{ (not commutative.)}$$

E.g.

$$\begin{array}{l} \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix} = \begin{bmatrix} 2 & 0 \\ 0 & 0 \end{bmatrix} \\ \begin{bmatrix} 0 & 0 \\ 2 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} = \begin{bmatrix} 0 & 0 \\ 2 & 2 \end{bmatrix} \end{array} \quad \left| \quad \begin{array}{l} A \times B \\ m \times n \quad n \times m \\ \hline A \times B \quad \text{is} \quad m \times m \\ B \times A \quad \text{is} \quad n \times n \end{array} \right.$$


Associativity: Properties of Matrices

$$\underline{3 \times 5 \times 2} \quad 3 \times (5+2) = (3 \times 5) \times 2$$

$3 \times 10 = 30 = 15 \times 2$

"Associative"

$$A \times (B \times C) \quad \leftarrow$$
$$(\underline{A \times B}) \times C \quad \leftarrow$$

$$A \times B \times C.$$

Let $D = B \times C$. Compute $A \times D$.

Let $E = A \times B$. Compute $E \times C$.

$A \times (B \times C)$
 $(A \times B) \times C$
Some
answer.

Identity : Properties of Matrices

Identity Matrix

Denoted I (or $I_{n \times n}$).

Examples of identity matrices:

$$\begin{bmatrix} 1 \end{bmatrix}$$

1×1

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$

2×2

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

3×3

$$\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}$$

4×4

Informally:

$$\begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{bmatrix}$$

For any matrix A ,

$$A \cdot I = I \cdot A = A$$

$\begin{matrix} \nearrow & \uparrow & \uparrow & \nwarrow & \nwarrow \\ m \times n & n \times n & m \times m & m \times n & m \times n \end{matrix}$

$$I_{n \times n}$$

Note:

$AB \neq BA$ in general

$$AI = \cancel{IA} IA \checkmark$$

1 is identity.

$$1 \times z = z \times 1 = z$$

for any z

Inverse: Properties of Matrix

1 = "identity."

$$3 \underbrace{(\underbrace{3^{-1}}_{\frac{1}{3}})} = 1$$

$$12 \times \underbrace{(12^{-1})}_{\frac{1}{12}} = 1$$

0 (0⁻¹) undefined

Not all numbers have an inverse.

Matrix inverse:

square matrix
(#rows = #columns)

A^{-1}

If A is an $m \times m$ matrix, and if it has an inverse,

$$\rightarrow \underline{A}(\underline{A^{-1}}) = \underline{A^{-1}}\underline{A} = \underline{I}.$$

E.g.

$$\underbrace{\begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix}}_{\substack{2 \times 2 \\ A}} \underbrace{\begin{bmatrix} 0.4 & -0.1 \\ -0.05 & 0.075 \end{bmatrix}}_{A^{-1}} = \underbrace{\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}}_{A^{-1}A} = I_{2 \times 2}$$

Matrices that **don't** have an inverse are "singular" or "degenerate"

Transpose: Properties of Matrix

Example:

$$\underline{A} = \begin{bmatrix} 1 & 2 & 0 \\ 3 & 5 & 9 \end{bmatrix} \quad \underline{B} = \underline{A}^T = \begin{bmatrix} 1 & 3 \\ 2 & 5 \\ 0 & 9 \end{bmatrix}$$

2×3 3×2

Let A be an $m \times n$ matrix, and let $B = A^T$.
Then B is an $n \times m$ matrix, and

$$\underline{B}_{ij} = \underline{A}_{ji}.$$

$$B_{12} = A_{21} = 2$$

$$B_{32} = 9$$

$$A_{23} = 9.$$

MULTIPLE FEATURES

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

MULTIPLE FEATURES

Size (feet ²)	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
...

$m = 47$

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.

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

$x_3^{(2)} = 2$

Representation of Hypothesis

$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n$$

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

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

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$\downarrow = 1$

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

θ^T
(n+1) x 1 matrix

x

Multivariate linear regression.

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

Parameters: $\theta_0, \theta_1, \dots, \theta_n$

Cost function:

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



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

Parameters: $\theta_0, \theta_1, \dots, \theta_n$

θ

n+1-dimensional vector

VECTOR

Cost function:

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

J(θ)

Gradient descent:

Repeat {

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

}

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



Gradient descent:

Repeat {

→ $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} \underline{J(\theta)}$

}

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

VECTOR

Gradient Descent

Previously ($n=1$):

Repeat {

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

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

(simultaneously update θ_0, θ_1)

}

➤ New algorithm ($n \geq 1$):

Repeat {

$\swarrow \frac{2}{2\theta_j} J(\theta)$

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

(simultaneously update θ_j for $j = 0, \dots, n$)

}

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

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

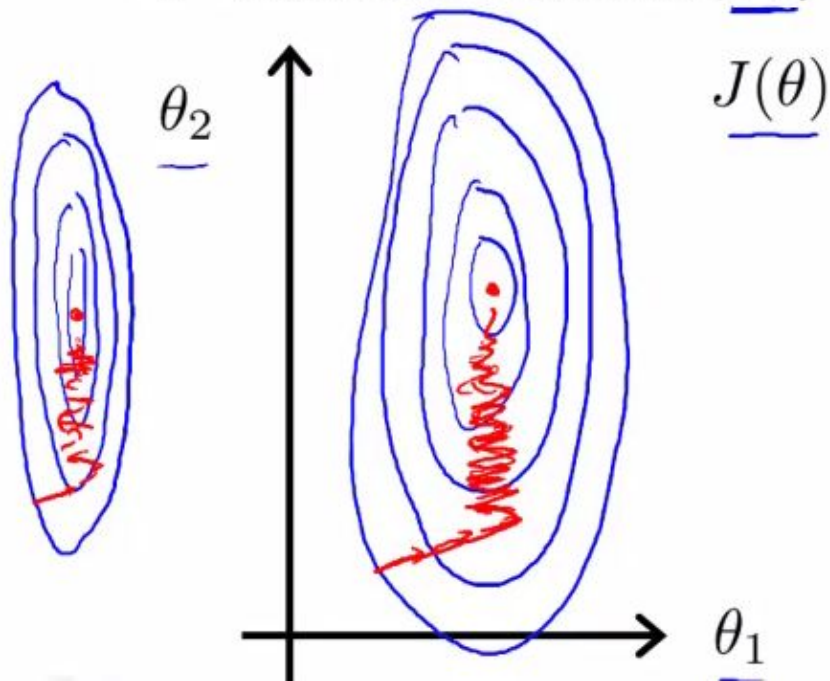
$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

Feature Scaling

Idea: Make sure features are on a similar scale.

E.g. $x_1 = \text{size (0-2000 feet}^2\text{)}$ ←

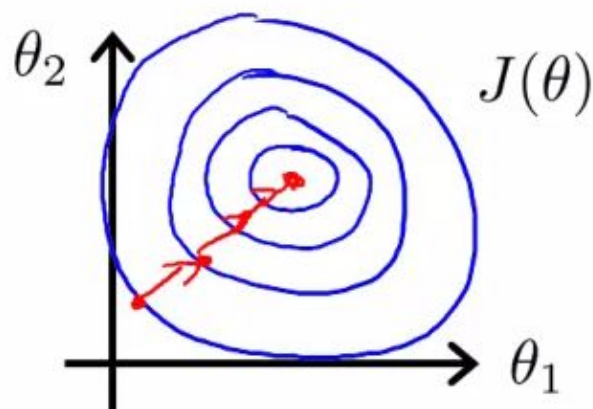
$x_2 = \text{number of bedrooms (1-5)}$ ←



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

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

$$0 \leq x_1 \leq 1 \quad 0 \leq x_2 \leq 1$$



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

Average size = 1000

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

1-5 bedrooms

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

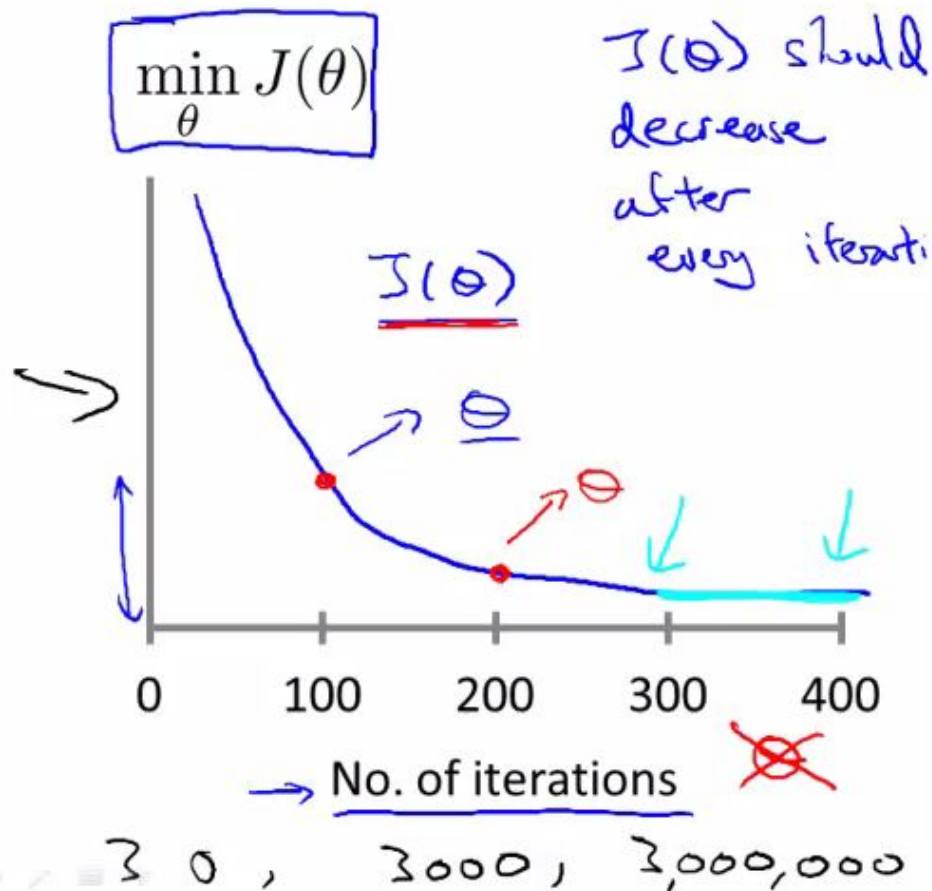
$$x_1 \leftarrow \frac{x_1 - \mu_1}{s_1}$$

← avg value of x_1 in training set

$$x_2 \leftarrow \frac{x_2 - \mu_2}{s_2}$$

range (max-min)
(or standard deviation)

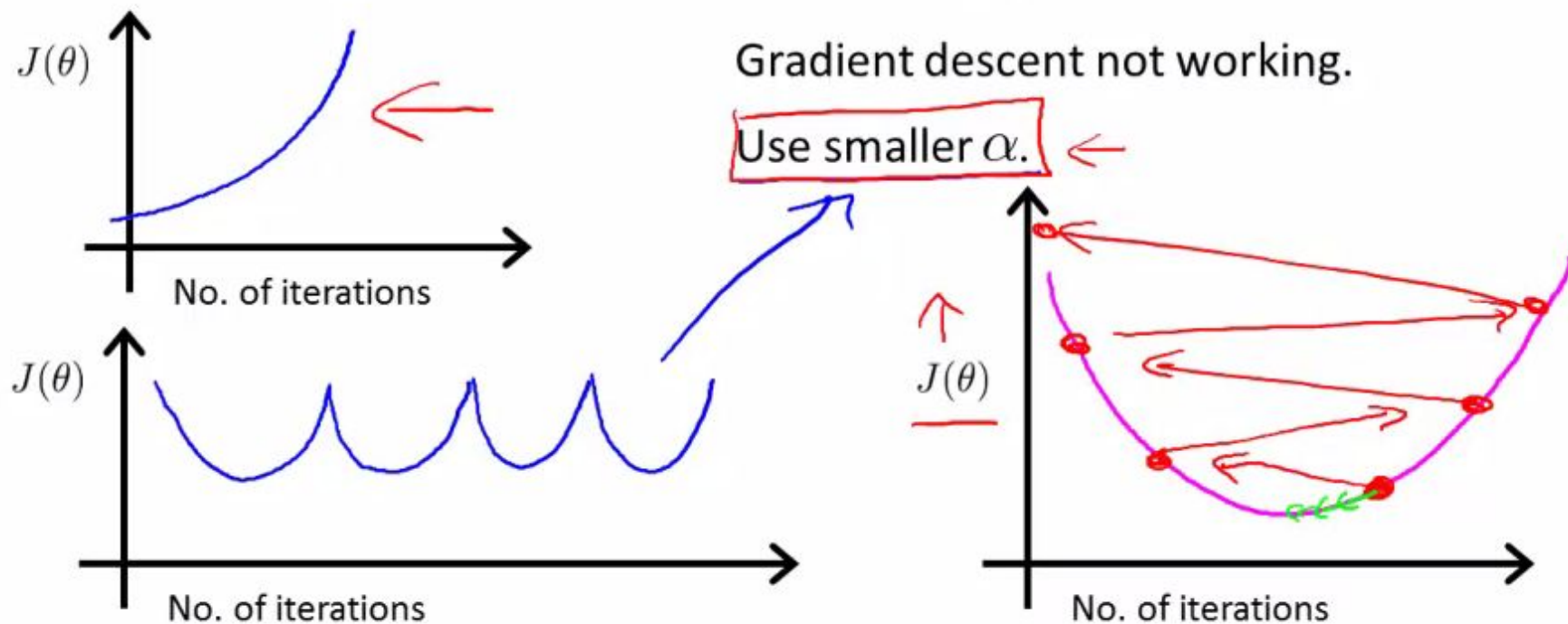
Making sure gradient descent is working correctly.



→ Example automatic convergence test:

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

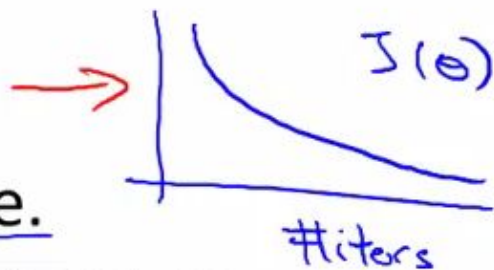
Making sure gradient descent is working correctly.



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

Summary:

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



To choose α , try

..., 0.001, 0.003, 0.01, 0.03, 0.1, 0.3, 1, ...

↑ 3x ~1x 3x ~3x ↑ ↑

Housing prices prediction

$$h_{\theta}(x) = \theta_0 + \theta_1 \times \underbrace{\text{frontage}}_{x_1} + \theta_2 \times \underbrace{\text{depth}}_{x_2}$$

Area

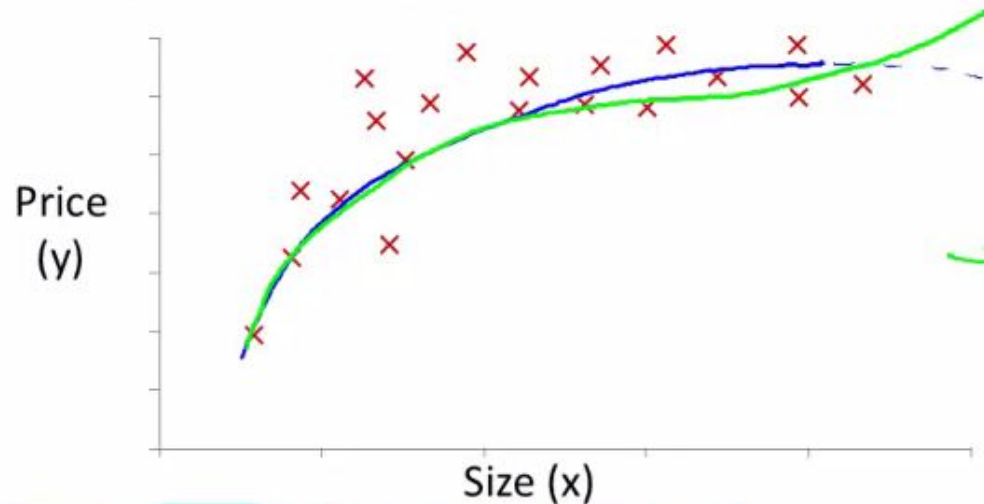
$$x = \underline{\text{frontage} * \text{depth}}$$

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

\nearrow land area



Polynomial regression



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

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

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

$$= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3$$

$$\rightarrow x_1 = (\text{size})$$

$$\rightarrow x_2 = (\text{size})^2$$

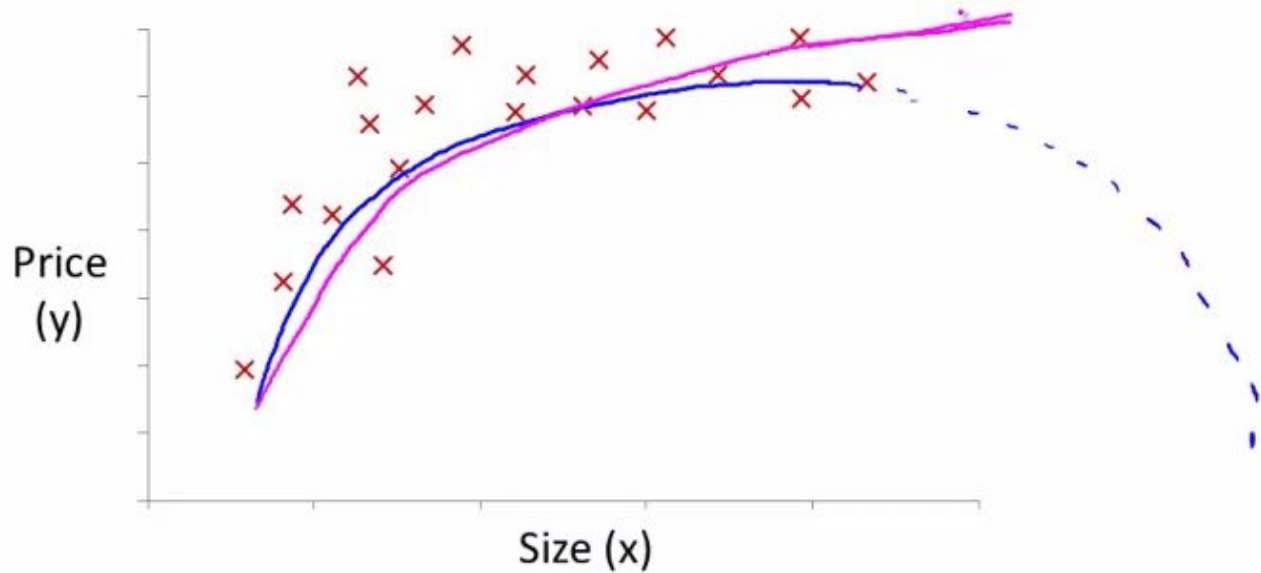
$$\rightarrow x_3 = (\text{size})^3$$

Size: 1-1000

Size²: 1-1,000,000

Size³: 1-10⁹

Choice of features



$$\rightarrow h_{\theta}(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2$$

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



NORMAL EQUATION

In $\min J(\theta_0, \theta_1)$, solve for θ_0, θ_1 exactly, without needing iterative algorithm (gradient descent).

Examples: $m = 4$.

	Size (feet ²)	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
x_0	x_1	x_2	x_3	x_4	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

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

$m \times (n+1)$

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

m -dimensional vector

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

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

$$\underline{x^{(i)}} = \begin{bmatrix} x_0^{(i)} \\ x_1^{(i)} \\ x_2^{(i)} \\ \vdots \\ x_n^{(i)} \end{bmatrix} \in \mathbb{R}^{n+1} \quad \bigg| \quad \begin{matrix} \times \\ \text{(design} \\ \text{matrix)} \end{matrix} = \begin{bmatrix} \text{---} (x^{(1)})^T \text{---} \\ \text{---} (x^{(2)})^T \text{---} \\ \vdots \\ \text{---} (x^{(m)})^T \text{---} \end{bmatrix}$$

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

$$\begin{matrix} m \times (n+1) \\ \begin{bmatrix} 1 & x_1^{(1)} \\ 1 & x_1^{(2)} \\ \vdots & \vdots \\ 1 & x_1^{(m)} \end{bmatrix} \end{matrix} \bigg| \begin{matrix} m \times 2 \\ \underline{y} = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \end{matrix}$$

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

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

Set $A = X^T X$

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

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

m training examples, n features.

Gradient Descent

- • Need to choose α .
- • Needs many iterations.
- Works well even when n is large.

Normal Equation

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

m training examples, n features.

Gradient Descent

- • Need to choose α .
- • Needs many iterations.
- Works well even when n is large.

- $n = 10^6$ (1 million features)

Normal Equation

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

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

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

- Redundant features (linearly dependent).
E.g. $x_1 = \text{size in feet}^2$
 $x_2 = \text{size in m}^2$
- Too many features (e.g. $m \leq n$).
 - Delete some features, or use regularization.

Hinton's Closing Prayer

Our father who art in n -dimensions

hallowed by the backprop,

thy loss be minimized,

thy gradients unvarnished,

on earth as it is in Euclidean space.

Give us this day our daily hyperparameters,

and forgive us our large learning rates,

as we forgive those whose parameters diverge,

and lead us not into discrete optimization,

but deliver us from local optima.

For thine are dimensions,

and the GPUs, and the glory,

forever and ever. Dropout.



From buZZrobot