



AI 330: Machine Learning

Fall 2023

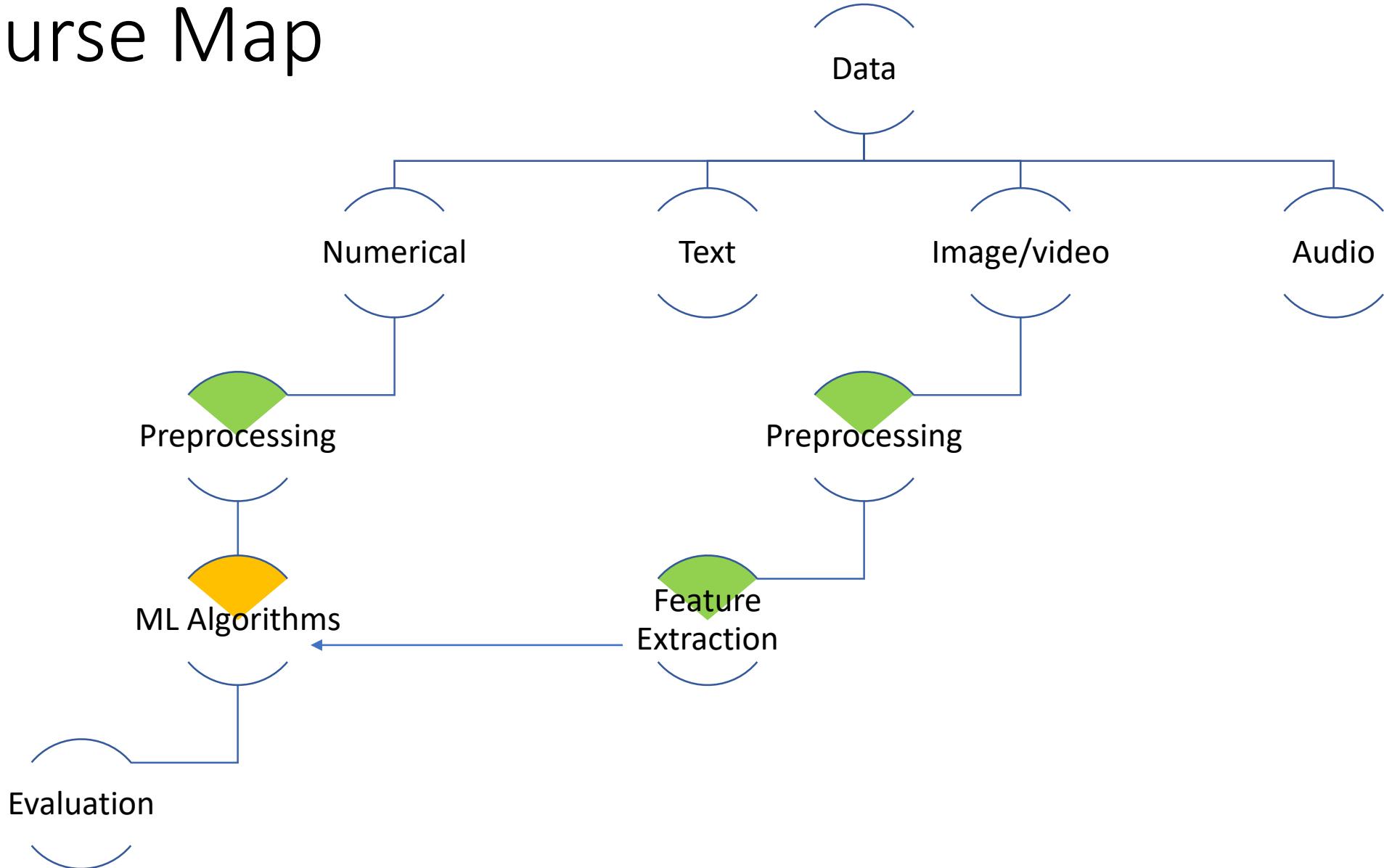
Dr. Wessam EL-Behaidy

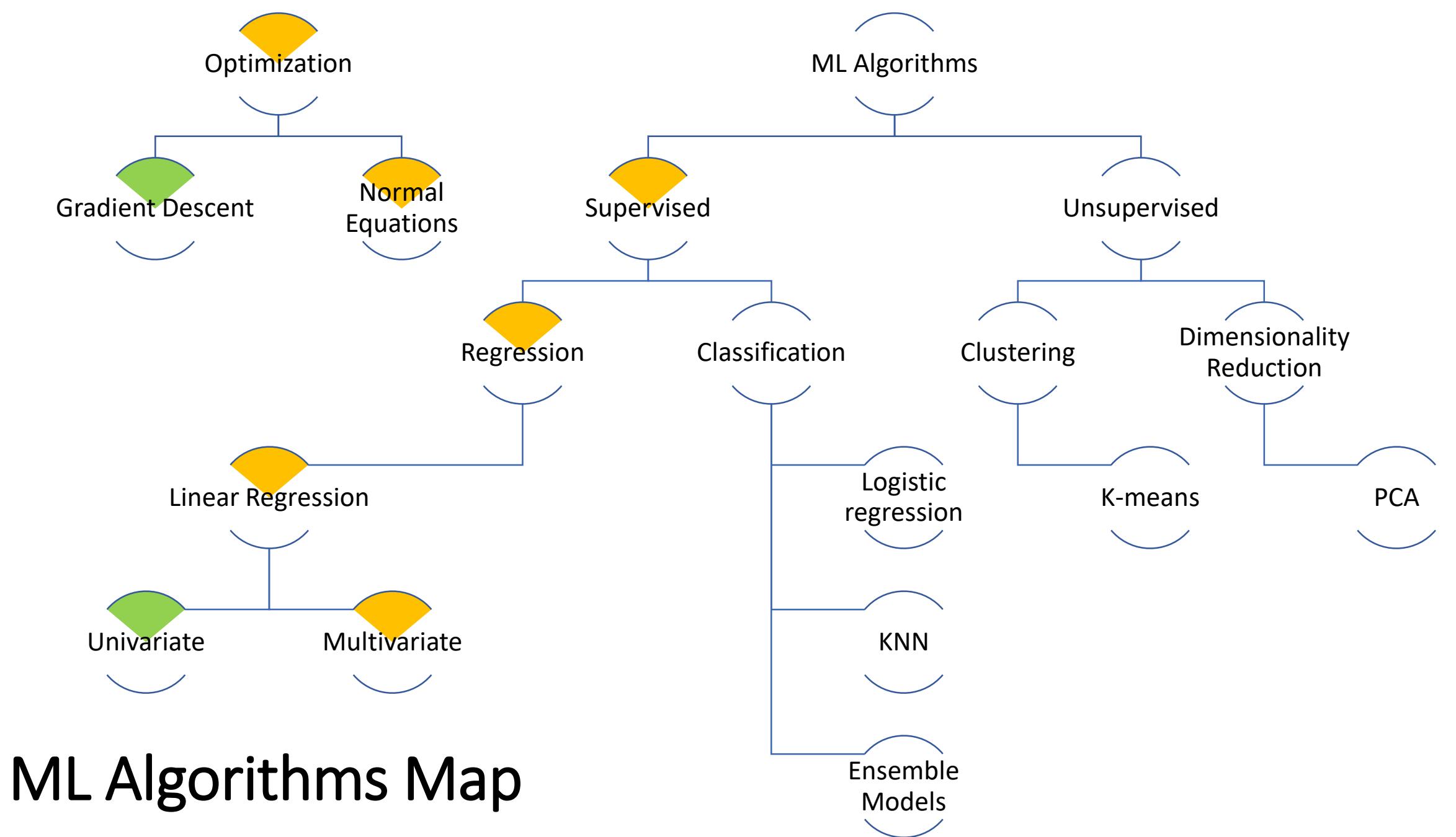
Associate Professor, Computer Science Department,
Faculty of Computers and Artificial Intelligence,
Helwan University.

Dr. Ensaif Hussein

Associate Professor, Computer Science Department,
Faculty of Computers and Artificial Intelligence,
Helwan University.

Course Map





Lecture 5

Multivariate Linear Regression

&

Normal Equation

Slides of:

Machine Learning Specialization <https://www.coursera.org/specializations/machine-learning-introduction>
at Stanford University (**Andrew Ng**)

Review: Matrices and vectors

Vector: Special case of the matrix

$n \times 1$ matrix 1 dimension
 1 vector

$$\bullet y = \begin{bmatrix} 4 \\ 50 \\ 3 \\ 25 \end{bmatrix} \in \mathbb{R}^4$$

4-dimensional vector $n = 4$

1-indexed vs 0-indexed:

$$y = \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix}$$

$$y = \begin{bmatrix} y_0 \\ y_1 \\ y_2 \\ y_3 \end{bmatrix}$$

y_i = i^{th} element.

$$y_2 = 50 \quad y_2 = 3$$

يساوى ٣x٢ فالمatrix 3×2 مكونة من ٣ columns و ٢ rows .
يساوى ٢x٥ فالمatrix 2×5 مكونة من ٢ columns و ٥ rows .
نجلب ٢matrix من المكتبات 3×2 و 2×5 .

Matrix Multiplication

Column's row wise

$$A \times x = y$$
$$\begin{bmatrix} 5 & 1 \\ 2 & 6 \\ 3 & 4 \end{bmatrix} \times \begin{bmatrix} 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 15 \\ 34 \\ 26 \end{bmatrix}$$

3×2 2×1 3×1 matrix \rightarrow 3 dim vector

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

- $y_1 = 5 \times 2 + 1 \times 5 = 15$
- $y_2 = 2 \times 2 + 6 \times 5 = 34$
- $y_3 = 3 \times 2 + 4 \times 5 = 26$

An $m \times n$ matrix multiplied by an $n \times 1$ vector results in an $m \times 1$ vector.

$$A \times B = C$$
$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \\ 2 & 1 \end{bmatrix} \times \begin{bmatrix} 1 \\ 0 \\ 5 \end{bmatrix} = \begin{bmatrix} 7 \\ 16 \\ 15 \end{bmatrix}$$

2×3 3×2 2×2 matrix

- $C_{11} = 1 + 0 + 6 = 7$
- $C_{21} = 4 + 0 + 12 = 16$
- $C_{12} = 2 + 10 + 3 = 15$
- $C_{22} = 8 + 25 + 6 = 39$

An $m \times n$ matrix multiplied by an $n \times o$ matrix results in an $m \times o$ matrix.

Matrix Multiplication Properties

- Matrices are **not commutative**: $A \times B \neq B \times A$
- Matrices are **associative**: $(A \times B) \times C = A \times (B \times C)$
- Identity Matrix I , or $I_{n \times n}$ for any A , $A \cdot I = I \cdot A = A$

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

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

①, all cols 1 w/e

Inverse and Transpose

- The **inverse** of a matrix A ($m \times m$) is denoted A^{-1} . Multiplying by the inverse results in the identity matrix. $A \times A^{-1} = 1$
A non square matrix does not have an inverse matrix.

$$\begin{bmatrix} 3 & 4 \\ 2 & 16 \end{bmatrix} \times \begin{bmatrix} 0.4 & -0.1 \\ -0.05 & 0.075 \end{bmatrix} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \text{ identity}$$

- The **transposition** $B_{ij} = A_{ji}$ of a matrix is like rotating the matrix 90° in clockwise direction and then reversing it.

- $A_{ij} = A_{ij}^T$

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix}^T = \begin{bmatrix} 1 & 4 \\ 2 & 5 \\ 3 & 6 \end{bmatrix}$$

View) 6mm 11/2
Anschrift

$m \times n$
 $n \times m$

Regression Problem: Univariate Linear Regression Example

Linear Regression with One Variable

Recap:

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

Parameters: θ_0, θ_1

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

Goal: $\underset{\theta_0, \theta_1}{\text{minimize}} J(\theta_0, \theta_1)$

Recap: Gradient Descent

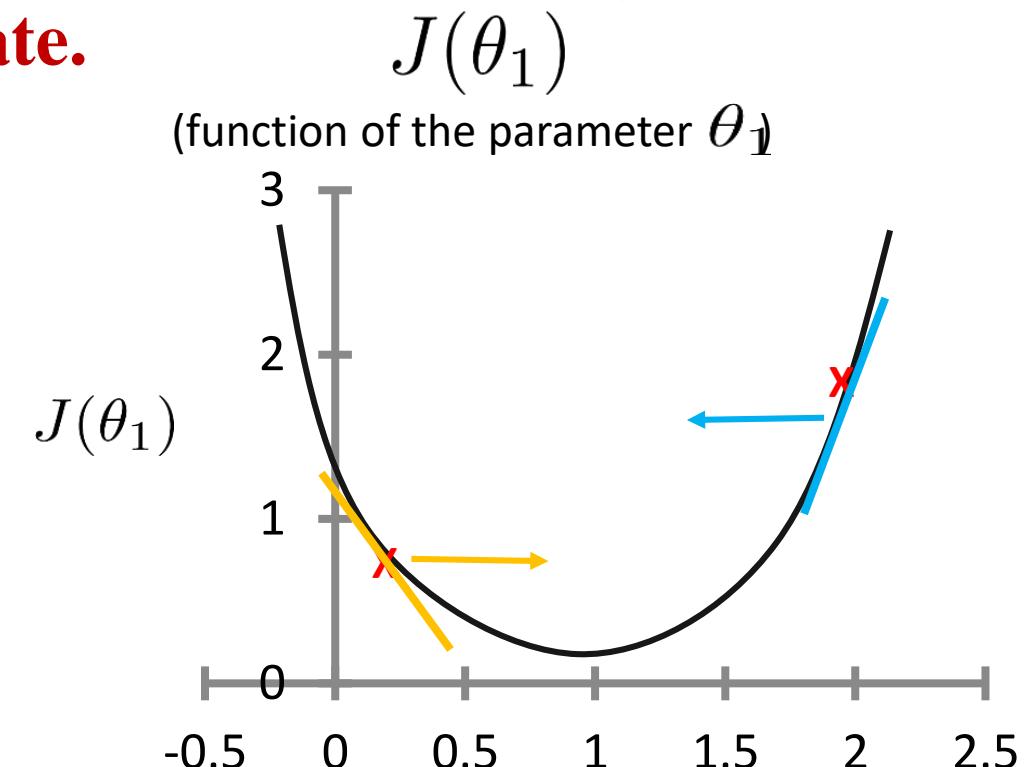
- We make steps down the cost function in the direction with the steepest descent, and the size of each step is determined by the parameter α , which is called the **learning rate**.
- The gradient descent algorithm is:
- Repeat until convergence {

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

(for $j = 0$ and $j = 1$)

}

Learning rate (step size)



Positive slope (positive number) $\rightarrow \Theta$ will decrease
Negative slope (negative number) $\rightarrow \Theta$ will increase

Recap: Gradient Descent for Linear Regression

- When specifically applied to the case of **linear regression**, a new form of the gradient descent equation can be derived. We can substitute our actual cost function and our actual hypothesis function and modify the equation to:

- Start by initializing the parameters θ_0, θ_1 randomly
- Repeat until convergence { // until error is small
 - Predicted values with linear regression hypothesis.
 - Calculate the cost function.
 - If cost is large, update parameters using GD

Should be done simultaneously

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \quad \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) :$$
$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)}) \cdot x^{(i)} \quad \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) :$$

}

Hypothesis using Matrix Product

Matrix is hypothesis function

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

← *Ans)*
uni variant

$$x = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} \in \mathbb{R}^2$$
$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \end{bmatrix} \in \mathbb{R}^2$$

For convenience of notation, define $x_0 = 1$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 = [\theta_0 \theta_1] \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \theta^T x$$

↑
transpose
↑
row vector
↑
Column
of θ

↓
(1). Out of 2 matrix multiplication
→ *Observe*

Linear Regression with one variable Example

- Suppose we are given a dataset, and we need to find the hypothesis of linear regression:

Experience (x)	Salary (y) (in lakhs)
2	3
6	10
5	4
7	13

<https://www.geeksforgeeks.org/mathematical-explanation-for-linear-regression-working/>

Iteration 1

$\theta_1 \rightarrow \theta_0$ I Vandom numbers bgrd

a) In the start, θ_0 and θ_1 values are randomly chosen.

- Let us suppose, $\theta_0 = 0$ and $\theta_1 = 0$.

b) Predicted values after iteration 1 with Linear regression hypothesis.

$$h_{\theta} = \theta_0 x_0 + \theta_1 x_1 = [\theta_0 \theta_1] \begin{bmatrix} x_0 \\ x_1 \end{bmatrix}$$

Initialising Parameters
Zero, Zero

$$h_{\theta} = [\theta_0 \theta_1] \begin{bmatrix} x_0 & x_0 & x_0 & x_0 \\ x_1 & x_2 & x_3 & x_4 \end{bmatrix}$$

With model fitting

$$h_{\theta} = [0 \ 0] \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 6 & 5 & 7 \end{bmatrix}$$

One training example

$$h_{\theta} = [0 \ 0 \ 0 \ 0]$$

Ans for all 4

Experience (x)	Salary (y) (in lakhs)
2	3
6	10
5	4
7	13

Iteration 1

c) Cost function Error (m=No. of samples = 4)

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_\theta(x_i) - y_i]^2$$

$$= \frac{1}{2 * 4} [(0 - 3)^2 + (0 - 10)^2 + (0 - 4)^2 + (0 - 13)^2]$$

hypothesis function

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m [h_\theta(x_i) - y^{(i)}]^2$$

$$= \frac{1}{2 * 4} [(0 - 3)^2 + (0 - 10)^2 + (0 - 4)^2 + (0 - 13)^2]$$

$$= 36.75$$

Experience (x)	Salary (y) (in lakhs)
2	3
6	10
5	4
7	13

Iteration 1

d) Gradient Descent – Updating θ_0 value

Here, $j=0$, $\alpha = 0.001$

$$\theta_0 = \theta_0 - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i)$$

$$= 0 - \frac{0.001}{4} [(0 - 3) + (0 - 10) + (0 - 4) + (0 - 13)]$$

$$= - \frac{0.001}{4} [(-3) + (-10) + (-4) + (-13)]$$

$$= - \frac{0.001}{4} [-30]$$

$$= 0.0075$$

new
 θ_0

Gradient desc to θ_0
Cost(2) \rightarrow θ_0
Min with 2 loops
2 min on m \rightarrow square min GD

Algorithm GD for linear regression:

repeat until convergence {

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

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

Experience (x)	Salary (y) (in lakhs)
2	3
6	10
5	4
7	13

Iteration 1

e) Gradient Descent – Updating θ_1 value

Here, $j=1$

$$\theta_1 = \theta_1 - \frac{\alpha}{m} \sum_{i=1}^m (h_\theta(x_i) - y_i) \cdot x_i$$

*on θ_1 Nj
 ∑ Nj xj use*

$$= 0 - \frac{0.001}{4} [(0 - 3)2 + (0 - 10)6 + (0 - 4)5 + (0 - 13)7]$$

$$= - \frac{0.001}{4} [(-6) + (-60) + (-20) + (-91)]$$

$$= - \frac{0.001}{4} [-177]$$

$$= 0.04425$$

Experience (x)	Salary (y) (in lakhs)
2	3
6	10
5	4
7	13

Algorithm GD for linear regression:

repeat until convergence {

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

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

Iteration 2

Iteration 2 – $\theta_0 = 0.0075$ and $\theta_1 = 0.04425$

- Predicting values after iteration 1 with linear regression hypothesis

$$h_{\theta} = [\theta_0 \quad \theta_1] \begin{bmatrix} x_0 & x_0 & x_0 & x_0 \\ x_1 & x_2 & x_3 & x_4 \end{bmatrix}$$

$$= [0.0075 \quad 0.04425] \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 6 & 5 & 7 \end{bmatrix}$$

$$= [0.096 \quad 0.273 \quad 0.22875 \quad 0.31725]$$

Experience (x)	Salary (y) (in lakhs)
2	3
6	10
5	4
7	13

Iteration 2

- Similar to iteration no. 1 performed above
 - we will again calculate Cost function and update θ_j values using Gradient Descent.
 - We will keep on iterating until Cost function doesn't reduce further (**converge**). At that point, model achieves best θ values.
 - Using these θ values in the model hypothesis will give the best prediction results.

Training II

Testing ←

- If we have the experience years and the hypotheses (h_θ), and you need to predict the salaries of these persons.

Experiences:

$$\boxed{-7 \ 3} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 21 & 31 & 45 & 11 \end{bmatrix} = h_\theta(x) = -7 + 3x$$

$[56 \ 86 \ 128 \ 26]$

21
31
45
11

Matrix **4x2**:

$$\begin{bmatrix} 1 & 21 \\ 1 & 31 \\ 1 & 45 \\ 1 & 11 \end{bmatrix} \times \begin{bmatrix} -7 \\ 3 \end{bmatrix} =$$

Prediction **4x1 matrix**: $h_\theta(21)$

$$\begin{bmatrix} -7 \times 1 + 3 \times 21 \\ -7 \times 1 + 3 \times 31 \\ -7 \times 1 + 3 \times 45 \\ -7 \times 1 + 3 \times 11 \end{bmatrix} = \begin{bmatrix} 56 \\ 86 \\ 128 \\ 26 \end{bmatrix}$$

2x1

Testing

- If we have the experiences and more than one hypotheses (h_θ):

$$h1_\theta(x) = -7 + 3x, h2_\theta(x) = 20 + x, h3_\theta(x) = -15 + 4x$$

available data points

Matrix **4x2**:

$3 \times 4 \leftarrow$ output

$$\begin{bmatrix} 1 & 21 \\ 1 & 31 \\ 1 & 45 \\ 1 & 11 \end{bmatrix} \times \begin{bmatrix} -7 & 20 & -15 \\ 3 & 1 & 4 \end{bmatrix} =$$

2x3

Prediction **4x3 matrix**:

$$\begin{bmatrix} 56 & 41 & 69 \\ 86 & 51 & 109 \\ 128 & 65 & 165 \\ 26 & 31 & 29 \end{bmatrix}$$

$h1_\theta$

Regression Problem: Multivariate Linear Regression

Linear Regression with multiple Variable

→ multiple Variables

more than one feature

Multiple variables (Features)

- Linear regression with multiple variables is also known as "**multivariate** linear regression".

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

<i>Size (feet)²</i> <i>x</i>	<i>Price (\$1000)</i> <i>y</i>
2104	460
1416	232
1534	315
852	178

Could the prediction be more accurate if we add #of rooms?



Multiple variables (Features)

<i>Size (feet)²</i>	<i>Number of bedrooms</i>	<i>Number of floors</i>	<i>Age of home (years)</i>	<i>Price (\$1000)</i>
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, number of training examples

n, number of features

$x^{(i)}$, the input (features) of the i^{th} training example

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

Row

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

Column

$$x_4^{(3)} = 30$$

Hypothesis

With one feature

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

Size (feet) ²	Number of bedrooms	Number of floors	Age of home (years)	Price (\$1000)
$\uparrow x_1$	$\uparrow x_2$	$\uparrow x_3$	$\downarrow x_4$	y
2104	5	1	45	460
1416	3	2	40	232
1534	3	2	30	315
852	2	1	36	178
...

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

With 3 θ_j instead of θ_0

$$h_{\theta}(x) = 80 - 0.1x_1 - 0.01x_2 - 3x_3 - 2x_4$$

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

For convenience of notation, define $x_0 = 1$

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

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

Hypothesis for Multiple Features

$$x = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

$$\theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \theta_3 \\ \vdots \\ \theta_n \end{bmatrix} \in \mathbb{R}^{n+1}$$

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

$$h_{\theta} = \sum_{j=0}^n \theta_j x_j = \theta^T x = [\theta_0 \quad \theta_1 \dots \quad \theta_n] \begin{bmatrix} x_0 \\ x_1 \\ \vdots \\ x_n \end{bmatrix}$$

Handwritten note: $\sum w_j x_j$

Cost Function for Multiple Variables

$$x_0 = 1$$

← always

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

θ is n + 1 dimension vector
From 1 ↘ *Hypothesis* ↗ *TVL*
Hypothesis function ↗ *Final result*

$$= \frac{1}{2m} \sum_{i=1}^m (\theta^T x^{(i)} - y^{(i)})^2$$
$$= \frac{1}{2m} \sum_{i=1}^m \left(\left(\sum_{j=0}^n \theta_j x_j^{(i)} \right) - y^{(i)} \right)^2$$

GD for Multiple Variables

- The gradient descent equation itself is generally the same form; we just have to repeat it for our ' n ' features

Repeat until convergence:{

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_j} J(\theta), \quad \text{simultaneously updates for every } j = 0, 1, 2, \dots, n$$

}

For $n=1$: one feature

Repeat {

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

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

(simultaneously update θ_0, θ_1)

}

For ($n > 1$):

Repeat until convergence:{

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

$\theta_1 := \dots$
 $\theta_2 := \dots$
 $\theta_3 := \dots$

, simultaneously updates θ_j
for every $j = 0, 1, 2, \dots, n$

}

(1) $\theta_0 := \dots$
 $x_0 \rightarrow \bar{x}_0$, $\theta_1 := \dots$
 $x_1 \rightarrow \bar{x}_1$

Gradient Descent in Practice I - Feature Scaling

- We can speed up gradient descent by having each of our input values in roughly the same range.
- Because θ will:
 - descend quickly on small ranges
 - descend slowly on large ranges, and
 - oscillate inefficiently down to the optimum when the variables are very uneven.
- The way to prevent this is to modify the ranges of our input variables so that they are all roughly the same.

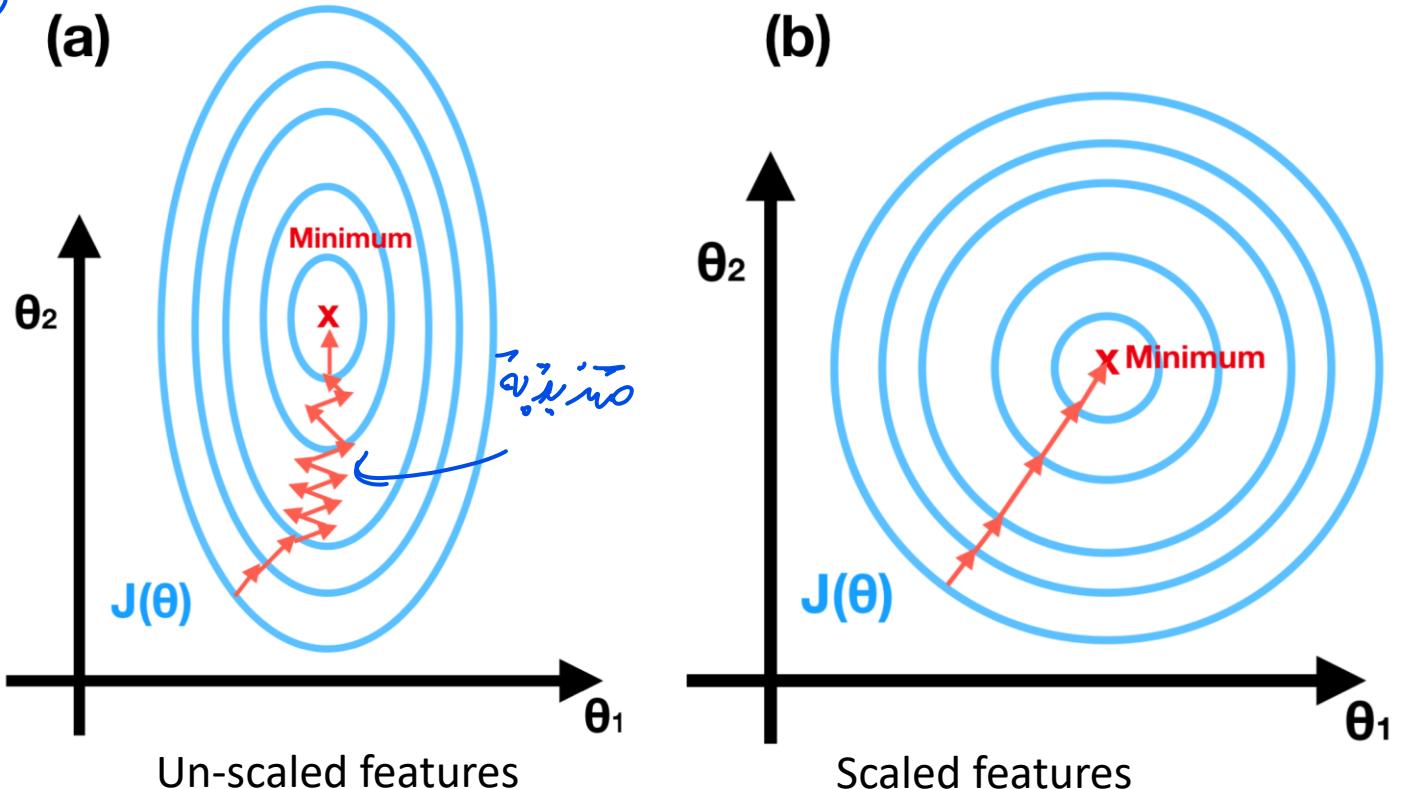
Normalization

Gradient Descent in Practice I - Feature Scaling

Make sure features are on same scale using:

- Mean Normalization
- Z-Score Normalization

(ML algorithm) will work better



Gradient Descent in Practice I - Feature Scaling

- Get every feature into approximately $-1 \leq x_j \leq 1$ range

(Do not apply to $x_0 = 1$)

$\{x_0\} \text{ and } (f_{\theta}(x)) \text{ do not change}$

$-1 \leq x_j \leq 1$ range
normalization
why? ω

Rule of thumb regarding acceptable ranges

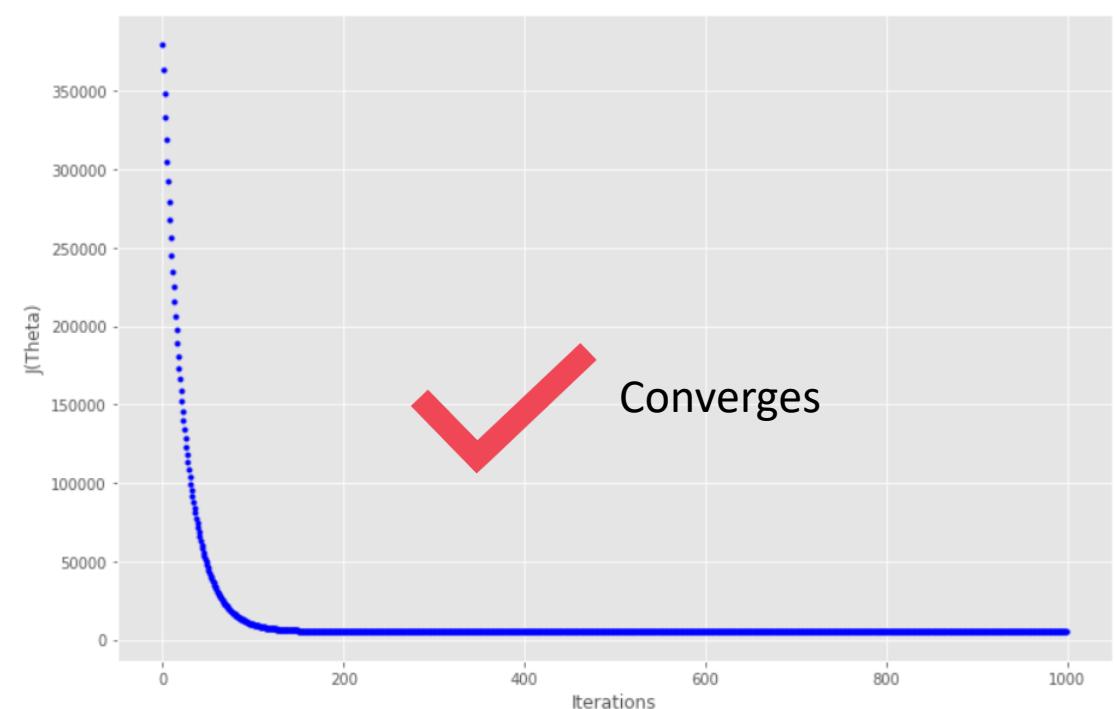
- 3 to +3 is generally fine - any bigger bad
- 1/3 to +1/3 is ok - any smaller bad

Gradient Descent in Practice II - Learning Rate



- Debugging gradient descent.
- Make a plot with number of iterations of gradient descent on the x – axis and the cost function $J(\theta)$ on the y – axis.
- $J(\theta)$ should decrease after each iteration else decrease α

Automatic convergence test. Declare convergence if $J(\theta)$ decreases by less than ε in one iteration, where ε is some small value such as 10^{-3} . However in practice it's difficult to choose this threshold value.

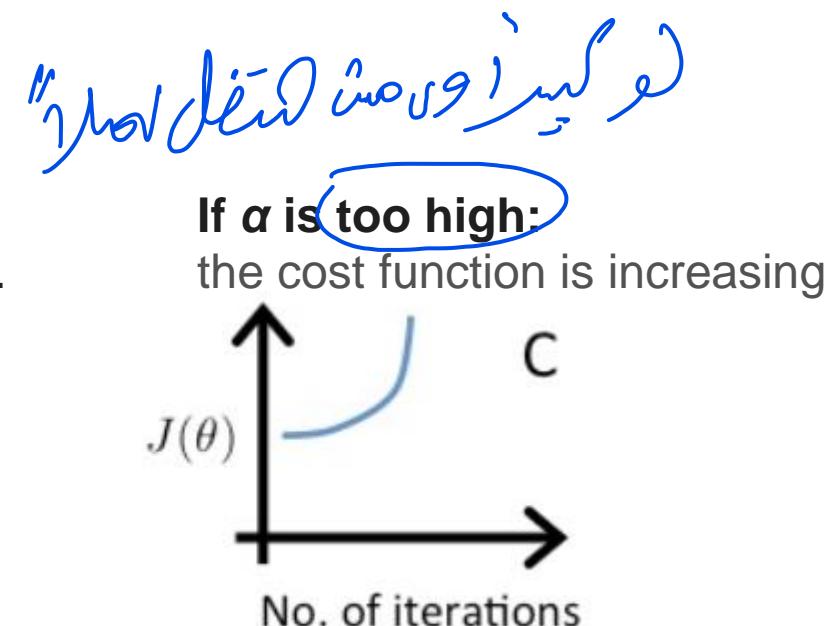
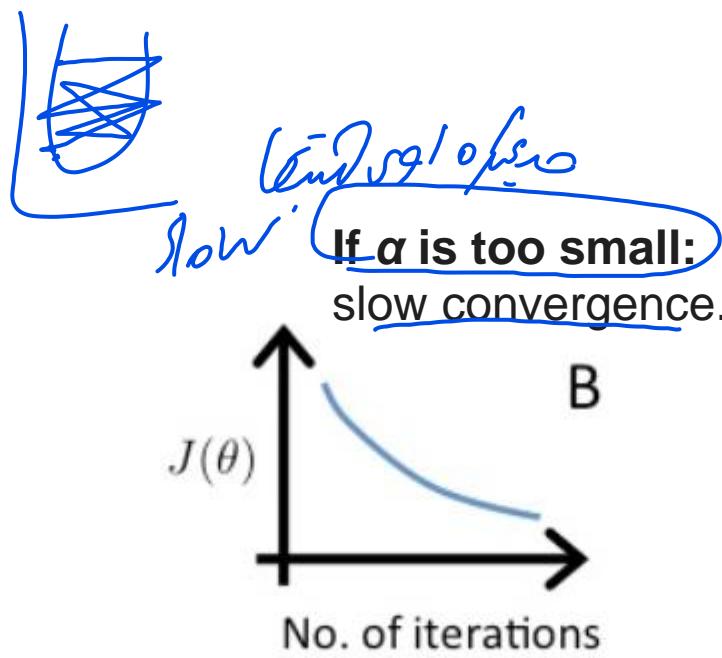
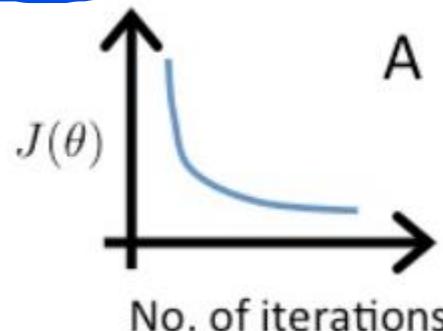


Gradient Descent in Practice II - Learning Rate

It has been proven that if learning rate α is sufficiently small, then $J(\theta)$ will decrease on every iteration. When Gradient Descent can't decrease the cost-function anymore and remains more or less on the same level, we say it has converged.

Note:

If you see in the plot that your learning curve is just going up and down, without really reaching a lower point, you also should try to decrease the learning rate.



Features choice /selection

- We can improve our features and the form of our hypothesis function.
- We can **combine** multiple features into one. For example, we can combine x_1 and x_2 into a new feature x_3 by taking $x_1 \cdot x_2$
Ex: x_1 is length and x_2 is width of a house, we can combine them into a new feature area= length x width

view feature into features interaction

Detecting Redundancies using Correlation Analysis

Vectors (Features)

- Consider two attributes A and B coming from two different data sources. The correlation between attributes A and B can be measured by the correlation coefficient $r_{A,B}$:

(This) relation between A and B , mean \bar{A} and \bar{B}

$$\text{Redundancies} \leftarrow r_{A,B} = \sum_{i=1}^n (A_i - \bar{A})(B_i - \bar{B}) / \{(n-1)\sigma_A\sigma_B\}$$

Where

$$\bar{A} = \left(\frac{\sum_{i=1}^n A_i}{n} \right)$$

$$\bar{B} = \left(\frac{\sum_{i=1}^n B_i}{n} \right)$$

$$\sigma_A = \sqrt{\frac{\sum_{i=1}^n (A_i - \bar{A})^2}{n-1}}$$

$$\text{and } \sigma_B = \sqrt{\frac{\sum_{i=1}^n (B_i - \bar{B})^2}{n-1}}$$

- Where n is the number of data items in A or B (assuming they have the same number of data items), \bar{A} and \bar{B} are the mean values of A and B , and σ_A and σ_B are the standard deviations of A and B

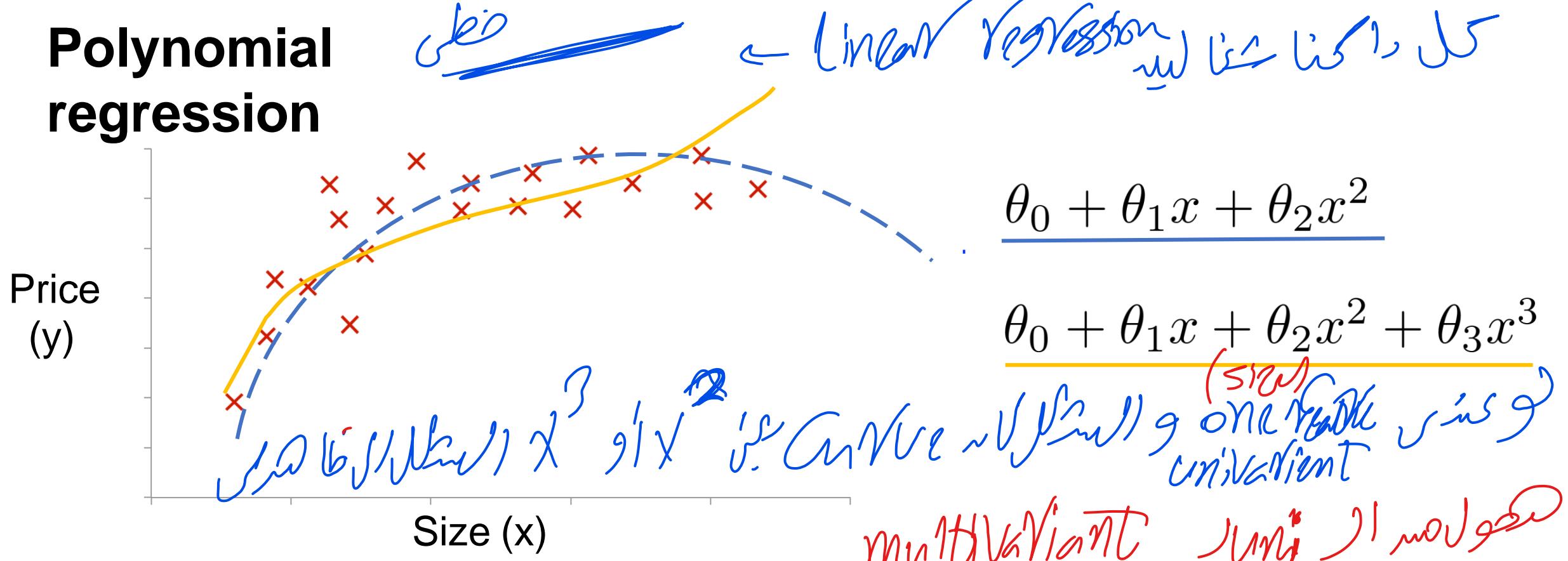
Detecting Redundancies using Correlation Analysis

- If the $r_{A,B} > 0$ → A and B are positively correlated
(دوسرا جھنڈا تک بڑے) *وہی تو*
→ the values of A increase as the values of B increase
→ **either A or B** may be removed as a redundancy but not both of them

- If the $r_{A,B} = 0$ → A and B are independent (no redundancy)
→ we cannot remove any of them

- If the $r_{A,B} < 0$ → A and B are negatively correlated
→ the values of one attribute increase as the values of the other decrease
وہی تو
→ we cannot remove any of the attributes

Polynomial regression



$$\begin{aligned} h_{\theta}(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2 + \theta_3(\text{size})^3 \end{aligned}$$

$x_1 = (\text{size})$

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

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

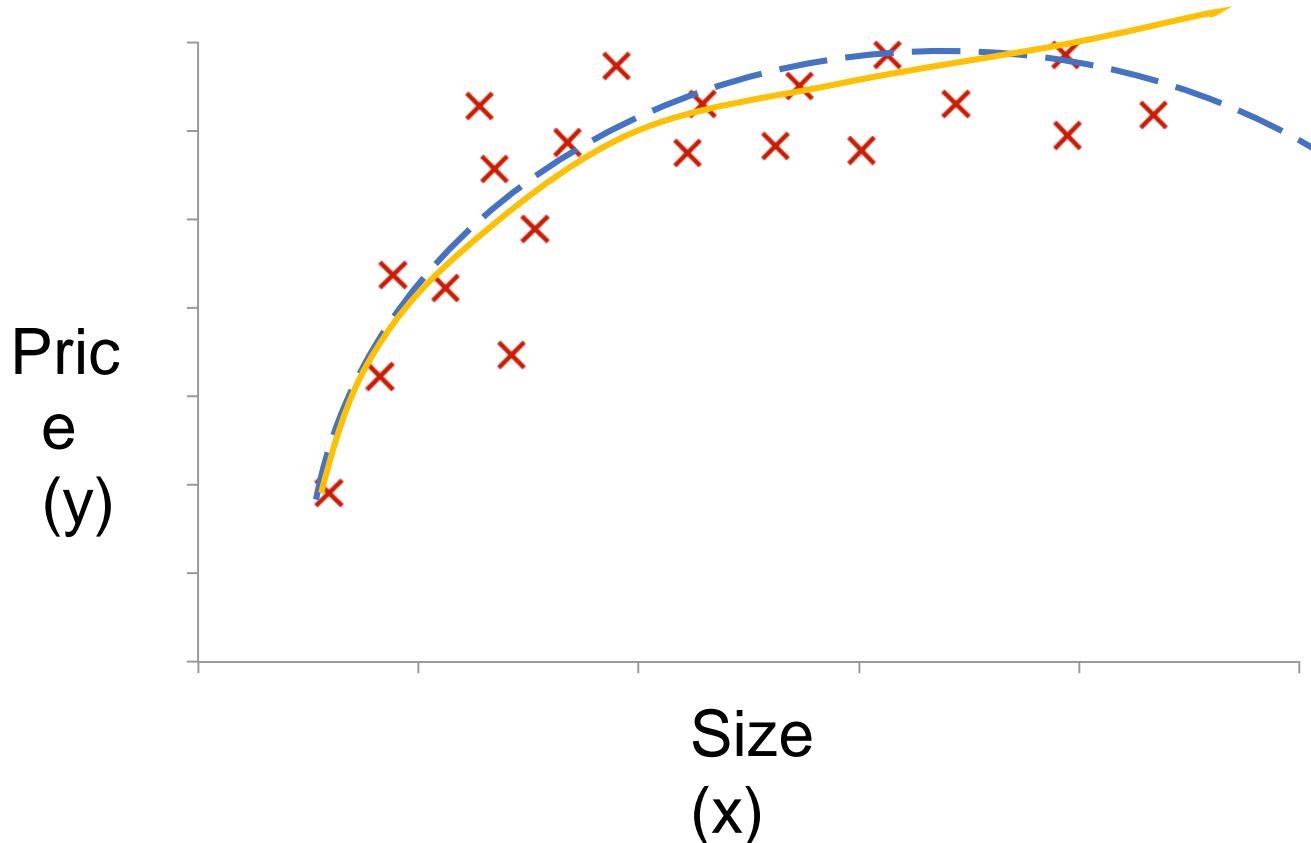
multivariate *feature scaling becomes very important.*

Size: 1- 1000

Size²: 1- 10⁶

Size³: 1- 10⁹

Choice of features



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

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

Optimization Technique: Normal Equations

Better values for θ

minimize the error

- To solve for θ analytically \rightarrow **normal equation**
- $\theta \in \mathbb{R}^{n+1}$

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

Minimize $J(\theta)$

$$\frac{\partial}{\partial \theta_j} J(\theta) = \dots = 0 \text{ (set to zero)}$$

for every j solve for $\theta_0, \theta_1, \dots, \theta_n$

Example: m=4

x_0	<i>Size (feet)²</i>	<i>Number of bedrooms</i>	<i>Number of floors</i>	<i>Age of home (years)</i>	<i>Price (\$1000)</i>
	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

(Features) Matrix \mathbf{X} $\in \mathbb{R}^{m \times n}$

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

Features matrix $m \times (n + 1)$

Output vector $\mathbf{y} \in \mathbb{R}^m$

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

$m - \text{dim vector}$

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

↳ has high complexity

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

m samples , n features

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

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

Design matrix
 $m \times (n+1)$

Example: one feature

$$x^{(i)} = \begin{bmatrix} 1 \\ x_1^{(i)} \end{bmatrix}$$

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

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

Gradient Descent Vs. Normal Equation

affiliation)) goes with this slide
and Jiggy

Gradient Descent	Normal Equation
Need to choose α	No need to choose α
Needs many iterations	No need to iterate
$O(kn^2)$ <i>slow</i>	$O(n^3)$ <i>high complexity</i>
Works well when n is large	Slow if n is very large

also
need α , many iterations

slow, high complexity

Example 2:

x_0	Age (x1)	Height in cm (x2)	Weight in Kg (y)
1	4	89	16
1	9	124	28
1	5	103	20

- What is X (design matrix) and y?

$$X = \begin{bmatrix} 1 & 4 & 89 \\ 1 & 9 & 124 \\ 1 & 5 & 103 \end{bmatrix}$$

$$y = \begin{bmatrix} 16 \\ 28 \\ 20 \end{bmatrix}$$

x_0 \hat{x}_1 \hat{x}_2

Check?

$$\hat{Y} = \theta_0 X_0 + \theta_1 X_1 + \dots + \theta_6 X_6 \quad \leftarrow \# \theta = 7$$

- Suppose you have $m = 25$ training examples with $n = 6$ features. The normal equation is $\theta = (X^T X)^{-1} X^T y$
 - For the given values of m and n , what are the dimensions of θ , X , and y in this equation?

$$x_0 \ x_1 \ x_2 \ \cdots \ x_6$$

$$\max(n+1)$$

951

$$= 25 \times 1$$

$$\theta = f_a$$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n$$

- Num of Searches \uparrow
 $\Rightarrow (n+1)$

Thanks