# Andrew Ng's

Supervised Learning - "Given right answer"

## Regression problem

| Training Set |

↓

Learning Algo

Size of house → $h$ → Estimated price

hypothesis

↳ $h$ maps from x's to y's

fitting line to data

↳ **Cost function:**

How to choose θ's?



$\theta_0, \theta_1$ → best to fit?

Idea: Choose $\theta_0$ $\theta_1$ so that $h_\theta(x)$ is close to y for our train. data $(x,y)$

## Notation

m - no. of examples (train)

x, y

$(x^i, y^i)$ → $i^{th}$ example

How to represent h?

$h_\theta(x) = $ ⓪$\theta_0$ + ⓪$\theta_1$ x ← parameters



$h(x) = \theta_0 + \theta_1 x$

Linear Regression with one variable
or
Univariate regression

$$\underset{\theta_0 \theta_1}{\text{minimize}} \frac{1}{2m}\sum_{i=1}^{m}\left(h_\theta(x^i) - y^i\right)^2$$

↳

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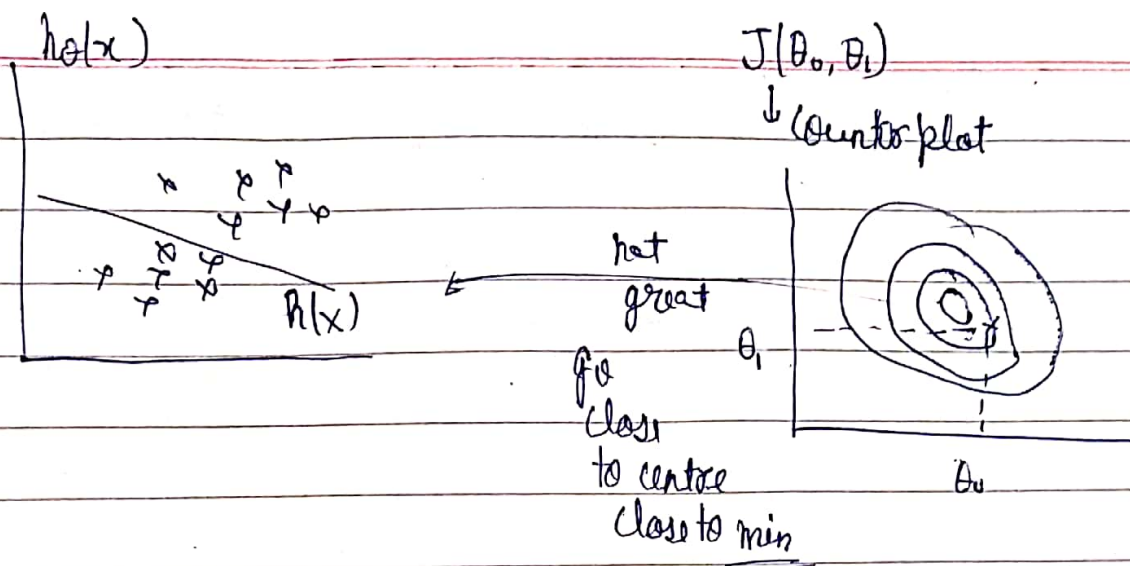
min. $J(\theta_0, \theta_1)$

Cost fn.
or
Squared error function



for $\theta_1 = 1$, $J(\theta) = 0$ best

$\theta_1 = 0.5$ $J(\theta) = 0.58$

$\boxed{\theta_1 = 1}$ → min $J(\theta_1)$

$h_\theta(x)$

$J(\theta_0, \theta_1)$
↓ Counter plot



$h(x)$

not
great

$\theta_1$

fo
closr
to centre
close to min

$\theta_0$

Algorithm to automatically find
values of $\theta_0, \theta_1$ to min $J(\theta_0, \theta_1)$

\# __Gradient Descent__ (to minimize cost fn)

- Have some function $J(\theta_0, \theta_1)$
- What min $J(\theta_0, \theta_1)$
- Outline:-
  Start with some $(\theta_0, \theta_1)$
  Keep changing $(\theta_0, \theta_1)$ to reduce $J(\theta_0, \theta_1)$
  until end up at minimum.

__algo:-__

Repeat until converge {

$$\theta_j := \theta_j - \alpha \boxed{\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)} \quad (\text{for } j = 0 \text{ and } j = 1)$$

}

learning rate (steps)     → derivative term (later)

| Correct: (Simultaneous update) | Incorrect |
|---|---|
| $temp0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1)$ | $temp 0 := \_\_\_\_\_$ |
| $temp1 := \theta_1 - \alpha \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1)$ | $(\theta_0:) = temp 0$ |
| $\theta_0 := temp 0$ | $temp 1 := \_\_\_ J(\theta_0, \theta_1)$ |
| $\theta_1 := temp1$ | $\theta_1 := temp1$  New value |

$$\theta_i = \theta_i - \alpha \frac{\partial}{\partial \theta_i} J(\theta_0, \theta_1) \qquad i = 0 \text{ and } 1$$
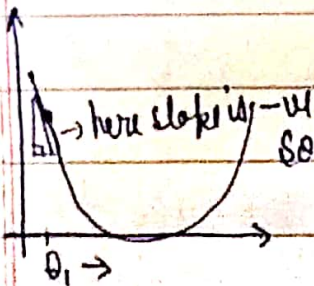
$J(\theta_1)$

$$\theta_i = \theta_1 - \alpha \frac{\partial}{\partial \theta} J(\theta_1)$$

what is slope

of line at this value of $\theta_1$

In this case slope is +ve

so

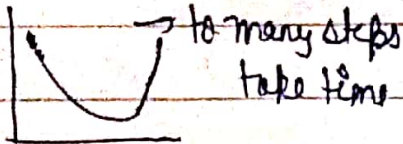$$\theta_1 = \theta_1 - \alpha (\text{positive value})$$

so $\theta_1$ decreases, close to min

→ here slope is −ve

so

$$\theta_i = \theta_1 - \alpha (\text{negative})$$

$\theta_1$ increases

← close to min.

$\theta_1 \rightarrow$

If $\alpha$ is too slow (small)

→ too many steps
take time

If $\alpha$ is too large, overshoot

never converge

$\theta_1$

slope = 0

$\theta_i$ local optima

$\theta_1$

$$\theta_i = \theta_1 - \alpha (0)$$
$$\theta_i = \theta_1 \text{ unchanged.}$$

As we approaches min
slope is less so change in $\theta_1$ less

∴ smaller steps

Gradient descent with cost function:-

G.d.

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

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

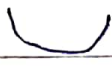$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^{m} (h_\theta x^i - y^i)^2$$

$$\frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum (h_\theta x^i - y^i)^2$$

or

$$\frac{\partial}{\partial \theta_j} \frac{1}{2m} \sum^m (\theta_0 + \theta_1 x^i - y^i)^2$$

for $\theta_0$

$$j = 0 : \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{1}{m} \sum (h_\theta(x)^i - y^i)$$

for $\theta_1$

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

Now,

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

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

] update simultaneously

Cost function for linear regression is always bow shaped $\smile$ means no local optima, converges at global optima.

"Batch Gradient Descent", the one we just done.
↳ "Batch" → Each step of G.d. uses all training examples.

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

## Normal Distribution :-

$$Eq^n$$

In 4D



Method to
Solve for $\theta$
analytically
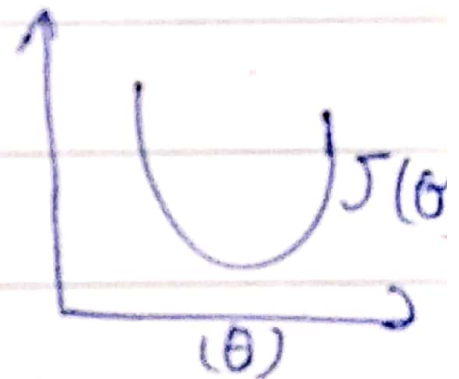
## Intuition ($\theta \in \mathbb{R}$)

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

To min. qued. $eq^n$

$$\frac{\partial}{\partial \theta} J(\theta) = 0$$



Solve for ($\theta$)

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

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

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

|  | G.d | Normal eqⁿ |
|---|---|---|
| | • Need to choose $\alpha$ | • No need |
| | • Need many epoch | • No need to iterate |
| | • Works well even when $n$ is large | • Need to compute $(X^TX)^{-1} \to n\times n$ |
| | | • Slow if $n$ is very large $O(n^3)$ |
| | | if $n > 1000$ |

## Example

$$m = 4$$

| Size | No. of bed | floor | Age | Price |
|---|---|---|---|---|
| 2104 | 5 | 1 | 45 | 460 |
| 1416 | 3 | 2 | 40 | 232 |
| 1534 | 3 | 2 | 30 | 315 |
| 852 | 2 | 1 | 36 | 178 |

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

$m \times (n+1)$
rows, features

m-d vector

$$\boxed{\theta = (X^TX)^{-1} - X^T y} \longrightarrow \text{Optimal value of } \theta \text{ that min. } J(\theta)$$

$$x^i = \begin{bmatrix} x_0^i \\ x_1^i \\ x_2^i \\ \vdots \\ x_n^i \end{bmatrix} \longrightarrow X = \begin{bmatrix} - (x^1)^T - \\ - (x^2)^T - \\ \vdots \\ - (x^n)^T - \end{bmatrix}$$

design matrix

$m \times n+1$

m examples
n features

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

$$X = \begin{bmatrix} 1 & x^{(1)} \\ 1 & x_2^{(1)} \\ -- & -- \end{bmatrix} \qquad y = \begin{bmatrix} y^1 \\ y^2 \\ \vdots \\ y^n \end{bmatrix}$$

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

In normal equation method, no need to do feature scaling.

# Normal Eqⁿ Noninvertibility :-

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

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

## In Octave

$$pinv(X' * X) * X' * y$$

{ pinv
→ inv

## Cause of invertibility ?

- Redundant features

Eg:- $x_1 =$ size in $feet^2$

$x_2 =$ size in $m^2$

$x_1 = (3.28)^2 \times x_2$

- Too many features (eg: $m \leq n$)
  - Delete some features or use regularization.