

## Revision

① 

$y = \theta_0 + \theta_1 x$

gradient descent:

hypothosis  
simplified

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

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

$y = \theta_0 + \theta_1 x$

$\theta_0$ : offset / intercept  
 $\theta_1$ : slope / rate of  $x$

target value / feature

Input var |

Training set

Testing set

cost func J

$y = \theta_0 + \theta_1 x$

$(x^{(i)}, y^{(i)}) \rightarrow$   $i^{\text{th}}$  training example

$(x^{(i)}, y^{(i)}) \rightarrow$  one training example

$\theta_0, \theta_1$

goal  $\rightarrow$  minimize  $J(\theta_0, \theta_1)$

$J(\theta_0, \theta_1)$

$$= \frac{1}{m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2 = \frac{1}{m} \sum_{i=1}^m h(x^{(i)}) - y^{(i)})^2$$

goal  $\rightarrow$  minimize  $J(\theta_0, \theta_1)$

nos of features

nos of records total nos of examples

★ hypothesis:  $h(x) = \theta_0 + \theta_1 x$

(for fixed  $\theta_1$ , this is a func of  $x$ )

$$\theta_1 = 0.5 \quad J(\theta_0, \theta_1)$$



$$\theta_1 = 0 \quad J(\theta_0, \theta_1)$$

No error

parameter estimation  $\rightarrow$  OLS or ordinary least sq.  
 of  $\theta_1$   $\rightarrow$  gradient descent  
 (optimization appro.)

① OLS J

$\hat{y} = \theta_0 + \theta_1 x$

$$\theta_1 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\sum_{i=1}^n (x_i - \bar{x})^2$$

$$\theta_0 = \bar{y} - \theta_1 \bar{x}$$

$$J(\theta_1) = \frac{1}{2} \sum_{i=1}^n [(1.5x_i - 1)^2 + (1.5x_i - 3)^2]$$

(func of the parameter  $\theta_1$ )

$$= 1.5x_1^2 - 2x_1 + 1.5x_2^2 - 3x_2$$

$$= 1.5x_1^2 - 2x_1 + 1.5x_2^2 - 3x_2$$



$$J(\theta_1) = \frac{1}{2} \sum_{i=1}^n [(1.5x_i - 1)^2 + (1.5x_i - 3)^2]$$

gradient descent intuition

J(Theta) (Theta R)

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n [y_i - (\theta_0 + \theta_1 x_i)]^2$$

If  $\alpha$  is too large gradient descent can overshoot the minimum

Diagram B shows a horizontal rectangle with a diagonal line from top-left to bottom-right. A curved arrow starts at the top-left corner, goes down along the left edge, then turns right along the diagonal line.

$\theta_1$   $\rightarrow \theta_1 : \theta_1 - \alpha$  (neg numb.)

Current value of  $\theta_1$

$$\theta_{L2} = \theta_L - \alpha \frac{d}{d\theta_L} J(\theta_L)$$

<u>gradient descent algo</u>	<u>linear regression model</u>
------------------------------	--------------------------------

repeat until convergence  
 $\theta_j^* = \frac{\partial L}{\partial \theta_j} (t(\theta_0, \theta_1))$  (simultaneously)  
 $\theta_0^* = 0 \quad \theta_1^* = 1$

repeat until convergence 3.  $h(x) = \theta_0 + \theta_1 x$

$$\frac{\partial \theta_j^o}{\partial \theta_1} = \theta_j^o - \alpha \sum_{i=1}^m (h_{\theta_0, \theta_1}(x^{(i)}) - y^{(i)})$$

9

Correct & Simultaneous update

$$f(180^\circ, \theta) \rightarrow 1$$

$$f_{\theta}(\cdot | \mathcal{D}_{t-1}, \theta_t) \rightarrow \frac{1}{Z} \sum_m \exp(\theta_m^T \mathbf{x}_t)$$

$$\text{temp}_0^{\circ} = \theta_0 - \alpha \int_{\theta_0}^{\theta} I(\theta, \theta_1) d\theta$$

$\theta_0 = \text{temp}_0$

$$\frac{\partial}{\partial \theta_2} J(\theta_0, \theta_1) \Rightarrow \frac{1}{2m} \sum_{i=1}^m 2 [ \theta_0 + \theta_1 x^{(i)} - y^{(i)} ]^2$$

$$\theta_{\text{obs}} = \text{temp}_0 - \frac{1}{\rho_0}$$

$$\theta_{1,i} = \text{Temp}_1$$

If  $\alpha$  is too small gradient descent can be slow

gradient descent algo.

repeat until convergence:

$$\theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}\phi(x^{(i)}) - y^{(i)})$$

$$\theta_1 := \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}\phi(x^{(i)}) - y^{(i)}) \cdot x_1^{(i)}$$

$$\theta_2 := \theta_2 - \alpha \frac{1}{m} \sum_{i=1}^m (\text{h}\phi(x^{(i)}) - y^{(i)}) \cdot x_2^{(i)}$$

As we approach local minimum, gradient descent will automatically take smaller steps. So no need to decrease  $\alpha$  over time.

★ Regression with multiple features (variables)

$x_0$	$x_1$	$x_2$	$x_3$	$x_4$	$x_5$	$x_6$	$x_7$	$x_8$	$x_9$
1	0	0	0	0	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0
1	2	0	0	0	0	0	0	0	0
1	3	0	0	0	0	0	0	0	0
1	4	0	0	0	0	0	0	0	0
1	5	0	0	0	0	0	0	0	0
1	6	0	0	0	0	0	0	0	0
1	7	0	0	0	0	0	0	0	0
1	8	0	0	0	0	0	0	0	0
1	9	0	0	0	0	0	0	0	0

$$\text{h}\phi(x) = \theta_0 x_0 + \theta_1 x_1 + \dots + \theta_n x_n$$

$$(\theta_0 = 1)$$

$$\text{h}\phi(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

$$\text{h}\phi(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + 0.01 x_3 + 3 x_4 - 2 x_5 \rightarrow \text{age}$$

$$\text{h}\phi(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n \rightarrow \text{multiple variable (overfit)}$$

$$\text{h}\phi(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n \rightarrow \text{overfit}$$

hypothesis<sup>o</sup>

$$\text{h}\phi(x) = \theta_0 + \theta_1 x_1$$

$$\text{h}\phi(x) = \theta^T x$$

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

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

$$\theta^T = \begin{bmatrix} \theta_0 & \theta_1 & \dots & \theta_n \end{bmatrix}$$

$$(n+1) \times 1$$

matrix

- notations:  $n = \text{number of features}$
- $x_i^{(j)} = \text{input feature of } j^{\text{th}} \text{ training example}$
- $\theta_j^{(i)} = \text{value of feature } j \text{ in } i^{\text{th}} \text{ training example}$

$$x^{(2)} = \begin{bmatrix} 1416 \\ 3 \\ 2 \\ 1 \\ 36 \\ 170 \end{bmatrix} \quad x_3^{(2)} = 2$$

\* cost drop as every iter.

$$\text{Hypothesis: } h(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \dots$$

parameters:  $\theta_0, \theta_1, \dots, \theta_m$   
 $(n+1)$  dimensional vector

$$\text{Cost func: } J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2$$

gradient descent,

$$\text{repeat } \left\{ \begin{array}{l} \boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \frac{\partial}{\partial \boldsymbol{\theta}} J(\boldsymbol{\theta}) \end{array} \right.$$

Simultaneously update for every  $j=0, 1, \dots, m$

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

$$\Rightarrow \frac{\partial}{\partial \theta_0} \left[ \frac{1}{2m} \sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)})^2 \right]$$

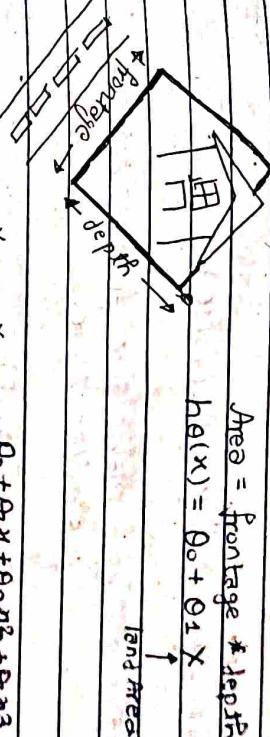
$$\frac{\partial}{\partial \theta_0} \left[ \frac{1}{2m} \sum_{i=1}^m ((\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_m x_m^{(i)}) - y^{(i)})^2 \right]$$

$$\frac{\partial}{\partial \theta_0} \left[ \frac{1}{2m} \sum_{i=1}^m ((\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_m x_m^{(i)}) - y^{(i)})^2 \right]$$

$$\frac{1}{m} \sum_{i=1}^m ((\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_m x_m^{(i)}) - y^{(i)})$$

$$\frac{1}{m} \sum_{i=1}^m [(\theta_0 x_0^{(i)} + \theta_1 x_1^{(i)} + \dots + \theta_m x_m^{(i)}) - y^{(i)}]$$

$$\theta_0 := \theta_0 - \alpha \frac{\partial}{\partial \theta_0} J(\boldsymbol{\theta})$$



$$h(x) = \theta_0 + \theta_1 \times \text{frontage} + \theta_2 \times \text{depth}$$

polynomial regression:

$$\text{need if house price prediction}$$

$$\theta_0 + \theta_1 - \alpha \sum_{i=1}^m (h(\mathbf{x}^{(i)}) - y^{(i)}) \mathbf{x}_i^{(i)}$$

$$S_{2e}(x)$$

$$h(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$\begin{aligned} S_{2e} &:= 1 - 1000 \\ &\rightarrow \theta_0 + \theta_1 (S_{2e}) + \theta_2 (S_{2e})^2 + \\ &\quad \theta_3 (S_{2e})^3 \\ &\quad S_{2e}^3 := 1 - 10^9 \end{aligned}$$

$$x_4 \rightarrow (S_{2e})$$

$$x_2 \rightarrow (S_{2e})^2$$

$$x_3 \rightarrow (S_{2e})^3$$

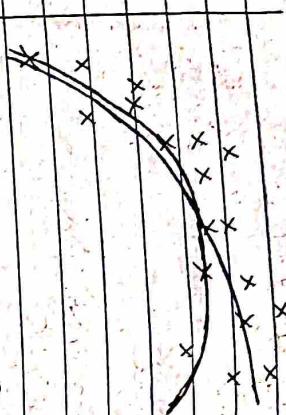
```
import pandas as pd
```

```
import numpy as np
```

```
df = pd.read_csv('placement.csv')
```

```
df = df.iloc[:, :-1]
```

	Cgpa	Package
0	6.89	3.26
1	5.12	1.98
2	3.02	3.25
3	7.42	3.67
4	6.94	3.53
5	0.00	0.00
6	6.22	2.33



$$h_0(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\text{size})^2.$$

$$h_0(x) = \theta_0 + \theta_1(\text{size}) + \theta_2(\sqrt{\text{size}})$$

feature scaling  $\circ$  makes sure features are on similar scale.

$X_1 \circ$  size (0-2000 feet<sup>2</sup>)  
 $X_2 \circ$  no. of bedrooms (1-5)

get every feature into approximately 2 range  $\circ$   $-1 \leq x_i \leq 1$

```
# feature scaling
def normalize_feature(feature):
    return (feature - np.mean(feature)) / np.std(feature)
```

```
= normalize_feature(df.iloc[:, 0:1].values)
```

```
= array([[-0.094213], [-1.75349666], [0.00000000], [1.253669871]]).
```

replace  $x_i$  with  $x_i - \mu$  to make features have approximately zero mean  
 $(\text{don't apply to } x_0 = 1)$

$$\frac{x_i - \bar{x}_i}{\text{range}} = \frac{x_i - 1000}{2000} \quad -0.5 \leq x_i \leq 0.5$$

$$y_2 = \text{bedrooms} - 2$$

$y$  avg. value of  $x_1$  in training set

$$x_1 \rightarrow x_1 - \bar{x}_1 \quad \text{range} \rightarrow x_2 - \bar{x}_2$$

Standard deviation

When we say "similar scale" we mean that the value of different features are within a comparable range. In the context of feature scaling, the goal is to make the values of each feature so that they proportionally comparable to each other. This ensures no single feature dominates the other due to differences in their magnitude.

### ① min-Max scaling (Normalization)

$$x_{\text{normalized}} = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$

range b/w 0 to 1

### ② Standardization (Z-score normalization)

$$\text{standardized } x = \frac{x - \mu}{\sigma}$$

This transforms the value to its mean ( $\mu$ ) and a standard deviation ( $\sigma$ ) of 1.

def hypothesis(theta0, theta1, x):  
 $h_\theta(x) = \theta_0 + \theta_1 x$

return theta0 + theta1 \* x

def cost\_function(theta0, theta1, x, y):

m = len(y)

$h = \text{hypothesis}(\theta_0, \theta_1, x)$

return (1 / (2 \* m)) \* np.sum((h - y) \*\* 2)

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

$(x_i, y_i) \rightarrow$  example of dataset  
 $(x^{(i)}, y^{(i)}) \rightarrow i^{\text{th}} \text{ example of dataset}$

def gradient\_descent(theta0, theta1, x, y,  
 learning\_rate=0.0001, iterations=1000,  
 convergence\_threshold=1e-4):

m = len(y)

costs = [] # to store cost value at each iteration

for iteration in range(Iterations):

$h = \text{hypothesis}(\theta_0, \theta_1, x)$

    delta-theta0 = (1/m) \* np.sum(h-y)

$$\frac{\partial J(\theta)}{\partial \theta_0} = \frac{1}{2m} \sum_{i=1}^m (h(x^{(i)}) - y^{(i)})^2$$
$$\Rightarrow \frac{1}{m} \sum_{i=1}^m [(h(x^{(i)}) - y^{(i)})]^2$$

    delta-theta1 = (1/m) \* np.sum((h-y)\*x)

$$\frac{\partial J(\theta)}{\partial \theta_1} = \frac{1}{m} \sum_{i=1}^m [x_i (h(x^{(i)}) - y^{(i)})]$$

    theta0 = learning-rate \* delta-theta0

    theta1 = learning-rate \* delta-theta1

repeat until convergence  $\xi$

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \quad (\text{simultaneously update } j=0 \text{ & } j=1)$$

# check for convergence based on gradients

if np.abs(delta-theta0) < convergence-threshold

and np.abs(delta-theta1) < convergence-threshold:

    print("Converged after  $\xi$  iterations")

Break # when meet minima

# Store cost for every iteration

cost = cost-function(theta0, theta1, x, y)

costs.append(cost)

# print cost for every 100 iteration

if iteration % 100 == 0:

    print(f"iteration {iteration}, cost: {cost}")

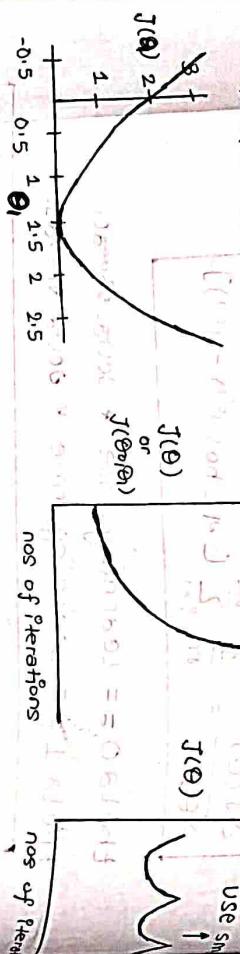
return theta0, theta1, costs

- Convergence criteria:
  - for convex functions, optimum occurs when  $\frac{\partial J(\theta)}{\partial \theta} = 0$

In practice, stop when  $\frac{\partial J(\theta)}{\partial \theta} \leq \epsilon \rightarrow 1e^{-3}, 1e^{-4}, 1e^{-5}$

As we approach a local minimum, gradient descent will automatically take smaller steps. So, we need to decrease ( $\alpha$ ) over time.

make sure gradient descent is working correctly



overshoot  
for sufficiently small  $\alpha$ ,  
should decrease on every iteration  
(large)

But, if  $\alpha$  is too small, gradient

descent can be slow to converge

# Initialize theta\_0 and theta\_1

theta\_0 = 1.0 or 0.01 (any value between 0 and 1)

theta\_1 = 0.1 or 0.01 (any value between 0 and 1)

# Run gradient-descent

final\_theta\_0, final\_theta\_1, costs =

gradient\_descent(theta\_0, theta\_1, alpha)

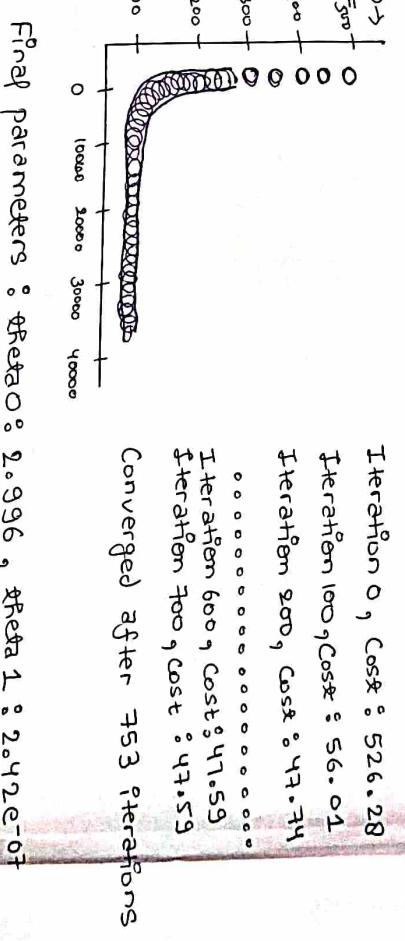
plotting the cost  $J(\theta)$  over iterations.

```
operations_range = range(0, len(costs)) * 100, 100)
plt.plot(operations_range, costs, marker='o')
```

```
plt.title('Cost over Iterations')
plt.xlabel('Iterations')
plt.ylabel('cost')
```

```
# print the final parameters
print(f'Final parameters : theta_0 = {final_theta_0},\ntheta_1 = {final_theta_1}')
```

```
Iteration 0, cost : 526.28
Iteration 100, cost : 56.01
Iteration 200, cost : 47.34
...
Iteration 600, cost : 41.59
Iteration 700, cost : 47.59
Converged after 753 iterations
```



$X = np.array([[1, 2, 3], [1, 3, 4], [1, 4, 5]])$   
 $\downarrow$   
 3x3 dimension  
 feature matrix

$x_1$	$x_2$	$x_3$	$y$
1	2	3	1
1	3	4	2
1	4	5	3

$\varphi = np.array([1, 2, 3])$

$$\begin{bmatrix} \varphi \\ \downarrow \end{bmatrix} \quad \text{(target variable)}$$

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

$$\text{hypothesis} \quad h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 + \dots + \theta_n x_n$$

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

$x^{(i)} : \text{input feature}$   
 of  $i^{\text{th}}$  training example

$$x^{(2)} = \begin{bmatrix} 1 \\ 3 \\ 4 \end{bmatrix}$$

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

for convenience of notation define  $x_0 = 1$  or  $x_0^{(i)} = 1$

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$

$$X = \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n+1} \quad \# \text{ add column of ones to array 'x'}$$

$$X = np.c_[np.ones(X.shape[0]), X]$$

- $np.ones(X.shape[0])$  : creates an array of ones with the same nos of rows as the original array 'X'.

$X.shape[0]$  represents nos of rows in array 'X'.

- $np.c_$  : concatenate array of ones with original array X.

$\Rightarrow$   $array([1, 1, 2, 3], [1, 2, 3, 4], [1, 1, 4, 5])$

$$X^2 = \begin{bmatrix} 1 \\ 1 \\ 3 \\ 4 \end{bmatrix}$$

$\circlearrowleft$   
 Now  
 $\downarrow$   
 3x4 dimensional

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

$$h_{\theta}(x) = \theta^T x$$

$$\begin{bmatrix} \theta_0 & \theta_1 & \theta_2 & \dots & \theta_n \end{bmatrix} \begin{bmatrix} x_0 \\ x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix}$$

```
def compute_cost(theta):
    repeat {
        J = C1 / C2 * np.sum((h(theta) - y) ** 2)
        return J.
```

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

simultaneously update for every  $\theta_j = 0, 1, 2, \dots, n$ .

$$\begin{aligned} m &= len(y) \\ r &= np.dot(X, theta) \\ (n+1) \times 1 & \text{matrix} \\ h &= np.dot(X, theta) \\ \theta_j &:= \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \end{aligned}$$

$$r = np.dot(X, theta)$$

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

$$\begin{aligned} \text{off} &\rightarrow \text{array}[[10, 10, 20, 30], \\ &[10, 10, 30, 40], \\ &[10, 10, 40, 50]] \end{aligned}$$

$$\text{cost func: } J(\theta_0, \theta_1, \theta_2, \dots, \theta_m) = \frac{1}{2m} \sum_{i=1}^m (h(\theta)x^{(i)}) - y^{(i)})^2$$

$$(n+1) \text{ dimensional vector}$$

$$J(\theta)$$

$$\theta = np.zeros(X.shape[1])$$

- The length of  $\theta$  is set to the numbers of columns in updated matrix  $X$ :

```
for i in range(num_iterations):
    r = np.dot(X, theta)
```

```
    errors = r - y
```

```
    gradient = (1/m) * np.dot(X.T, errors)
```

```
    theta = theta - alpha * gradient
```

```
J = compute_cost(X, y, theta)
```

```
J_history.append(J)
```

If print-cost and i% print-every == 0:

print f "cost after iteration {i}

return theta, J-history

h = np.dot(x, theta)  $\xrightarrow{\text{initially} = 0}$  array([0, 0, 0])

errors = h - y  $\xrightarrow{\text{array}} [-1, -2, -3]$

gradient = (1/m) \* np.dot(x.T, errors)

initially  $x_0 \xrightarrow{\begin{bmatrix} 1 \\ 1 \\ 1 \\ 0 \end{bmatrix} [-1 -2 -3]} \leftarrow \theta_0$

as in loop we get others  $\theta_1^T \theta_4, \theta_2^T \theta_2$  and many

model performance evaluation:

Loss function  $L(y, h_\theta(x))$

$L$  (Actual Value, Predicted Value)

Examples:

(assuming loss for underpredicting = overpredicting)

Absolute Error  $\rightarrow L(y, h_\theta(x))$

$$\Rightarrow |y - h_\theta(x)|$$

Mean Absolute Error  $\rightarrow L(y, h_\theta(x))$

$$\Rightarrow (y - h_\theta(x))^2$$

### \* Training Error Computation:

① Define a Loss function  $L(y, h_\theta(x))$

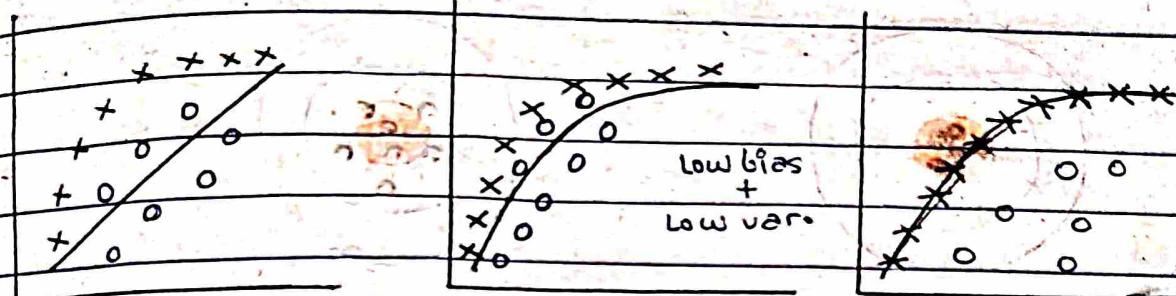
e.g. squared error, absolute error

② Training error = avg. loss on house in training set

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n L(y^{(i)}, h_\theta(x^{(i)}))$$

## Bias - Variance - Tradeoff

low testing error



Under-fitting

Appropriate-fitting

Over-fitting

$$\theta_0 + \theta_1 x_1$$

$$\theta_0 + \theta_1 x_1 + \theta_2 x_1^2$$

x : training data

o : testing data

- Training error is high (High bias)
- Training error is less (Low biased)
- Training error is zero (Low biased)

poor testing on unseen data (low variance)

Excellent performance on training set but poor performance with unseen data.  
(High variance)

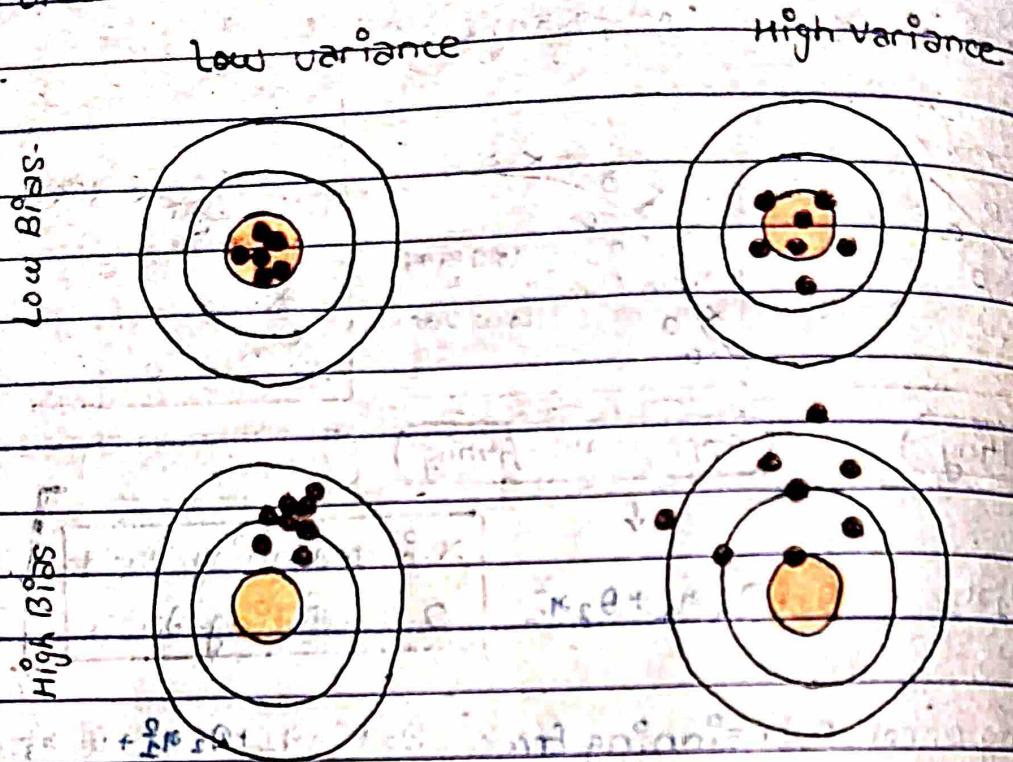
Training error is zero (Low biased)

→ If the algorithm is too simple (hypothesis with linear equation) then it may be on high bias and low variance condition and thus is error prone.

→ If algorithm fits too complex (hypothesis with high degree equation) then it may be on high variance and low bias.

→ In the latter condition the new entries will not perform well. Well, there's something b/w both of these conditions known as Bias-Variance Tradeoff.

There are four possible combination of bias and variance —



① Low bias, Low variance : The combination of low bias & low variance shows an ideal model. However, it is not possible practically.

② Low Bias, High Variance : model predictions are inconsistent and accurate on average. This occur when a model learns with a large nos of parameters and hence lead to an overfitting.

③ High Bias, Low Variance : predictions are consistent but inaccurate on average. This case occur when model learns with a large nos of parameters and hence lead to underfitting. (few)

④ High bias, High Variance : With high variance & high bias prediction are inconsistent and inaccurate on average.

# ★ Regularization in Machine Learning :

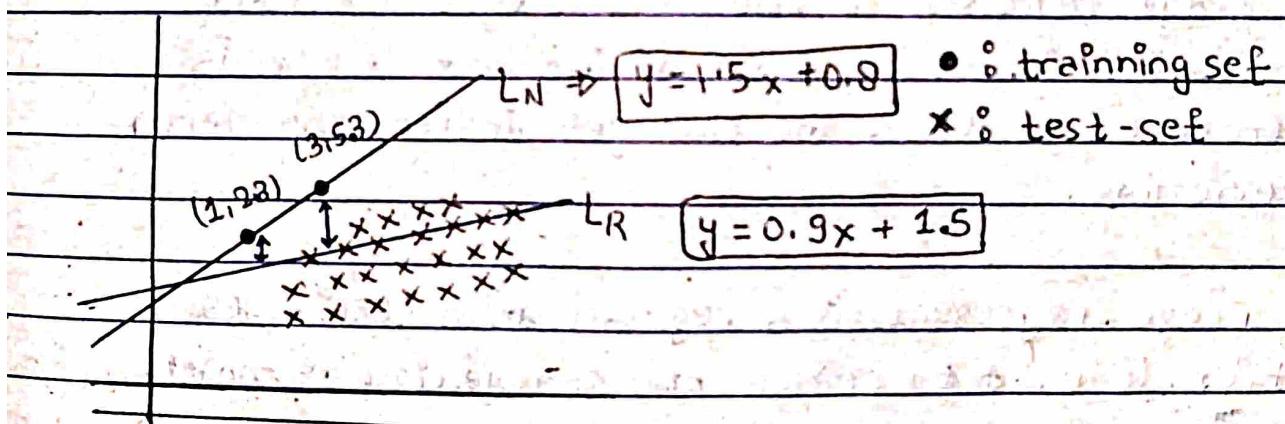
While developing ML models you must have encountered a situation in which the training accuracy of model is low but the testing accuracy is too low. This is the case of Overfitting.

Regularization helps us to solve the problem of overfitting, which reduces error by fitting the function appropriately on the given training set.

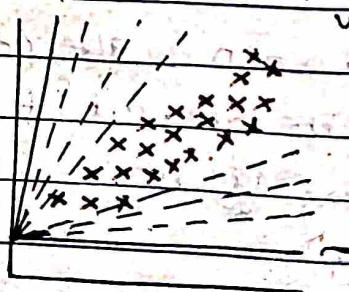
Regularization techniques →

- (1) Lasso regularization - L1 regularization
- (2) Ridge regularization - L2 regularization
- (3) Elastic-Net regularization - L1 & L2 regularization

→ Ridge Regularization



Note: → Overfitting ( $m \rightarrow \infty$ )  $y = m x + b$   
very much depend on 'm'. on 'm' : scale factor to each input value

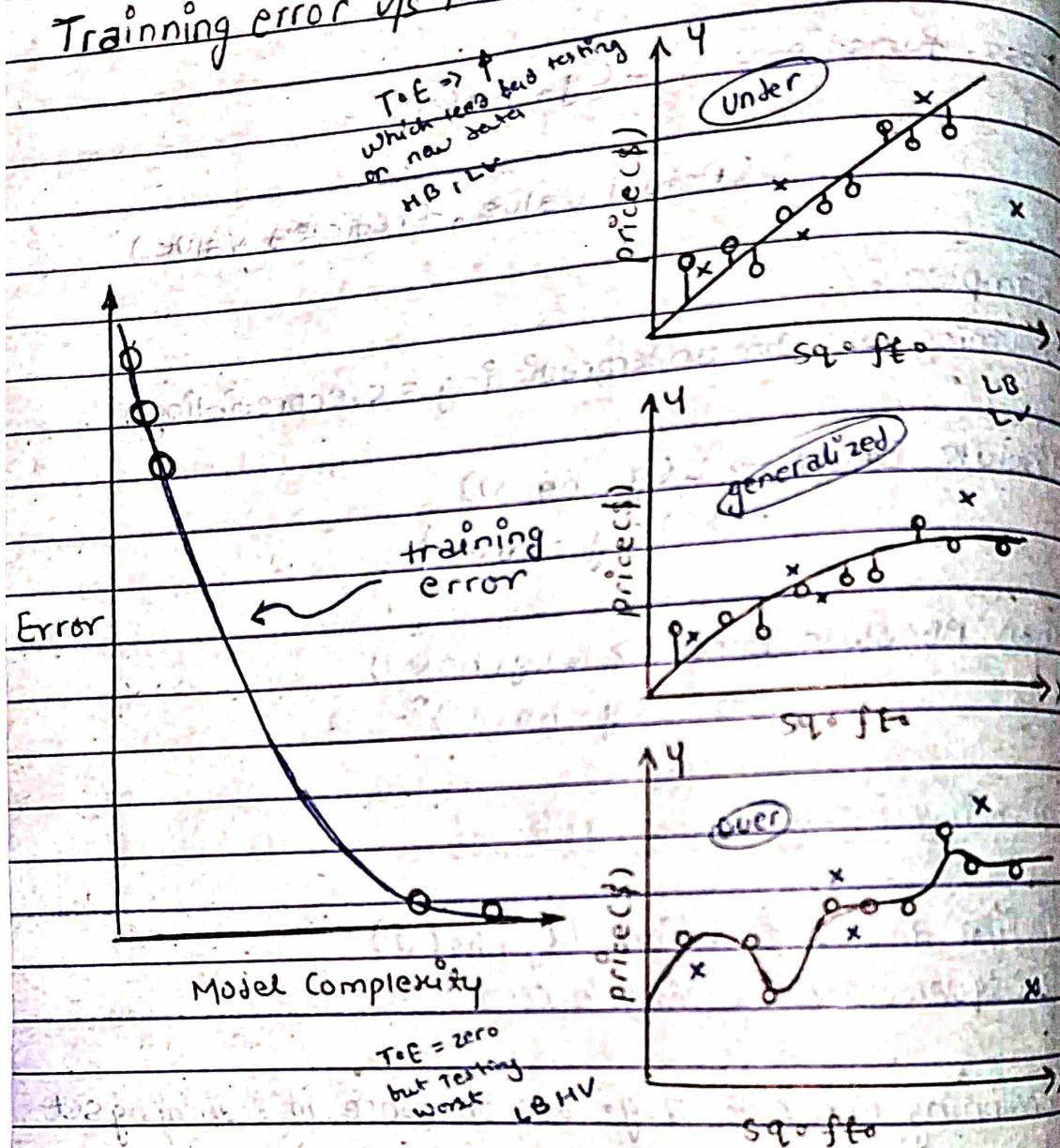


→ Underfitting ( $m \rightarrow 0$ )

target value is not so depending on 'm'

means, Overfitting करते के लिए slope 'm' को बढ़ावा दें।

## Training error v/s model complexity

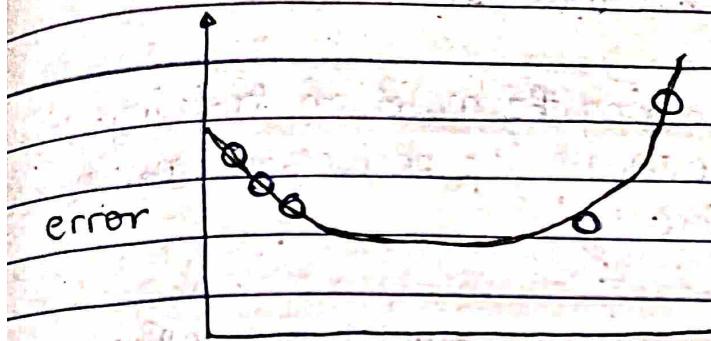


Is training Error a good measure of predictive performance?

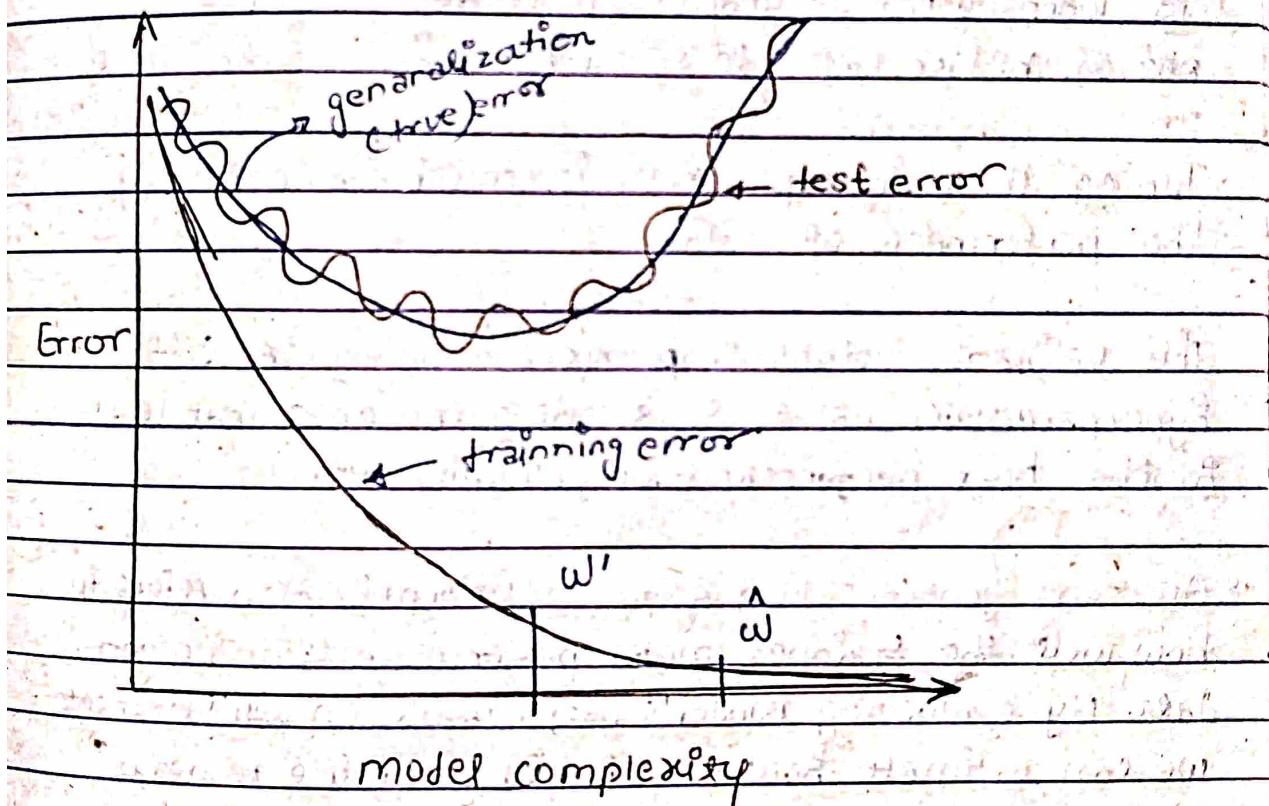
- Training Error is overly optimistic because  $\theta$  was fit to training data
- small training error  $\nrightarrow$  good predictions
- unless Training data include everything you might ever see

Date \_\_\_\_\_  
Page \_\_\_\_\_

## Generalization error Training, true and test error vs model complexity



model complexity



model complexity

Overfitting if  $\hat{w}$  if there exist a model with estimated parameters  $\hat{w}$

① training error ( $\hat{w}$ ) < training error ( $w'$ )

② true error ( $\hat{w}$ ) > true error ( $w'$ )

## regularization

$$h_{\theta}(x) = \theta_0 + \theta_1 x_1 + \theta_2 x_2^2 + \theta_3 x_3^3 + \dots + \theta_n x_n^n$$

$\theta_2, \theta_3 \rightarrow$  unimportant

- ① remove manually
- ② use model selection algo
- ③ make weight  $\theta_j^*$  small ( $\lambda$ )  $\downarrow$  regularization

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

repeat till converge

$$\theta_j^* = \theta_j^* - \alpha \frac{\partial J(\theta_0, \theta_1, \theta_2, \dots, \theta_n)}{\partial \theta_j}$$

update simultaneously  
( $j = 0, 1, 2, \dots, n$ )

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

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

$$\theta_j^* = \theta_j^* - \alpha \left[ \frac{1}{m} \sum_{i=1}^m [h_{\theta}(x^{(i)}) - y^{(i)}] \theta_j^{(i)} - \frac{\lambda}{m} \sum_{j=1}^n \theta_j^* \right]$$

$$\theta_j^* \left( 1 - \alpha \frac{\lambda}{m} \right) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \theta_j^{(i)}$$