	Jigar Siddhpura
()	6000 421015 5
1	Experiment 2
	e xposiment 2
	Micro of To and to a fail of
	Airo: To perform theor & find error associated in the model.
	(acoustile facilitates) for tool
	Theory:
	1 Springly 2 x 1 mu c h I what som -
	Linear regression is a supervised ML algorithm . It is
	used for predictions for continuous freel or mumeric
	Vosibles (like salary, sales, etc.). It shows
bla mil	linear soldienship between a dependent variable (y) & or more independent var (x). It shows how
	dependent variable changes acc to independent var.
	dependent variable changes acc to independent var.
	Methemodically IR is frepr as: y = bo + b, n + E y - dependent vor X - independent vor bo - intercept of line b, > 12- co-efficient.
	y - dependent vor X - independent vor
	bo - intercept of line 1) - 12 co-efficient.
1 00 1/1	E - com (forth as a good to the) see look
C.	Different values of x' & y variables, weights or co-efficients
	of line (b., b) gives diff. line of regression & ost
	of time (b. b.) gives diff. line of regression & cost function is used to estimate best fit ine. It
	is used to find acc of mapping function, which
	maps input var to output var. Here we used
2.470	N A 173C WHENE
	MSE = 1 2 4: - (b: x; -b.) 2.
	en at the Not in the second of hard -
	where, N = total absention
	y: = actual value
	yo = bixi + bo = predicted val.
WGVVGD	AMOSTADUGE 903 FOR EDUCATIONAL USE
NGVKAR	

0.816	Tiper Siddle
	S. Marchand
	C largery !
- 1	D 10 1
A1 6	Prod Procedure 3
	Park and I are 1
	Past 10 (without stilledon)
	- import of all libraries
7 -17	- Here dataset of 2 var (x & Y) (array) - Reshape the lists:
airemun	- Using townslar for a calculating analysis &
	- Using tormulas for g calculating gradients &
61	Part 1.2 (with sklears)
Lond	- " wed knook aperson library from skeam linear mode
. 300	- fit the data using fit () function.
	- fit the data using fit() function. - accuracy - used Score() Anction
-	Part 2-1 (without skleger, new dataset)
· tui	De wed the came steps as in 11
	Part 2:2 (with Stilenon, new dataset) Here, using new dataset 'Student Performance! Perform one hot encoding
1	- Poll - hat a dataset simeni reparraying
La	- the library of in dra
1	The library as in 102
	Part 3-1 (without stillean, dataset 2)
4.	- used the same step as in +1
	Part 3.2 (with library, datact 2)
	- Used the salvey dates that only has I columns
	- Here X = TensExperience & Y = Salony.
	- Used the library from have as in 1-2
	Ando las ar sade
	silve letter . A
	. In hallburg : get skill a ge
Newcol	FOR EDUCATIONAL USE
NGVKGR	TOX EDUCATION AE USE

	Observations: for part 1.1, $\omega = [[1.8853]]$ $b = -0.2708$ $MSE = 4.1559$
	for part 1.2, $w = [[1.9763]]$ score = 0.1873 b = -0.4873
	for part 2-1 (with formula), w = [[27762]] b = 0.3872 MSE = 3.3126
	for past 2.2 = 5:0re = 0.1967 for past 3.1, w = [[(2318.527]] b = 6481.917
	MSE = 23/15 + 88 1. 111
	Conclusions & Hence, we beared how to apply the model on a linear dependent dataset wing est statistical & method of wing library, Also, learned how to cale. Cost function i.e. MSE & R
NAVKAR	FOR EDUCATIONAL USE

Describe the procedure that is used to perform Linear regression using Least Square Method carry out the
experiment step-by-step for simple linear regression for following dataset with and without using scikit library.
Describe every line of code with the proper interpretation of the output.

Х	2	3	4	5	6	7	8	9	10	
Υ	1	3	6	9	11	13	15	17	20	

```
"""### 1. without sklearn"""
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
X = [i for i in range(2,11)]
Y = [1,3,6,9,11,13,15,17,20]
df = pd.DataFrame({'X':X,'Y':Y})
X = df['X'].values.reshape(-1,1)
Y = df['Y'].values.reshape(-1,1)
learning_rate=0.01
no_of_itr=100
y_pred_arr = []
m,n = X.shape
w = np.zeros((n,1))
b = 0
def predict(X):
   return X.dot(w) + b
def update_wt():
   global w
   global b
   y_pred = predict(X)
   y_pred_arr.append(y_pred)
   dw = -(X.T).dot(Y - y_pred)/m
   db = -np.sum(Y - y_pred)/m
   w = w - learning_rate * dw
   b = b - learning_rate * db
for i in range(no_of_itr):
 update_wt()
print(w)
print(b)
MSE = np.square(np.subtract(Y,y_pred_arr)).mean()
print(MSE)
```

OUTPUT:

```
[[1.88536856]]
-0.27035591489738453
4.155960303200094
```

WITH SKLEARN:

```
"""### 1.2 with *sklearn*""
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

X = df['X'].values.reshape(-1,1)
Y = df['Y'].values.reshape(-1,1)

model = LinearRegression()
model.fit(X,Y)
print(model.score)
```

OUTPUT:



2. DATASET: STUDENT PERFORMANCE

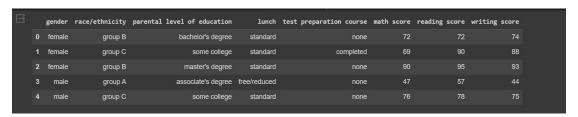
WITH SKLEARN:

```
"""### 2.1 with dataset with lib"""
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import LabelEncoder

df = pd.read_csv('/content/gdrive/MyDrive/DMW/datasets/StudentsPerformance.csv')
df.head()
df['final_score'] = df.apply(lambda x: (x['math score'] + x['reading score'] +
x['writing score'])/3, axis=1)

df2 = df
df2 = df2.drop(['math score', 'reading score', 'writing score'],axis=1)

df2.head()
```



```
df2 = pd.get_dummies(df2, columns=['gender','lunch','parental level of
education','race/ethnicity','test preparation course'])
df2.head()
```

```
| Parental level of Final_score | gender_female | gender_male | lunch_free/reduced | lunch_standard | education_associate's | education_baschelor's |
```

```
# multi-variate
y = df2['final_score']
x = df2.drop(['final_score'],axis=1)

xtrain, xtest, ytrain, ytest = train_test_split(x,y,test_size=0.25,random_state=10)

model = LinearRegression()
model.fit(xtrain,ytrain)
score = model.score(xtest,ytest)
print(score)
ypred = model.predict(xtest)
```

OUTPUT:

0.19674412629893356

WITHOUT SKLEARN:

```
"""### 2.2 with formula"""
X = df['X'].values.reshape(-1,1)
Y = df['Y'].values.reshape(-1,1)
learning_rate=0.01
no_of_itr=100
y_pred_arr = []
m,n = X.shape
w = np.zeros((n,1))
b = 0
def predict(X):
   return X.dot(w) + b
def update_wt():
   global w
   global b
   y_pred = predict(X)
   y_pred_arr.append(y_pred)
   dw = -(X.T).dot(Y - y_pred)/m
   db = -np.sum(Y - y_pred)/m
   w = w - learning_rate * dw
   b = b - learning_rate * db
for i in range(no_of_itr):
 update_wt()
```

```
print(w)
print(b)
MSE = np.square(np.subtract(Y,y_pred_arr)).mean()
print(MSE)
```

OUTPUT:

```
[[2.776264564763]]
0.3872355565633
3.312643565737
```

3. DATASET: SALARY

WITHOUT SKLEARN:

```
"""### 3.1 without lib"""

df = pd.read_csv('/content/gdrive/MyDrive/BDI/salary_data.csv')
df.head()
```

∃		YearsExperie	nce	Salary
	0		1.1	39343.0
	1		1.3	46205.0
	2		1.5	37731.0
	3		2.0	43525.0
	4		2.2	39891.0

```
X = df['YearsExperience'].values.reshape(-1,1)
Y = df['Salary'].values.reshape(-1,1)
learning_rate=0.01
no_of_itr=100
y_pred_arr = []

m,n = X.shape
w = np.zeros((n,1))
b = 0
def predict(X):
    return X.dot(w) + b

def update_wt():
    global w
    global b
    y_pred_arr.append(y_pred)

# calculating gradients
```

```
dw = -(X.T).dot(Y - y_pred)/m

db = -np.sum(Y - y_pred)/m

# updating weights

w = w - learning_rate * dw

b = b - learning_rate * db

for i in range(no_of_itr):
    update_wt()

print(w)

print(b)

MSE = np.square(np.subtract(Y,y_pred_arr)).mean()
print(MSE)
```

OUTPUT:

```
[[12315.52743012]]
6481.917498721374
237548657.81185108
```

WITH SKLEARN:

```
"""# 3.2"""

model = LinearRegression()
model.fit(X,Y)
ypred = model.predict([[4]])
print(ypred)
```

OUTPUT:

```
[[63592.04948449]]
```