

22

22

March

Week 13 Day 188-268

Tuesday

MARCH 2022

SUN	MON	TUE	WED	THU	FRI	SAT
1	2	3	4	5	6	7
8	9	10	11	12	13	14
15	16	17	18	19	20	21
22	23	24	25	26	27	28
29	30	31				

## 09:00 LINEAR REGRESSION

## 10:00 Simple LINEAR REGRESSION -

$$y = \beta_0 + \beta_1 x$$

↑      ↑  
target Y    parameters  
              constant

13:00 To find the parameters we need to find ERROR function.

$$\ell_i = \hat{y}_i - y_i$$

↓      ↓  
real    predicted  
target   outcome

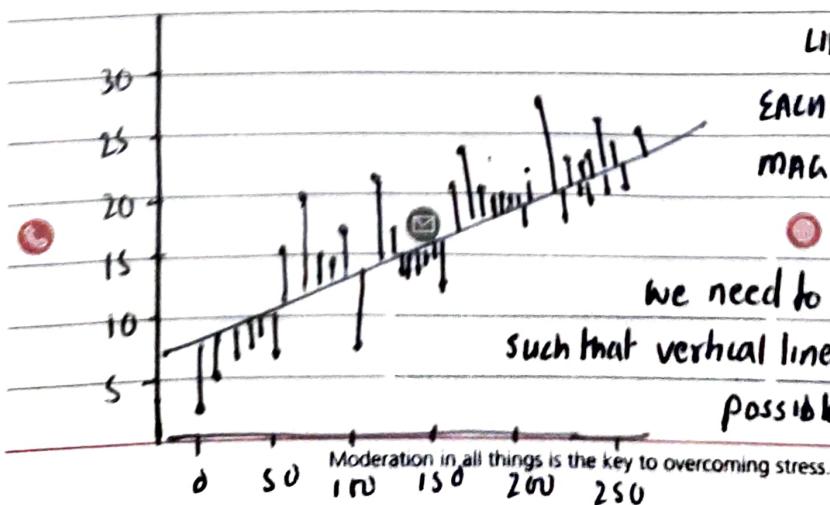
17:00 FOR LINEAR REGRESSION WE MINIMISES SUM OF SQUARED ERRORS.

## 18:00 ERROR VISUALIZATION -

- DOTS REPRESENT DATA
- - LINE IS FITTED STRAIGHT LINE

EACH VERTICAL LINE IS MAGNITUDE OF ERROR

we need to fit the straight line such that vertical lines are small as possible.



Moderation in all things is the key to overcoming stress.

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7
8	9	10				
11	12	13	14	15	16	17

11 12 13 14 15 16 17 18 19 20 21 22 23 24  
25 26 27 28 29 30

March

22

Wednesday

Week 13 Day (082-283)

23

09.00

WHY WE SQUARE THE ERRORS?

10.00

- ERRORS CAN BE POSITIVE OR NEGATIVE (ABOVE OR BELOW THE LINE)

11.00

↳ so the results could be positive or negative.

12.00

↳ if we did not square the errors, we would be adding bunch of negative errors and reduce the sum of errors.

13.00

↳ it would trick us in thinking we are fitting a good straight line, but in reality we are not.

14.00

- Negative values could decrease the error

15.00

- Penalize the errors.

16.00

Finding Coefficient -

17.00

$$\beta_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$\bar{x}$  - mean of independent variable

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

$\bar{y}$  - mean of the target.



To do list



24

'22

March

Week 13 Day (083-202)

Thursday

MARCH 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MAR	APR	MAY	JUN	JUL	AUG	SEP	OCT
1	2	3	4	5	6	7	8	9	10	11	12	13		
14	15	16	17	18	19	20	21	22	23	24	25	26	27	
28	29	30	31											

09:00

## FINDING THE COEFFICIENT IN PRACTICE - + PYTHON

10:00

`reg = LinearRegression ()` // initialize the model

11:00

`reg. fit (X, y)` // fit on X and Y, then we can retrieve both intercept and coeff.

12:00

`intercept = reg. intercept_ [0]``coeff = reg. coef_ [0] [0]`

13:00

Estimate the relevance of coefficients by

14:00

• p-value : quantify statistical significance

15:00

• Tells us if we can reject the null hypothesis or not

16:00

In Python, we can analyse pvalue for each coefficient

like this —

`exog = sm. add_constant (X)``est = sm. OLS (y, exog). fit ()``print (est. summary ())`

// we statistical model package that allows us to use summary of the model.

sm model provides

- statistical test

- data exploration

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7
8	9	10				
11	12	13	14	15	16	17

18 19 20 21 22 23 24  
25 26 27 28 29 30

March

'22

25

Friday

Week 13 Day (084-281)

## 09.00 NULL HYPOTHESIS -

10.00	coef.	std err	t	P> t	[0.025	0.915]
	7.0326	0.458	15.360	0.00	6.130	7.435
11.00	0.0475	0.003	11.668	0.00	0.042	0.053

↑  
 actually it is not zero  
 but it is so small, that it  
 appears to be zero.

## 14.00 ONCE WE KNOW ARE PARAMETERS ARE RELAVANT

WE MUST ACCESS THE MODEL ITSELF

## 16.00 Accuracy of Linear Model -

17.00 Residual Standard Error (RSE) =  $\sqrt{\frac{1}{n-2} RSS} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$

- smaller the value - the better it is

since

\*\*\* difference b/w actual and predicted is small



26

22

March

Week 13 Day (085-280)

Saturday

MARCH 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	13	
14	15	16	17	18	19	20	21	22	23	24	25	26	27
28	29	30	31										

09.00

$R^2$  Value - It measures the proportion of variability explained by a feature  $x$ .

10.00

As it approaches 1 - it means that we are explaining a lot of variability in our target.

11.00

$$R^2 = \frac{TSS - RSS}{TSS} = 1 - \frac{RSS}{TSS}, TSS = \sum (y_i - \hat{y}_i)^2$$

12.00

- Measure of proportion of variability explained by feature  $X$

14.00

In Python, we can use the same process as we did before to find the p-value using the same package.

15.00

`exog = sm.add_constant(X)`

`est = sm.OLS(y, exog).fit()`

16.00

`print(est.summary())`

Sunday 27



APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
11	12	13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30								

March

22

Monday

Week 14 Day (087-278)

09.00 ACCURACY OF LINEAR MODEL

10.00 OLS REGRESSION RESULTS

11.00	Dep. Variable	sales	R-squared	0.612
	Model	OLS	Adj. R-squared	0.610
12.00	Method	Least Square	F-statistic	312.1
	Date	Sat, 09 May 20	Prob (F-statistic)	1.47e-42
13.00	Time	15:19:51	Log-Likelihood	-514.05
	No. Obs.	200	AIC	1042
14.00	Df Residuals	198	BIC	1049
	Df Model	1		
15.00	Covariance Type	nonrobust		

16.00	coef	std err	t	P> t	[ 0.025	0.975 ]
	const	7.0326	0.458	15.360	0.00	6.130 7.935
17.00	TV	0.0475	0.003	11.668	0.00	0.042 0.053

18.00	Omnibus :	0.531	Durbin-Watson	1.935
	Prob (Omnibus)	0.767	Jarque-Bera (JB)	0.669
	Skew :	-0.089	Prob (JB)	0.716
	Kurtosis :	2.779	Cond. No	338.



29

22

March

Week 14 Day 1088-2771

Tuesday

MARCH 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	13	
14	15	16	17	18	19	20	21	22	23	24	25	26	27
28	29	30	31										

09:00

## MULTIPLE LINEAR REGRESSION

10:00

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p$$

req = LinearRegression()

req.fit(X, y)

beta\_0 = req.intercept\_[0]

beta\_1 = req.intercept\_[0][0]

beta\_2 = req.intercept\_[0][1]

beta\_3 = req.intercept\_[0][2]

TO ACCESS THE MODEL OF MULTIPLE LINEAR REGRESSION

↳ we use 'F' statistics.

$$F = \frac{\frac{TSS - RSS}{P}}{\frac{RSS}{(n-p-1)}}$$

P - no. of predictors

n - no. of data points



APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10	11	12	13	14
15	16	17	18	19	20	21	22	23	24	25	26	27	28

March 30  
Wednesday

'22  
Week 14 Day (089-276)

09:00

In Python, we would use same package for  $F^2$  statistics

10:00

★ Whole Result Report as previous.

11:00

- How RELEVANT IS THE MODEL

12:00

→  $F \gg 1$ , there is a strong relationship

13:00

→ for small dataset, F has to be way larger than 1

14:00

Hands On -

15:00

- On Google Colab / Could also be done on Jupyter Notebook

16:00

```
import pandas as pd
```

```
import numpy as np
```

```
import matplotlib.pyplot as plt
```

```
from sklearn.linear_model import
```

```
LinearRegression
```

```
import statsmodels.api
```

.matplotlib inline

31

22

March

Thursday

09:00

## READING THE CSV FILE

10:00

11:00

12:00

```
data = pd.read_csv('path or url')
```

13:00

## SIMPLE LINEAR REGRESSION

14:00

```
plt.figure(figsize=(16, 8))
```

15:00

```
plt.scatter(data['TV'], data['sales'],
            c='black')
```

16:00

```
plt.xlabel('Money Spent on TV Ads ($)')
```

17:00

```
plt.ylabel('Sales (k$)')
```

18:00

```
plt.show()
```



MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

April 02

Saturday

Week 14 Day (092-273)

09.00

## SPECIFYING OUR FEATURE

10.00

```
X = data['TV'].values.reshape(-1,1)  
y = data['sales'].values.reshape(-1,1)
```

11.00

```
req = LinearRegression()  
req.fit(X,y)
```

12.00

```
print(f"The Linear Model is \n")
```

13.00

```
y = (req.intercept_[0]) +  
(req.coef_[0][0]*TV")
```

14.00

15.00

16.00

## OUTPUT-

17.00

The Linear Model is:

$$Y = (\text{req.intercept}_[0]) + (\text{req.coef}_[0][0] * TV)$$

Sunday 03

→ NEXT STEP IS PLOTTING THE LINE ←



04

'22

April

Week 15 Day (094-271)

Monday

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10				
11	12	13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30								

09.00

## # PLOTTING THE LINE

10.00

predictions = reg.predict(X)

11.

plt.figure(figsize=(16, 8))

12.

plt.scatter(X, y, c="black")

13.

plt.plot(X, predictions, c="blue",  
linewidth=2)

14.

plt.xlabel('Money Spent on TV Ads')

15.

plt.ylabel('Sales (K\$)')

16.

plt.show()



calling predict method

on X

black is color of dots

Y in this case is

predictions

blue is line color

17.00

## Access the Quality of Our Model

18.00

# CATCH - FIRST IN ORDER TO USE ANY API, IT SHOULD  
BE PROPERLY EXPORTED.

In our case, import statsmodels.api as sm - ✓

import statsmodels.api - X



MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

April 05  
Tuesday

Week 15 Day (095-270)

09.00

## ACCESS THE QUALITY OF OUR MODEL

10.00

```
import statsmodels.api as sm
```

11.00

```
X = data['TV']
```

12.00

```
y = data['Sales']
```

13.00

```
exog = sm.add_constant(X)
```

14.00

```
est = sm.OLS(y, exog).fit()
```

15.00

```
print(est.summary())
```

### OLS Regression Results

16.00

`X = data['TV']` - selects the TV column of the dataset and assigns it to the variable 'X'

17.00

`y = data['Sales']` - select the sales column of the dataset and assigns it to variable 'y'

`exog = sm.add_constant(X)` - adds column of ones to the X variable to serve as intercept term and assign the results to exog

`'est = sm.OLS(y, exog).fit()'` - Fits an ordinary least square(OLS) regression model using 'y' as the response variable and 'exog' as the predictor variables, then assign result to 'est'.

06

22

April

Week 15 Day (096-269)

Wednesday

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
				1	2	3	4	5	6	7	8	9	10
11	12	13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30								

09.00

print (est.summary ()) - print a summary of the regression results, which include information such as

10.00

coefficient, standard errors, t-values, p-values, of predictor variables, as well as goodness of fit statistics like R-squared and F-statistic

11.00

12.00

13.00

Overall, this code is fitting a simple linear regression model to the data and printing out summary of results.

14.00

## Mul Linear Regression

15.00

```
Xs = data.drop(['Sales'], axis=1)
y = data['Sales'].values.reshape(-1,1)
```

16.00

```
req = LinearRegression()
req.fit(Xs, y)
```

```
print(f"The Linea Model is: \n Y =
{req.intercept_[0]} + {req.coef_[0][0]} *
TV + {req.coef_[0][1]} * radio +
{req.coef_[0][2]} * newspaper")
```



MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
					1	2	3	4	5	6	7	8	
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

April 07  
Thursday

22  
Week 15 Day (097-268)

09:00 LINEAR REGRESSION CLASS FROM SCIKIT LEARN

10:00 `Xs = data.drop(['sales'], axis=1)` - select all columns of the dataset except for 'sales' and assigns them to the variable 'Xs'

11:00 `y = data['sales'].values.reshape(-1,1)` - select the sales column of dataset and convert it to a two-dimensional Numpy array with one column using the 'reshape()' method. This is assigned to the variable 'y'.

12:00 `reg = LinearRegression()` - create an instance of the LR class and assign it to variable 'reg'.

13:00 `reg.fit(Xs,y)` - fits the linear regression model using 'Xs' as the predictor variable and 'y' as the responsible variable

14:00 `print(f'The Linear — newspaper")` - prints the equation of linear model, including intercept term and coefficient of the predictor variable

15:00

Overall, code is performing a multiple LR using TV, radio and newspaper column of dataset as predictor variable and sales column as the response variable.

It then prints out the equation of Linear model that was fitted



08

22

April

Week 15 Day (098-267)

Friday

APRIL 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
1	2	3	4	5	6	7	8	9	10				
11	12	13	14	15	16	17	18	19	20	21	22	23	24
25	26	27	28	29	30								

09.00

This code fits a multilinear regression model to data, with 'TV', 'radio', and 'newspaper' as predictor variables and 'sales' as the response variable.

10.00

```
X = np.column_stack((data['TV'],
12.00          data['newspaper']))
```

```
y = data['sales'].values.reshape(-1,1)
```

```
exog = sm.add_constant(X)
```

```
est = sm.OLS(y, exog).fit()
```

```
print(est.summary())
```

1

2

3

4

5

16.00

17.00

18.00

1 - creates matrix X with 3 columns, corresponding to the TV, radio, and newspaper variables.

The line stacks the columns horizontally using np.column\_stack() function from NumPy.

2 - creates a column vector y containing the sales variables, reshaped to have one column and many rows as original data.

MAY 2022

MON	TUE	WED	THU	FRI	SAT	SUN	MON	TUE	WED	THU	FRI	SAT	SUN
9	10	11	12	13	14	15	16	17	18	19	20	21	22
23	24	25	26	27	28	29	30	31					

April

22

09

Week 15 Day (099-266)

Saturday

09.00 The line uses the `e.values`' attribute of pandas DataFrame to get the the underlying data as Numpy array and  
10.00 then applies the `'reshape()'` method to convert it to the desired shape.

11.00  
12.00 3) `exog = sm.add_constant(X)` adds a constant column to the matrix X using the `'add_constant()'` function from statsmodels. The constant term needed to fit  
13.00 the intercept term in regression model.

14.00 4) `est = sm.OLS(y, exog).fit()` - fits a multilinear regression model to the data, using the ordinary least square (OLS) method from statsmodel.

15.00 The line creates new object 'est' of the OLS class and applies its `'fit()'` method to the data.

16.00 This method estimates the regression coefficient by minimizing the sum of squared residuals b/w the Sunday 10 predicted values and actual values and returns an object of the 'RegressionResults' class that contains various statistics and diagnostic about the model.

5) print summary of regression results to the console.

