

# SALARY DATA STATISTICS

```
In [4]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

```
In [6]: import warnings
warnings.filterwarnings('ignore')
```

```
In [22]: salary = pd.read_csv(r"E:\Data Science & AI\Dataset files\Salary_Data.csv")
```

```
In [24]: salary.head()
```

```
Out[24]:
```

	YearsExperience	Salary
0	1.1	39343
1	1.3	46205
2	1.5	37731
3	2.0	43525
4	2.2	39891

```
In [26]: salary.shape
```

```
Out[26]: (30, 2)
```

```
In [28]: salary.columns
```

```
Out[28]: Index(['YearsExperience', 'Salary'], dtype='object')
```

```
In [30]: salary.isnull().sum()
```

```
Out[30]: YearsExperience    0
Salary                    0
dtype: int64
```

```
In [32]: salary.dtypes
```

```
Out[32]: YearsExperience    float64
Salary                    int64
dtype: object
```

```
In [38]: x = salary.drop('YearsExperience',axis=1)
y = salary['Salary']
```

```
In [48]: from sklearn.model_selection import train_test_split
x_train, x_test, y_train, y_test=train_test_split(x,y,random_state=1234,test_siz
```

```
In [50]: x_train.shape,x_test.shape
```

Out[50]: ((21, 1), (9, 1))

```
In [52]: y_train.shape,y_test.shape
```

Out[52]: ((21,), (9,))

```
In [56]: salary.shape
```

Out[56]: (30, 2)

```
In [58]: x_train
```

Out[58]:

	Salary
--	--------

13	57081
----	-------

22	101302
----	--------

24	109431
----	--------

0	39343
---	-------

2	37731
---	-------

27	112635
----	--------

26	116969
----	--------

18	81363
----	-------

5	56642
---	-------

16	66029
----	-------

25	105582
----	--------

11	55794
----	-------

9	57189
---	-------

17	83088
----	-------

29	121872
----	--------

20	91738
----	-------

12	56957
----	-------

21	98273
----	-------

6	60150
---	-------

19	93940
----	-------

15	67938
----	-------

```
In [60]: y_train
```

```
Out[60]: 13      57081
          22     101302
          24     109431
           0      39343
           2      37731
          27     112635
          26     116969
          18      81363
           5      56642
          16     66029
          25     105582
          11      55794
           9      57189
          17     83088
          29     121872
          20      91738
          12      56957
          21      98273
           6      60150
          19      93940
          15      67938
          Name: Salary, dtype: int64
```

```
In [62]: x_test
```

```
Out[62]:
```

	Salary
7	54445
10	63218
4	39891
1	46205
28	122391
8	64445
3	43525
23	113812
14	61111

```
In [64]: y_test
```

```
Out[64]: 7      54445
          10     63218
           4     39891
           1     46205
          28     122391
           8     64445
           3     43525
          23     113812
          14     61111
          Name: Salary, dtype: int64
```

```
In [66]: x_train.ndim
```

Out[66]: 2

```
In [68]: from sklearn.linear_model import LinearRegression
LR = LinearRegression()
LR.fit(x_train,y_train)
```

Out[68]:

▼ LinearRegression ⓘ ?

LinearRegression()

```
In [70]: y_predications = LR.predict(x_test)
```

```
In [72]: y_predications
```

Out[72]: array([ 54445., 63218., 39891., 46205., 122391., 64445., 43525.,  
113812., 61111.])

```
In [74]: y_test.shape,y_predications.shape
```

Out[74]: ((9,), (9,))

```
In [76]: x_test
```

Out[76]:

	Salary
7	54445
10	63218
4	39891
1	46205
28	122391
8	64445
3	43525
23	113812
14	61111

```
In [80]: x_test.iloc[0]
x_test.iloc[0].values
```

Out[80]: array([54445], dtype=int64)

```
In [86]: LR.predict([x_test.iloc[0].values,x_test.iloc[1].values])
```

Out[86]: array([54445., 63218.])

```
In [88]: ip1 = [5]
LR.predict([ip1])
```

Out[88]: array([5.])

```
In [92]: x_test.shape,y_test.shape,y_predications.shape
```

```
Out[92]: ((9, 1), (9,), (9,))
```

```
In [96]: test_data=x_test
test_data['y_actual']=y_test
test_data['y_predictions']=y_predications
test_data
```

```
Out[96]:
```

	Salary	y_actual	y_predictions
7	54445	54445	54445.0
10	63218	63218	63218.0
4	39891	39891	39891.0
1	46205	46205	46205.0
28	122391	122391	122391.0
8	64445	64445	64445.0
3	43525	43525	43525.0
23	113812	113812	113812.0
14	61111	61111	61111.0

```
In [102...] print(y_test.values[:5])
```

```
[ 54445  63218  39891  46205 122391]
```

```
In [106...] print(y_predications[:5])
```

```
[ 54445.  63218.  39891.  46205. 122391.]
```

```
In [108...] from sklearn.metrics import r2_score,mean_squared_error
```

```
In [114...] R2=r2_score(y_test,y_predications)
MSE=mean_squared_error(y_test,y_predications)
#MSE**(1/2)
RMSE=np.sqrt(MSE)
#accuracy_score(y_test,y_predictions) # it is a regression tech
print("R-sqaure:",R2)
print("MSE:",MSE)
print("RMSE:",RMSE)
```

```
R-sqaure: 1.0
```

```
MSE: 4.117521271375071e-23
```

```
RMSE: 6.41679146565873e-12
```

```
In [120...] s=0
for i in range(len(y_test)):
    v1=y_test.values[i]-y_predications[i]
    v2=v1**2
    s=s+v2
print(s/len(y_test))
```

```
4.117521271375071e-23
```

```
In [122... LR.coef_  
print("The coeffiecnt of Years_of_experience is:",LR.coef_)
```

The coeffiecnt of Years\_of\_experience is: [1.]

```
In [124... LR.intercept_
```

```
Out[124... 1.4551915228366852e-11
```

```
In [126... x_train.columns
```

```
Out[126... Index(['Salary'], dtype='object')
```

```
In [138... from sklearn.feature_selection import VarianceThreshold  
vt=VarianceThreshold(threshold=0)  
# Threshold variance value  
# we want to drop the feaure based on threshold  
vt.fit(salary)
```

```
Out[138... ▼ VarianceThreshold ⓘ ?  
VarianceThreshold(threshold=0)
```

```
In [140... dir(vt)
```

```
Out[140... ['__abstractmethods__',
             '__annotations__',
             '__class__',
             '__delattr__',
             '__dict__',
             '__dir__',
             '__doc__',
             '__eq__',
             '__format__',
             '__ge__',
             '__getattr__',
             '__getstate__',
             '__gt__',
             '__hash__',
             '__init__',
             '__init_subclass__',
             '__le__',
             '__lt__',
             '__module__',
             '__ne__',
             '__new__',
             '__reduce__',
             '__reduce_ex__',
             '__repr__',
             '__setattr__',
             '__setstate__',
             '__sizeof__',
             '__sklearn_clone__',
             '__str__',
             '__subclasshook__',
             '__weakref__',
             '_abc_impl',
             '_build_request_for_signature',
             '_check_feature_names',
             '_check_n_features',
             '_doc_link_module',
             '_doc_link_template',
             '_doc_link_url_param_generator',
             '_get_default_requests',
             '_get_doc_link',
             '_get_metadata_request',
             '_get_param_names',
             '_get_support_mask',
             '_get_tags',
             '_more_tags',
             '_parameter_constraints',
             '_repr_html_',
             '_repr_html_inner',
             '_repr_mimebundle_',
             '_sklearn_auto_wrap_output_keys',
             '_transform',
             '_validate_data',
             '_validate_params',
             'feature_names_in_',
             'fit',
             'fit_transform',
             'get_feature_names_out',
             'get_metadata_routing',
             'get_params',
             'get_support',
```

```
'inverse_transform',  
'n_features_in_',  
'set_output',  
'set_params',  
'threshold',  
'transform',  
'variances_']
```

```
In [142... vt.variances_  
# 300 is first column variance (T)  
# 1.25 is second column variance (T)  
# 30 is column variance (T)  
# 0 is fourth column variance (F)
```

```
Out[142... array([7.78515556e+00, 8.46600000e+04])
```

```
In [144... vt.get_support()
```

```
Out[144... array([ True,  True])
```

```
In [146... vt.get_support()
```

```
Out[146... array([ True,  True])
```

```
In [148... vt.get_params()
```

```
Out[148... {'threshold': 0}
```

```
In [150... vt.threshold
```

```
Out[150... 0
```

```
In [154... cols=vt.get_feature_names_out()  
# the above syntax gives the column names  
# These feature only we want include  
salary[cols]
```



Out[154...

	YearsExperience	Salary
0	1.1	39343
1	1.3	46205
2	1.5	37731
3	2.0	43525
4	2.2	39891
5	2.9	56642
6	3.0	60150
7	3.2	54445
8	3.2	64445
9	3.7	57189
10	3.9	63218
11	4.0	55794
12	4.0	56957
13	4.1	57081
14	4.5	61111
15	4.9	67938
16	5.1	66029
17	5.3	83088
18	5.9	81363
19	6.0	93940
20	6.8	91738
21	7.1	98273
22	7.9	101302
23	8.2	113812
24	8.7	109431
25	9.0	105582
26	9.5	116969
27	9.6	112635
28	10.3	122391
29	10.5	121872

In [158...

```
salary=pd.read_csv(r"E:\Data Science & AI\Dataset files\Salary_Data.csv")
salary.head()
from sklearn.feature_selection import VarianceThreshold
```

```
vt=VarianceThreshold(threshold=0)
### Make sure before fitting the dataframe , do not include output column
x=salary.drop('YearsExperience',axis=1)
# x is self a data frame
vt.fit(x)
vt.variances_
vt.get_support()
cols=vt.get_feature_names_out()
x[cols]
```

Out[158...

	Salary
0	39343
1	46205
2	37731
3	43525
4	39891
5	56642
6	60150
7	54445
8	64445
9	57189
10	63218
11	55794
12	56957
13	57081
14	61111
15	67938
16	66029
17	83088
18	81363
19	93940
20	91738
21	98273
22	101302
23	113812
24	109431
25	105582
26	116969
27	112635
28	122391
29	121872

In [160...

```
from statsmodels.api import OLS
OLS(y_train,x_train).fit().summary()
```

Out[160...

OLS Regression Results

<b>Dep. Variable:</b>		Salary	<b>R-squared (uncentered):</b>		1.000
<b>Model:</b>		OLS	<b>Adj. R-squared (uncentered):</b>		1.000
<b>Method:</b>		Least Squares	<b>F-statistic:</b>		1.583e+32
<b>Date:</b>		Fri, 20 Sep 2024	<b>Prob (F-statistic):</b>		1.82e-310
<b>Time:</b>		21:48:30	<b>Log-Likelihood:</b>		479.24
<b>No. Observations:</b>		21	<b>AIC:</b>		-956.5
<b>Df Residuals:</b>		20	<b>BIC:</b>		-955.4
<b>Df Model:</b>		1			
<b>Covariance Type:</b>		nonrobust			
	<b>coef</b>	<b>std err</b>	<b>t</b>	<b>P&gt; t </b>	<b>[0.025 0.975]</b>
<b>Salary</b>	1.0000	7.95e-17	1.26e+16	0.000	1.000 1.000
<b>Omnibus:</b>		1.354	<b>Durbin-Watson:</b>		0.163
<b>Prob(Omnibus):</b>		0.508	<b>Jarque-Bera (JB):</b>		1.139
<b>Skew:</b>		0.522	<b>Prob(JB):</b>		0.566
<b>Kurtosis:</b>		2.539	<b>Cond. No.</b>		1.00

Notes:

[1]  $R^2$  is computed without centering (uncentered) since the model does not contain a constant.

[2] Standard Errors assume that the covariance matrix of the errors is correctly specified.

In [162...

```
import pickle
pickle.dump(LR,open('YearsExperience_model.pkl','wb'))
```

In [164...

```
# Loading model to compare the result
model=pickle.load(open('YearsExperience_model.pkl','rb'))
model
```

Out[164...

▼ LinearRegression ⓘ ?  
LinearRegression()

In [ ]: