

```
import numpy as np
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
plt.tight_layout()
import os
import statsmodels.formula.api as sm
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import train_test_split
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV
import warnings
```

<Figure size 640x480 with 0 Axes>

```
warnings.simplefilter(action='ignore', category=FutureWarning)
```

```
data_path = "/content/Advertising.csv"
df = pd.read_csv(data_path)
```

+ Code

+ Text

```
df.head()
```

| | Unnamed: 0 | TV | Radio | Newspaper | Sales |
|---|------------|-------|-------|-----------|-------|
| 0 | 1 | 230.1 | 37.8 | 69.2 | 22.1 |
| 1 | 2 | 44.5 | 39.3 | 45.1 | 10.4 |
| 2 | 3 | 17.2 | 45.9 | 69.3 | 9.3 |
| 3 | 4 | 151.5 | 41.3 | 58.5 | 18.5 |
| 4 | 5 | 180.8 | 10.8 | 58.4 | 12.9 |

```
df.columns
```

```
Index(['Unnamed: 0', 'TV', 'Radio', 'Newspaper', 'Sales'], dtype='object')
```

```
df.rename(columns={'Unnamed: 0': 'Index'}, inplace=True)
```

```
df.shape
```

```
(200, 5)
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 200 entries, 0 to 199
Data columns (total 5 columns):
#   Column      Non-Null Count  Dtype
---  ---
0   Index       200 non-null    int64
1   TV          200 non-null    float64
2   Radio       200 non-null    float64
3   Newspaper   200 non-null    float64
4   Sales       200 non-null    float64
dtypes: float64(4), int64(1)
memory usage: 7.9 KB
```

```
df
```

| | Index | TV | Radio | Newspaper | Sales |
|---|-------|-------|-------|-----------|-------|
| 0 | 1 | 230.1 | 37.8 | 69.2 | 22.1 |
| 1 | 2 | 44.5 | 39.3 | 45.1 | 10.4 |
| 2 | 3 | 17.2 | 45.9 | 69.3 | 9.3 |

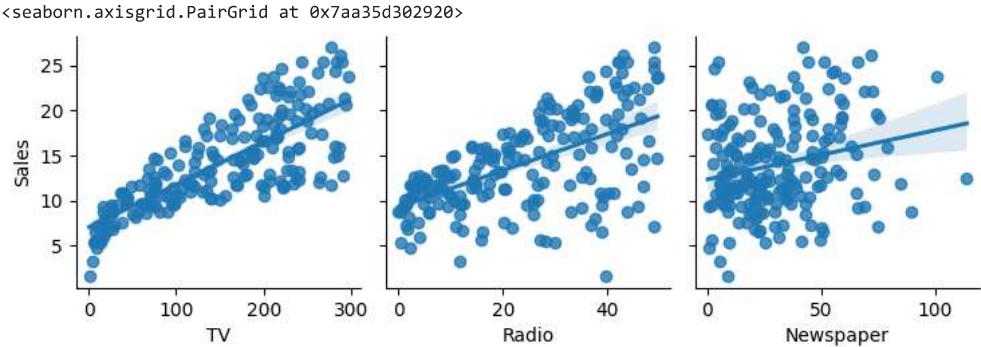
```
df.describe()
```

| | Index | TV | Radio | Newspaper | Sales |
|-------|------------|------------|------------|------------|------------|
| count | 200.000000 | 200.000000 | 200.000000 | 200.000000 | 200.000000 |
| mean | 100.500000 | 147.042500 | 23.264000 | 30.554000 | 14.022500 |
| std | 57.879185 | 85.854236 | 14.846809 | 21.778621 | 5.217457 |
| min | 1.000000 | 0.700000 | 0.000000 | 0.300000 | 1.600000 |
| 25% | 50.750000 | 74.375000 | 9.975000 | 12.750000 | 10.375000 |
| 50% | 100.500000 | 149.750000 | 22.900000 | 25.750000 | 12.900000 |
| 75% | 150.250000 | 218.825000 | 36.525000 | 45.100000 | 17.400000 |
| max | 200.000000 | 296.400000 | 49.600000 | 114.000000 | 27.000000 |

```
df.isnull().values.any()
df.isnull().sum()
```

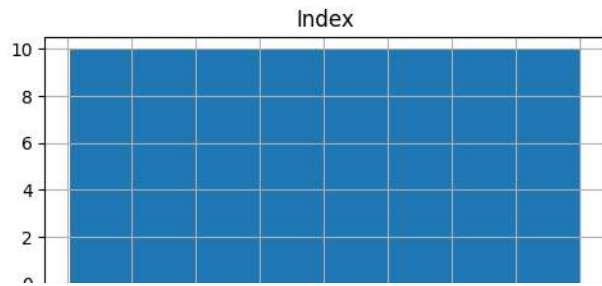
```
Index      0
TV         0
Radio      0
Newspaper  0
Sales      0
dtype: int64
```

```
sns.pairplot(df, x_vars=["TV", "Radio", "Newspaper"], y_vars="Sales", kind="reg")
```

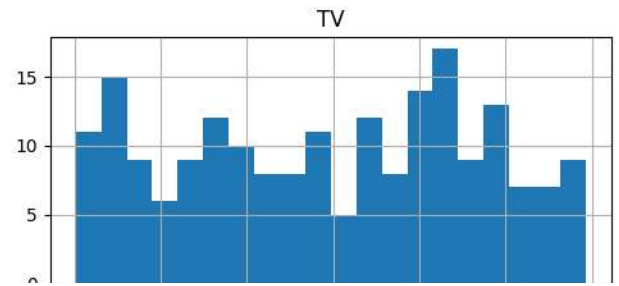


```
df.hist(bins=20, figsize=(13, 9))
```

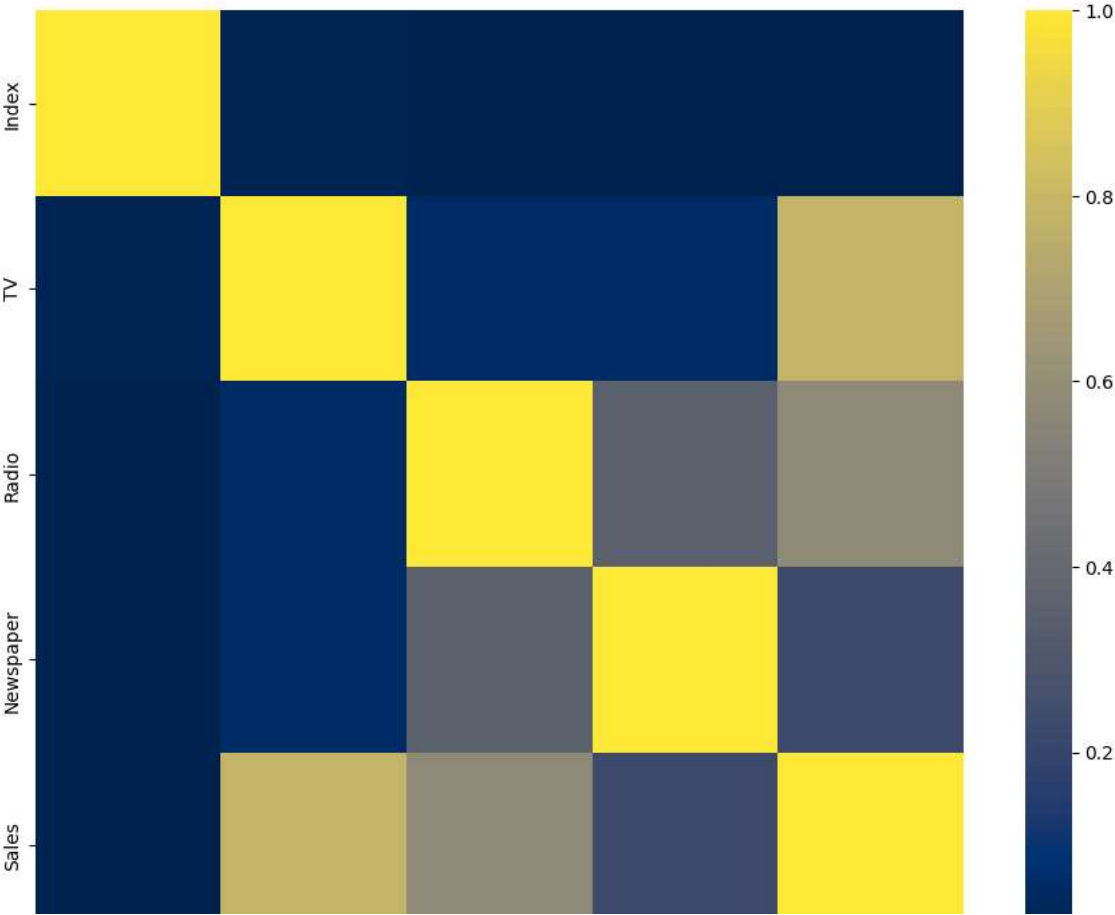
```
array([[<Axes: title={'center': 'Index'}>,
        <Axes: title={'center': 'TV'}>],
       [<Axes: title={'center': 'Radio'}>,
        <Axes: title={'center': 'Newspaper'}>],
       [<Axes: title={'center': 'Sales'}>, <Axes: >]], dtype=object)
```



```
sns.lmplot(x='TV', y='Sales', data=df)
sns.lmplot(x='Radio', y='Sales', data=df)
sns.lmplot(x='Newspaper', y='Sales', data=df)
```



```
corrmat = df.corr()
f, ax = plt.subplots(figsize=(12, 9))
sns.heatmap(corrmat, vmin=0, vmax=1, square=True, cmap="cividis", ax=ax)
plt.show()
```



```
X = df.drop('Sales', axis=1)
y = df[["Sales"]]

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.20, random_state=40)

lin_model = sm.ols(formula="Sales ~ TV + Radio + Newspaper", data=df).fit()

print(lin_model.params, "\n")

Intercept    2.938889
TV            0.045765
Radio         0.188530
Newspaper    -0.001037
dtype: float64

print(lin_model.summary())
```

| OLS Regression Results | | | | | | |
|------------------------|------------------|---------------------|----------|-------|--------|--------|
| ===== | | | | | | |
| Dep. Variable: | Sales | R-squared: | 0.897 | | | |
| Model: | OLS | Adj. R-squared: | 0.896 | | | |
| Method: | Least Squares | F-statistic: | 570.3 | | | |
| Date: | Fri, 15 Sep 2023 | Prob (F-statistic): | 1.58e-96 | | | |
| Time: | 05:47:22 | Log-Likelihood: | -386.18 | | | |
| No. Observations: | 200 | AIC: | 780.4 | | | |
| Df Residuals: | 196 | BIC: | 793.6 | | | |
| Df Model: | 3 | | | | | |
| Covariance Type: | nonrobust | | | | | |
| ===== | | | | | | |
| | coef | std err | t | P> t | [0.025 | 0.975] |
| ----- | | | | | | |
| Intercept | 2.9389 | 0.312 | 9.422 | 0.000 | 2.324 | 3.554 |
| TV | 0.0458 | 0.001 | 32.809 | 0.000 | 0.043 | 0.049 |
| Radio | 0.1885 | 0.009 | 21.893 | 0.000 | 0.172 | 0.206 |
| Newspaper | -0.0010 | 0.006 | -0.177 | 0.860 | -0.013 | 0.011 |
| ===== | | | | | | |
| Omnibus: | 60.414 | Durbin-Watson: | 2.084 | | | |
| Prob(Omnibus): | 0.000 | Jarque-Bera (JB): | 151.241 | | | |
| Skew: | -1.327 | Prob(JB): | 1.44e-33 | | | |
| Kurtosis: | 6.332 | Cond. No. | 454. | | | |
| ===== | | | | | | |

Notes:

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

```
models = [('LinearRegression', LinearRegression())]

for name, model in models:
    model.fit(X_train, y_train)
    y_pred = model.predict(X_test)
    rmse = np.sqrt(mean_squared_error(y_test, y_pred))
    print(f"{name}: RMSE = {rmse:.2f}")

    LinearRegression: RMSE = 2.36

new_data_1 = pd.DataFrame({'TV': [100], 'Radio': [50], 'Newspaper': [25]})
predicted_sales_1 = lin_model.predict(new_data_1)
print("Predicted Sales (Data 1):", predicted_sales_1)

    Predicted Sales (Data 1): 0      16.915917
    dtype: float64

new_data_2 = pd.DataFrame({'TV': [25], 'Radio': [63], 'Newspaper': [80]})
predicted_sales_2 = lin_model.predict(new_data_2)
print("Predicted Sales (Data 2):", predicted_sales_2)

    Predicted Sales (Data 2): 0      15.877397
    dtype: float64
```

