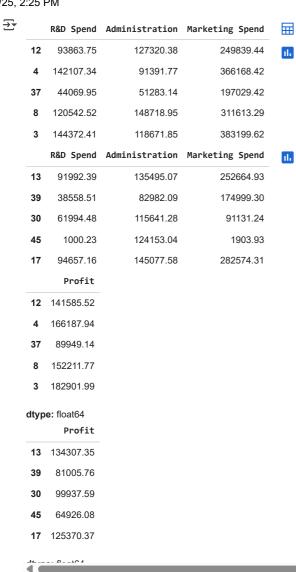
```
import pandas as pd
import numpy as np
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean absolute error, mean squared error, r2 score
from google.colab import files
# Use files.upload() to prompt the user to upload files
uploaded = files.upload()
data = pd.read_csv(next(iter(uploaded))) # e.g., "companies.csv"
print(data.head())
Choose Files 50_Startup...mpanies.csv
     • 50_Startups_companies.csv(text/csv) - 1964 bytes, last modified: 6/10/2025 - 100% done
     Saving 50_Startups_companies.csv to 50_Startups_companies.csv
       R&D Spend Administration Marketing Spend
     0 165349.20
                    136897.80
                                         471784.10 192261.83
     1 162597.70
                       151377.59
                                         443898.53 191792.06
                                         407934.54 191050.39
383199.62 182901.99
    2 153441.51
                      101145.55
118671.85
     3 144372.41
     4 142107.34
                        91391.77
                                         366168.42 166187.94
features = ['R&D Spend', 'Administration', 'Marketing Spend']
target = 'Profit'
X = data[features]
y = data[target]
display(X.head())
display(y.head())
        R&D Spend Administration Marketing Spend
      0 165349.20
                                          471784.10
                         136897.80
                                                       ıl.
      1 162597.70
                         151377.59
                                          443898.53
      2 153441.51
                         101145.55
                                          407934.54
      3 144372.41
                         118671.85
                                          383199.62
      4 142107.34
                          91391.77
                                          366168.42
           Profit
      0 192261.83
      1 191792.06
     2 191050.39
      3 182901.99
      4 166187.94
     dtype: float64
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
display(X_train.head())
display(X_test.head())
display(y_train.head())
display(y_test.head())
```



```
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_test_scaled = scaler.transform(X_test)
```

## Model-1: Linear Regression

## Model 2: Decision Tree Regression

```
from sklearn.tree import DecisionTreeRegressor
dt = DecisionTreeRegressor(random_state=42)
dt.fit(X_train_scaled, y_train)
y_pred_dt = dt.predict(X_test_scaled)
```

```
print("Decision Tree Metrics:")
print(f"MAE: {mean_absolute_error(y_test, y_pred_dt):.2f}")
\label{eq:print}  \texttt{print}(\texttt{f"RMSE: } \{\texttt{mean\_squared\_error}(\texttt{y\_test}, \ \texttt{y\_pred\_dt}) :. 2\texttt{f}\}") 
print(f"R2: {r2_score(y_test, y_pred_dt):.4f}")
→ Decision Tree Metrics:
     MAE: 13755.66
     RMSE: 400026479.25
     R2: 0.5060
Model 3: Random Forest Regression
from sklearn.ensemble import RandomForestRegressor
rf = RandomForestRegressor(random_state=42)
rf.fit(X_train_scaled, y_train)
y pred rf = rf.predict(X test scaled)
print("Random Forest Metrics:")
print(f"MAE: {mean_absolute_error(y_test, y_pred_rf):.2f}")
print(f"RMSE: {mean_squared_error(y_test, y_pred_rf):.2f}")
print(f"R2: {r2_score(y_test, y_pred_rf):.4f}")
Random Forest Metrics:
     MAE: 6437.50
     RMSE: 72625008.62
     R<sup>2</sup>: 0.9103
Model 4: Support Vector Regression (SVR)
from sklearn.svm import SVR
svr = SVR()
svr.fit(X_train_scaled, y_train)
y_pred_svr = svr.predict(X_test_scaled)
print("SVR Metrics:")
print(f"MAE: {mean_absolute_error(y_test, y_pred_svr):.2f}")
print(f"RMSE: \{mean\_squared\_error(y\_test, y\_pred\_svr):.2f\}")
print(f"R2: {r2_score(y_test, y_pred_svr):.4f}")

→ SVR Metrics:
     MAF: 22846.73
     RMSE: 955620367.28
     R2: -0.1801
results = {
    "Model": ["Linear Regression", "Decision Tree", "Random Forest", "SVR"],
    "MAE": [6979.15, 13755.66, 6437.50, 22846.73],
    "RMSE": [8995.91, 400026479.25, 72625008.62, 955620367.28],
    "R<sup>2</sup>": [0.9001, 0.5060, 0.9103, -0.1801]
results_df = pd.DataFrame(results)
print("Model Comparison:")
print(results_df.sort_values(by="R2", ascending=False))
→ Model Comparison:
                     Model
                                 MAE
                                               RMSF
     2
            Random Forest
                            6437.50 7.262501e+07 0.9103
                            6979.15 8.995910e+03 0.9001
       Linear Regression
            Decision Tree 13755.66 4.000265e+08 0.5060
                      SVR 22846.73 9.556204e+08 -0.1801
best model = max(results df.to dict('records'), key=lambda x: x['R2'])
print(f"Best Model: \{best\_model['Model']\} \ (R^2 = \{best\_model['R^2']:.4f\})")

    Best Model: Random Forest (R² = 0.9103)
import pandas as pd
import sklearn
import numpy as np
import xgboost
import matplotlib
print(f"pandas: {pd. version }")
print(f"scikit-learn: {sklearn.__version__}}")
```

print(f"numpy: {np.\_\_version\_\_}")
print(f"xgboost: {xgboost.\_\_version\_\_}")
print(f"matplotlib: {matplotlib.\_\_version\_\_}")

pandas: 2.2.2 scikit-learn: 1.6.1 numpy: 2.0.2 xgboost: 2.1.4 matplotlib: 3.10.0