

## Linear Regression

```
In [27]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

In [6]: df1 = pd.read_csv("D:\manish\Python\GROW AI\python\Electric_Vehicle_Population_Data.csv")

Out[6]:
```

VIN (1-10)	County	City	State	Postal Code	Model Year	Make	Model	Electric Vehicle Type	Clean Alternative Fuel Vehicle (CAVF) Eligibility	Electric Range	Legislative District	DOL	Vehicle ID	Vehicle Location	Electric Utility	2020 Census Tract
0 5YJYGDDE8L	Thurston	Tumwater	WA	98501.0	2020	TESLA	MODEL Y	Battery Electric Vehicle (BEV)	Clean Alternative Fuel Vehicle Eligible	291.0	35.0	124633715	POINT	PUGET SOUND ENERGY INC	5.306701e+10	
1 5YXCAE2XJ	Snohomish	Bothell	WA	98021.0	2018	TESLA	MODEL X	Battery Electric Vehicle (BEV)	Clean Alternative Fuel Vehicle Eligible	238.0	1.0	474826075	POINT	PUGET SOUND ENERGY INC	5.306105e+10	
2 5YJ3E1EBXK	King	Kent	WA	98031.0	2019	TESLA	MODEL 3	Battery Electric Vehicle (BEV)	Clean Alternative Fuel Vehicle Eligible	220.0	47.0	280302733	POINT	PUGET SOUND ENERGY INC	5.303303e+10	
3 7SAYGDEEAT	King	Issaquah	WA	98027.0	2026	TESLA	MODEL Y	Battery Electric Vehicle (BEV)	Eligibility unknown as battery range has not b...	0.0	41.0	280786565	POINT	PUGET SOUND ENERGY INC	5.303302e+10	
4 WAUUPBFF9G	King	Seattle	WA	98103.0	2016	AUDI	A3	Plug-in Hybrid Electric Vehicle (PHEV)	Not eligible due to low battery range	16.0	43.0	19898891	POINT	CITY OF SEATTLE - (WA) CITY OF TACOMA - (WA)	5.303300e+10	
... ... ... ... ... ...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
270257 1C4RJXN60R	Pierce	Joint Base Lewis McChord	WA	98433.0	2024	JEEP	WRANGLER	Plug-in Hybrid Electric Vehicle (PHEV)	Not eligible due to low battery range	21.0	28.0	266021122	POINT	PUGET SOUND ENERGY INC	5.305307e+10	
270258 1C4J3XR66N	Mason	Hoodsport	WA	98548.0	2022	JEEP	WRANGLER	Plug-in Hybrid Electric Vehicle (PHEV)	Not eligible due to low battery range	22.0	35.0	282429398	POINT	BONNEVILLE POWER ADMINISTRATION	5.304596e+10	
270259 7SAYGDEEXP	Pierce	Tacoma	WA	98406.0	2023	TESLA	MODEL Y	Battery Electric Vehicle (BEV)	Eligibility unknown as battery range has not b...	0.0	27.0	228485085	POINT	BONNEVILLE POWER ADMINISTRATION	5.305306e+10	
270260 5YJYGDDE2M	Snohomish	Bothell	WA	98021.0	2021	TESLA	MODEL Y	Battery Electric Vehicle (BEV)	Eligibility unknown as battery range has not b...	0.0	1.0	282699217	POINT	PUGET SOUND ENERGY INC	5.306105e+10	
270261 JN1BF0BA5P	Chelan	Wenatchee	WA	98801.0	2023	NISSAN	ARIYA	Battery Electric Vehicle (BEV)	Eligibility unknown as battery range has not b...	0.0	12.0	261475224	POINT	PUD NO 1 OF CHELAN COUNTY	5.300796e+10	

270262 rows × 16 columns

How can we use Linear Regression to predict the Electric Range of a vehicle?

```
In [8]: df = df[['Model Year', 'Electric Range', 'Legislative District']].dropna()
y = df['Electric Range']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
model = LinearRegression().fit(X_train, y_train)
y_pred = model.predict(X_test)

print("Coefficients: ", model.coef_)
print("Intercept: ", model.intercept_)
print("R^2: ", r2_score(y_test, y_pred))
print("Predictions: ", y_pred)

Coefficients: [-14.14570503  0.02866648]
Intercept: 28641.659603274988
R^2: 0.2958970015663
Predictions: [ 12.04246216 40.36269818 25.24323268 ... -16.22012193 -2.21774932
 39.30703825]
```

```
In [10]: df = df[['Model Year', 'Electric Range', 'Legislative District']].dropna()
y = df['Electric Range']

# Train model
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3, random_state=42)
model = LinearRegression().fit(X_train, y_train)
y_pred_train = model.predict(X_train)
y_pred_test = model.predict(X_test)

# Create visualizations
fig, axes = plt.subplots(2, 2, figsize=(15, 12))
fig.suptitle('Linear Regression: Predicting Electric Range', fontsize=16)

# 1. Actual vs Predicted (Training)
axes[0,0].scatter(y_train, y_pred_train, alpha=0.7, color='blue')
axes[0,0].plot([y_train.min(), y_train.max()], [y_train.min(), y_train.max()], 'r--', lw=2)
axes[0,0].set_xlabel('Actual Electric Range')
axes[0,0].set_ylabel('Predicted Electric Range')
axes[0,0].set_title('Training: R^2 = %r' % (r2_score(y_train, y_pred_train),))

# 2. Actual vs Predicted (Testing)
axes[0,1].scatter(y_test, y_pred_test, alpha=0.7, color='green')
axes[0,1].plot([y_test.min(), y_test.max()], [y_test.min(), y_test.max()], 'r--', lw=2)
axes[0,1].set_xlabel('Actual Electric Range')
axes[0,1].set_ylabel('Predicted Electric Range')
axes[0,1].set_title('Testing: R^2 = %r' % (r2_score(y_test, y_pred_test),))

# 3. Residuals Plot
residuals = y_test - y_pred_test
axes[1,0].scatter(y_test, residuals, alpha=0.7, color='orange')
axes[1,0].plot([-100, 100], [0, 0], 'r--', linestyle='--')
axes[1,0].set_xlabel('Actual Electric Range')
axes[1,0].set_ylabel('Residuals')
axes[1,0].set_title('Testing: R^2 = %r' % (r2_score(y_test, y_pred_test),))

# 4. Model Year vs Electric Range with regression line
scatters = axes[1,1].scatter(df['Model Year'], df['Electric Range'],
                           c=df['Legislative District'], s=100, alpha=0.7)
model_year_line = axes[1,1].scatter(df['Model Year'], df['Electric Range'],
                                    c=df['Legislative District'], s=100, alpha=0.7)
model = LinearRegression().fit(df['Model Year'], df['Electric Range'])
model_year_line = axes[1,1].plot(df['Model Year'], model.predict(df['Model Year']),
                                 color='red', s=100, marker='x', label='Predictions')
axes[1,1].set_xlabel('Model Year')
axes[1,1].set_ylabel('Electric Range')
axes[1,1].set_title('Model Year vs Electric Range')
axes[1,1].legend()

plt.tight_layout()
plt.show()
```



C:\Users\LENOVO\anaconda3\lib\site-packages\c:\Python\3.7\lib\site-packages\ipython\core\pylabtools.py:170: UserWarning: Creating legend with loc="best" can be slow with large amounts of data.

fig.canvas.print\_figure(fig.get\_file\_name(), \*\*kw)

Linear Regression: Predicting Electric Range

Training: R<sup>2</sup> = 0.297

Testing: R<sup>2</sup> = 0.296

Actual Electric Range

Predicted Electric Range

Residuals Plot

Model Year vs Electric Range

Model Year

Electric Range

Model Coefficients: Model Year=-14.15, District=0.03

Intercept: 28641.66

R<sup>2</sup>: 0.296

RMSE: 66.52

What independent variables (features) can be used to predict Electric Range? (e.g., Model Year, Base MSRP, Make)

```
In [11]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
from sklearn.linear_model import LinearRegression
from sklearn.ensemble import RandomForestRegressor
from sklearn.preprocessing import LabelEncoder
```

# Potential numeric features for 'Electric Range prediction'

categorical\_features = ['Model Year', 'Legislative District', 'Clear Alternative Fuel Vehicle (CAVF) Eligibility']

# Create feature matrix (encode categorical, handle missing values)

df\_numeric = df[numeric\_features + ['Electric Range']].dropna()

X\_numeric = df\_numeric[numerical\_features]

# Encode categorical features (one-hot for small dataset)

df\_cat = pd.get\_dummies(df[categorical\_features + ['Electric Range']].dropna(), drop\_first=True)

X\_cat = df\_cat.drop(['Electric Range'], axis=1)

print("Available Numeric Features: ", numeric\_features)

print("Available Categorical Features: ", categorical\_features)

Available Numeric Features: ['Model Year', 'Legislative District', 'Postal Code']

Available Categorical Features: ['Make', 'Model', 'Electric Vehicle Type', 'Clean Alternative Fuel Vehicle (CAVF) Eligibility']

Dataset shape after cleaning: (269608, 4)

How do we handle categorical variables like Make and Model in regression analysis?

```
In [14]: import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score
```

# Clean up categories

df\_cat = df\_cat[['Model Year', 'Make', 'Model', 'Electric Vehicle Type', 'Clean Alternative Fuel Vehicle (CAVF) Eligibility']]

df\_cat['Model Year'] = df\_cat['Model Year'].replace('2018', '2019')

df\_cat['Model'] = df\_cat['Model'].replace('2018', '2019')

df\_cat['Electric Vehicle Type'] = df\_cat['Electric Vehicle Type'].replace('2018', '2019')

df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'] = df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'].replace('2018', '2019')

df\_cat['Legislative District'] = df\_cat['Legislative District'].replace('2018', '2019')

df\_cat['Postal Code'] = df\_cat['Postal Code'].replace('2018', '2019')

df\_cat['Model Year'] = df\_cat['Model Year'].replace('2018', '2019')

df\_cat['Make'] = df\_cat['Make'].replace('2018', '2019')

df\_cat['Model'] = df\_cat['Model'].replace('2018', '2019')

df\_cat['Electric Vehicle Type'] = df\_cat['Electric Vehicle Type'].replace('2018', '2019')

df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'] = df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'].replace('2018', '2019')

df\_cat['Legislative District'] = df\_cat['Legislative District'].replace('2018', '2019')

df\_cat['Postal Code'] = df\_cat['Postal Code'].replace('2018', '2019')

df\_cat['Model Year'] = df\_cat['Model Year'].replace('2018', '2019')

df\_cat['Make'] = df\_cat['Make'].replace('2018', '2019')

df\_cat['Model'] = df\_cat['Model'].replace('2018', '2019')

df\_cat['Electric Vehicle Type'] = df\_cat['Electric Vehicle Type'].replace('2018', '2019')

df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'] = df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'].replace('2018', '2019')

df\_cat['Legislative District'] = df\_cat['Legislative District'].replace('2018', '2019')

df\_cat['Postal Code'] = df\_cat['Postal Code'].replace('2018', '2019')

df\_cat['Model Year'] = df\_cat['Model Year'].replace('2018', '2019')

df\_cat['Make'] = df\_cat['Make'].replace('2018', '2019')

df\_cat['Model'] = df\_cat['Model'].replace('2018', '2019')

df\_cat['Electric Vehicle Type'] = df\_cat['Electric Vehicle Type'].replace('2018', '2019')

df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'] = df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'].replace('2018', '2019')

df\_cat['Legislative District'] = df\_cat['Legislative District'].replace('2018', '2019')

df\_cat['Postal Code'] = df\_cat['Postal Code'].replace('2018', '2019')

df\_cat['Model Year'] = df\_cat['Model Year'].replace('2018', '2019')

df\_cat['Make'] = df\_cat['Make'].replace('2018', '2019')

df\_cat['Model'] = df\_cat['Model'].replace('2018', '2019')

df\_cat['Electric Vehicle Type'] = df\_cat['Electric Vehicle Type'].replace('2018', '2019')

df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'] = df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'].replace('2018', '2019')

df\_cat['Legislative District'] = df\_cat['Legislative District'].replace('2018', '2019')

df\_cat['Postal Code'] = df\_cat['Postal Code'].replace('2018', '2019')

df\_cat['Model Year'] = df\_cat['Model Year'].replace('2018', '2019')

df\_cat['Make'] = df\_cat['Make'].replace('2018', '2019')

df\_cat['Model'] = df\_cat['Model'].replace('2018', '2019')

df\_cat['Electric Vehicle Type'] = df\_cat['Electric Vehicle Type'].replace('2018', '2019')

df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'] = df\_cat['Clean Alternative Fuel Vehicle (CAVF) Eligibility'].replace('2018', '2019')

df\_cat['Legislative District'] = df\_cat['Legislative District'].replace('2018', '2019')

df\_cat['Postal Code'] = df\_cat['Postal Code'].replace('2018', '2019')

df\_cat['Model Year'] = df\_cat['Model Year'].replace('2018', '2019')

df\_cat['Make'] = df\_cat['Make'].replace('2018', '2019')

df\_cat['Model'] = df\_cat['Model'].replace('2018', '2019')

df\_cat['Electric Vehicle Type'] = df\_cat['Electric Vehicle Type'].replace('2018', '2019')