

SC310005 Artificial Intelligence

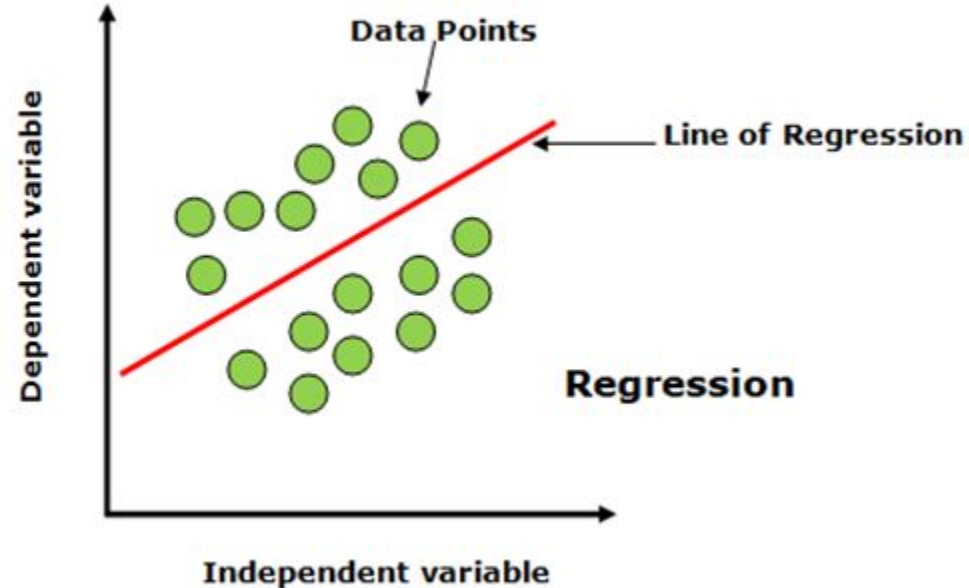
Lecture 7: Regression

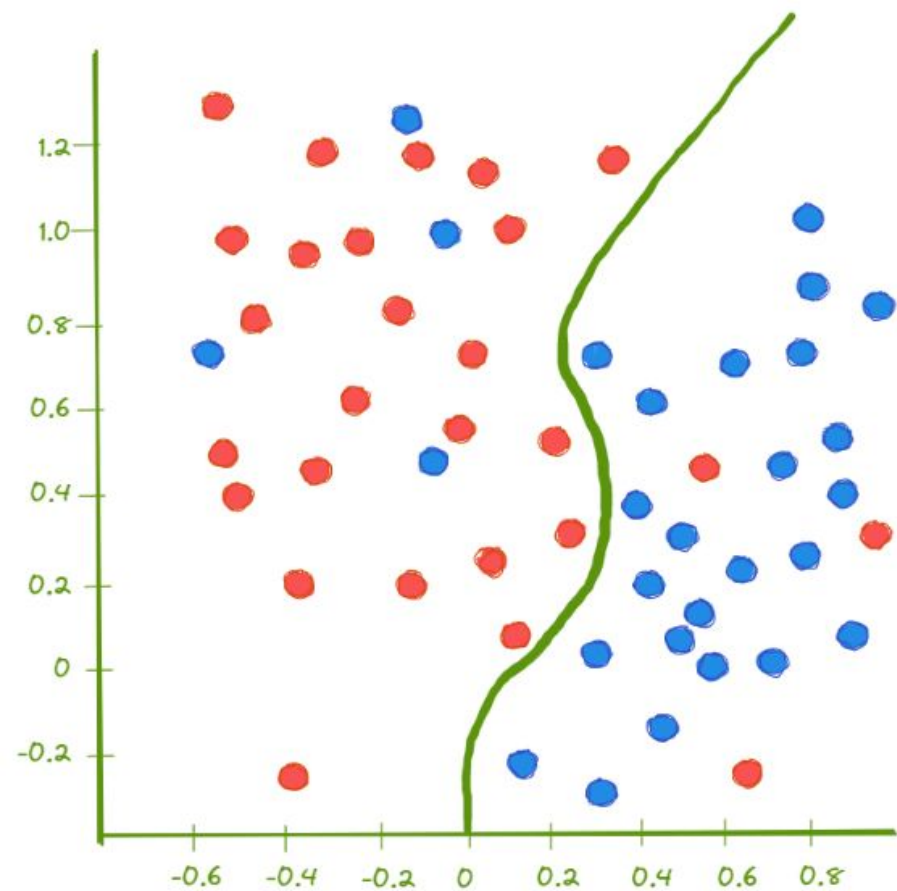
teerapong.pa@chula.ac.th

Reference

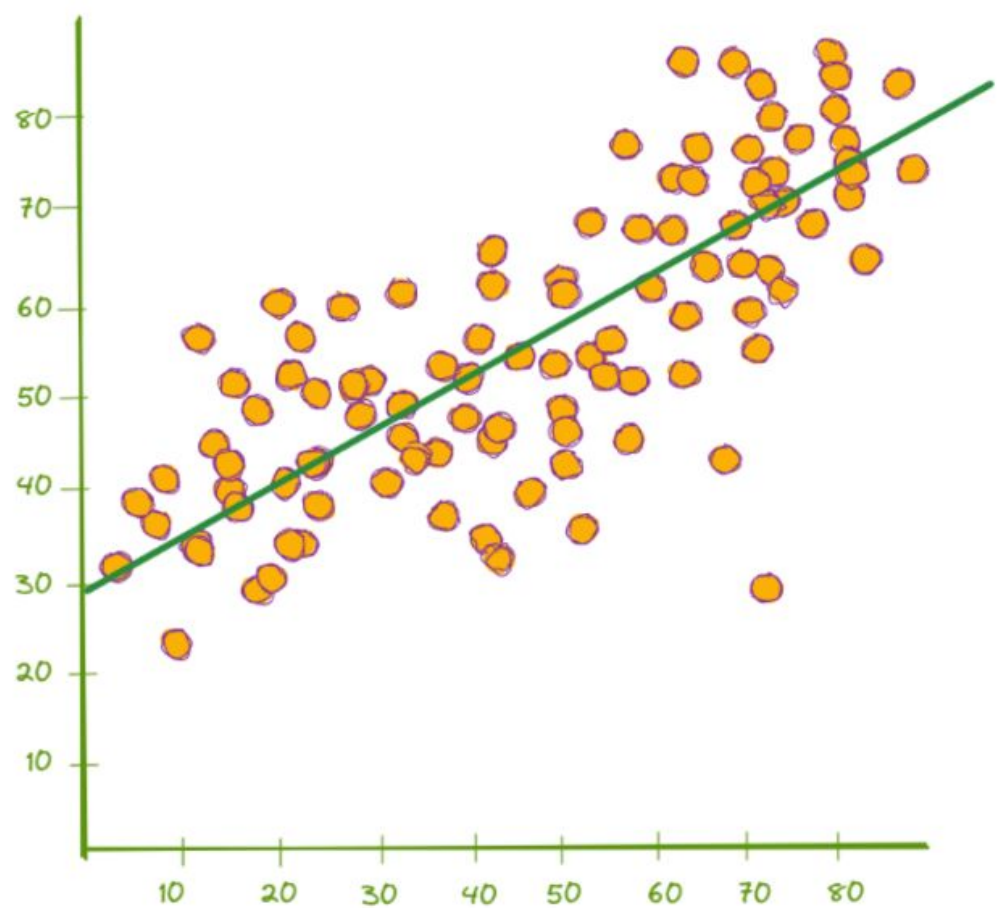
- <https://www.linkedin.com/pulse/understanding-linear-regression-machine-learning-tool-jadhav>
- <https://medium.com/@rndayala/linear-regression-a00514bc45b0>
- <https://medium.com/@polanitzer/the-minimum-mean-absolute-error-mae-challenge-928dc081f031>
- <https://www.investopedia.com/terms/p/positive-correlation.asp>
- <https://www.linkedin.com/pulse/correlation-vs-causality-krishnakumar-ramanathan>

Understanding Linear Regression in Machine Learning: A Fundamental Tool for Predictive Modeling



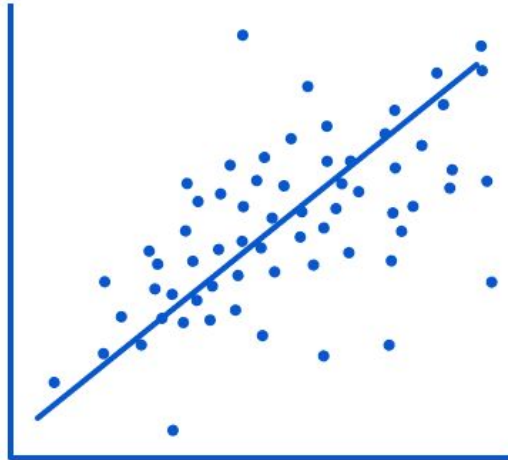


classification

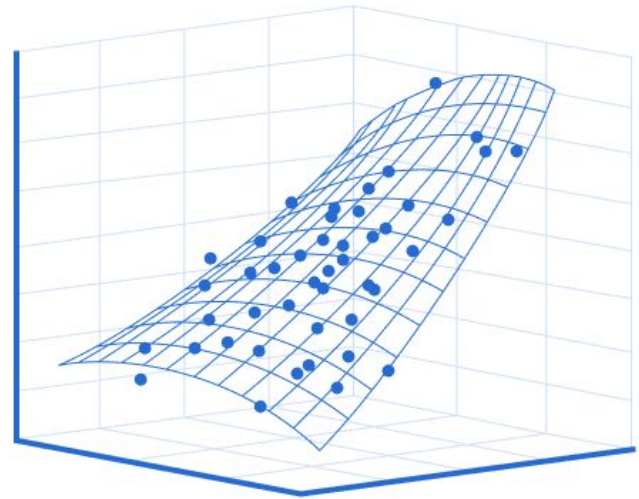


regression

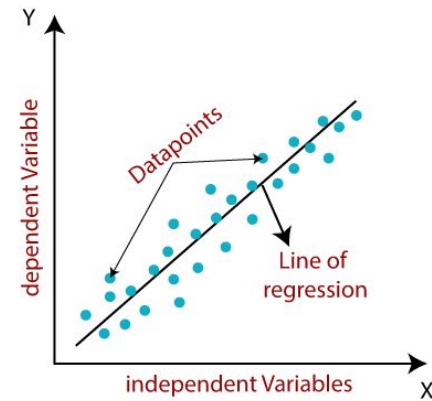
Simple Linear Regression



Multiple Linear Regression



Linear Regression (Equation)



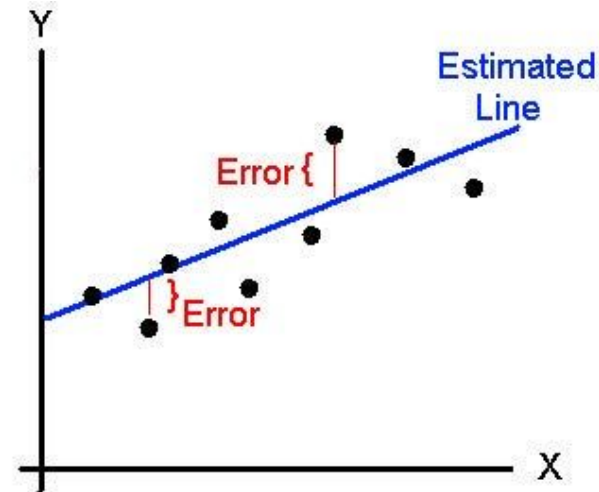
Estimated
(or predicted)
Y value for
observation i

Estimate of
the regression
intercept

Estimate of the
regression slope

Value of X for
observation i

$$\hat{Y}_i = b_0 + b_1 X_i$$



single value of dependent variable

slope

single value of independent variable

y-intercept

$$y = mx + b$$

all observed values for dependent variable

y-intercept a.k.a "bias"

slope a.k.a. "coefficient"

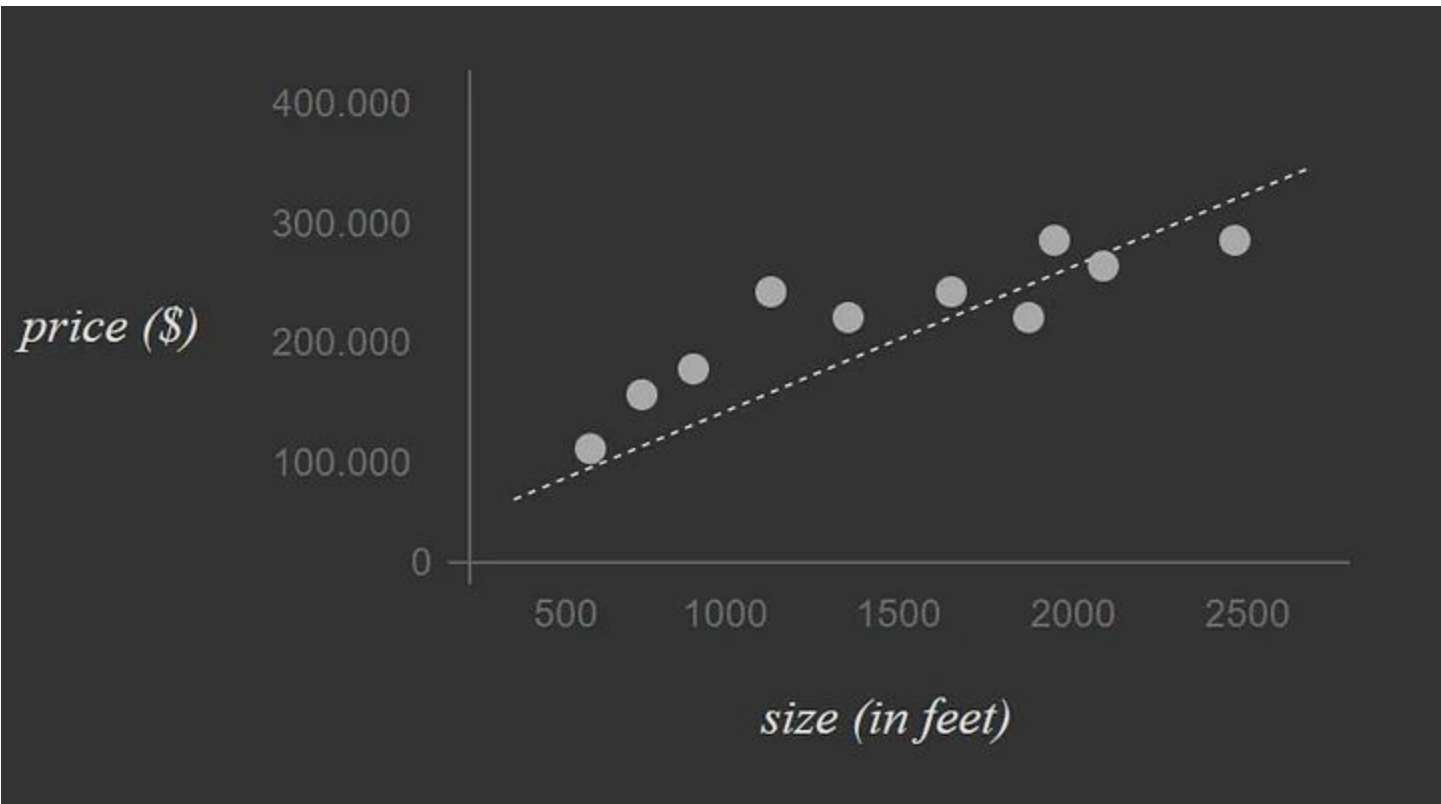
all observed values of independent variable

error*

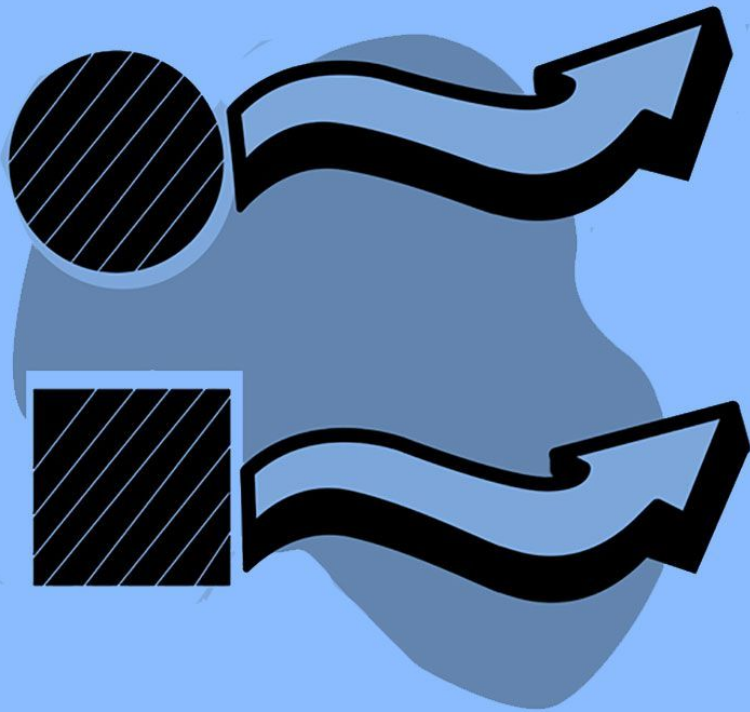
$$Y = \beta_0 + \beta_1 X + \epsilon$$

* additional term

α



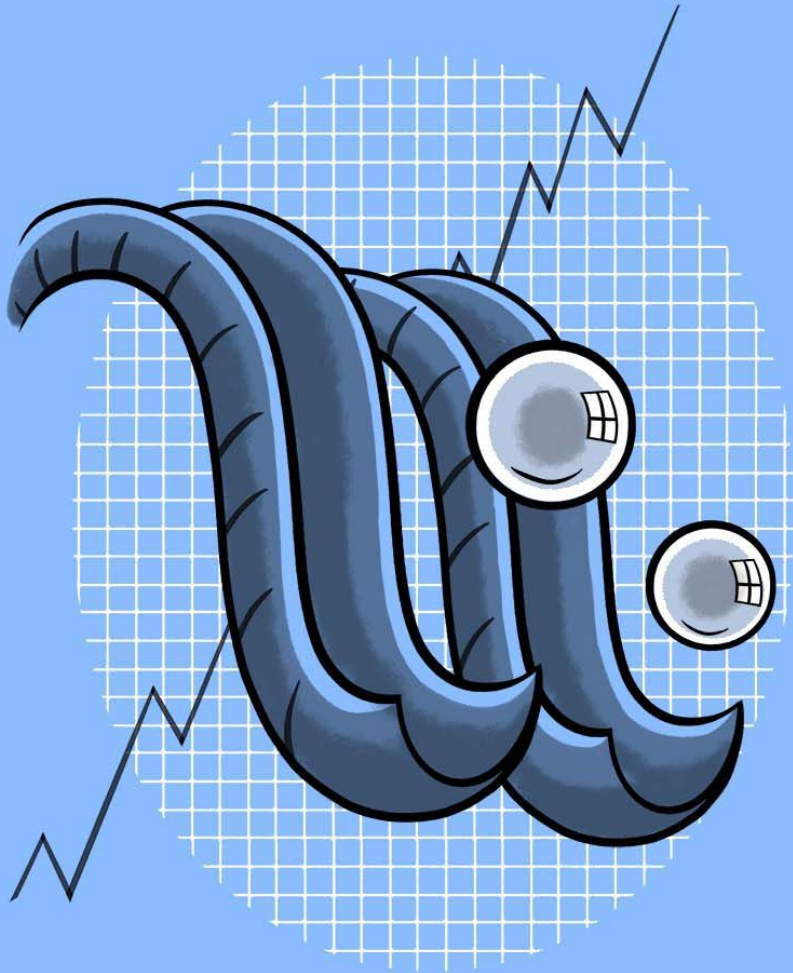
size (feet ²) (x)	price (\$) (y)
815	165,000
1510	310,000
2100	410,000
...	...



Correlation

[, kor-ə-'lā-shən]

A statistic that measures the degree to which two securities move in relation to each other.

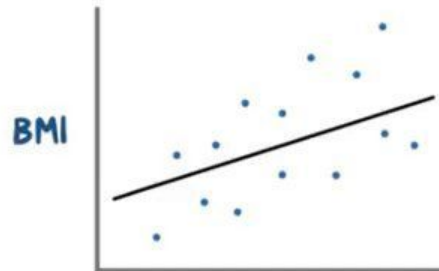


Positive Correlation

["pä-zə-tiv ,kôr-ə-'lā-shən]

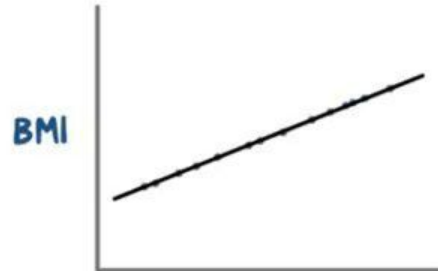
A relationship between two variables that move in the same direction.

SCATTERPLOT



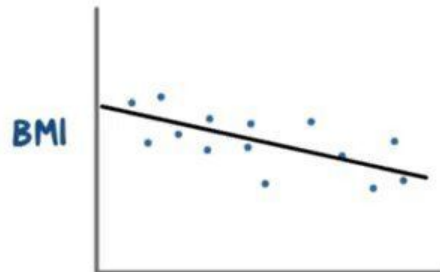
**POSITIVE
CORRELATION**

 **NUMBER of
BEVERAGES**



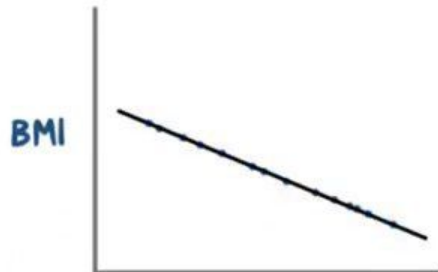
**PERFECT POSITIVE
CORRELATION**

 **NUMBER of
BEVERAGES**



**NEGATIVE
CORRELATION**

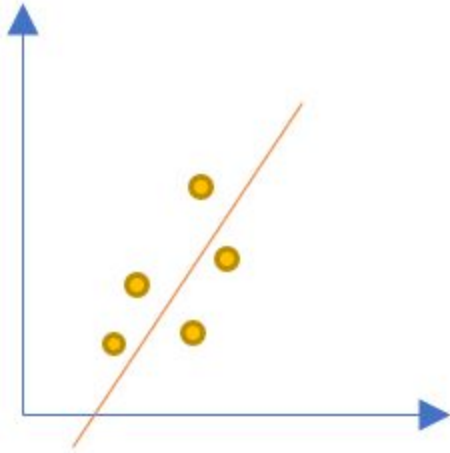
 **NUMBER of
BEVERAGES**



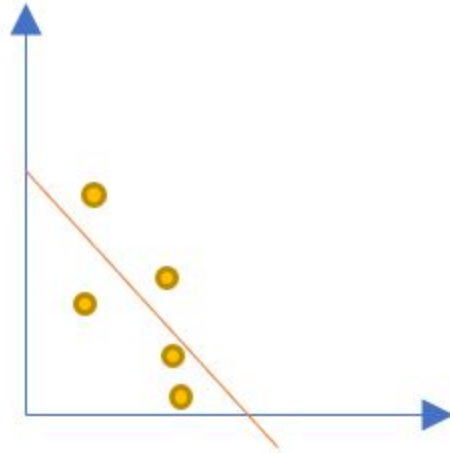
**PERFECT NEGATIVE
COR**

 **NUMBER of
BEVERAGES**

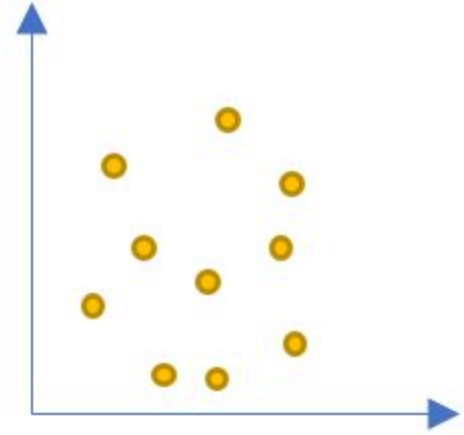
Types of Correlation Coefficients



Positive



Negative



No Correlation

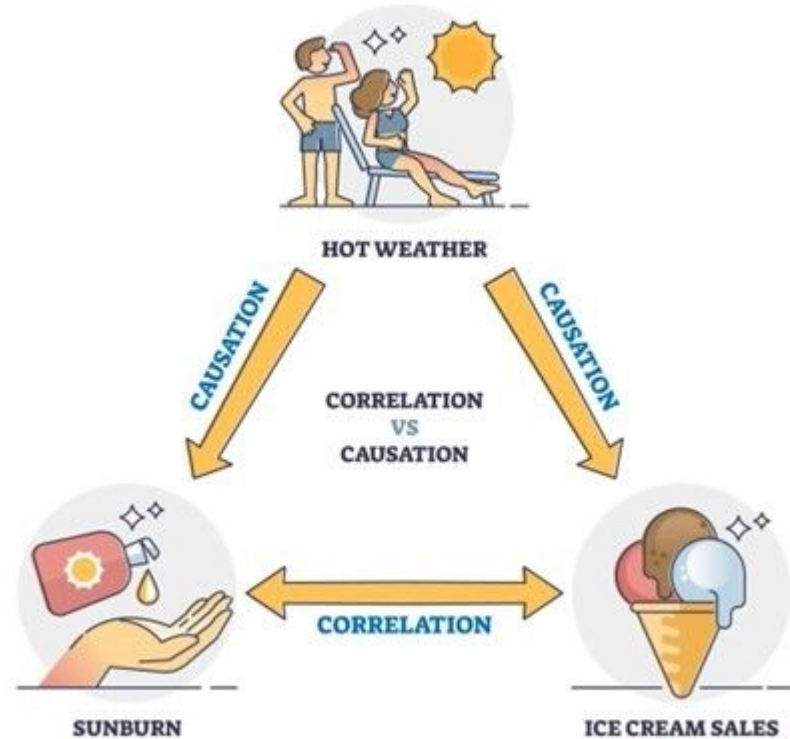
Correlation Vs. Causality!

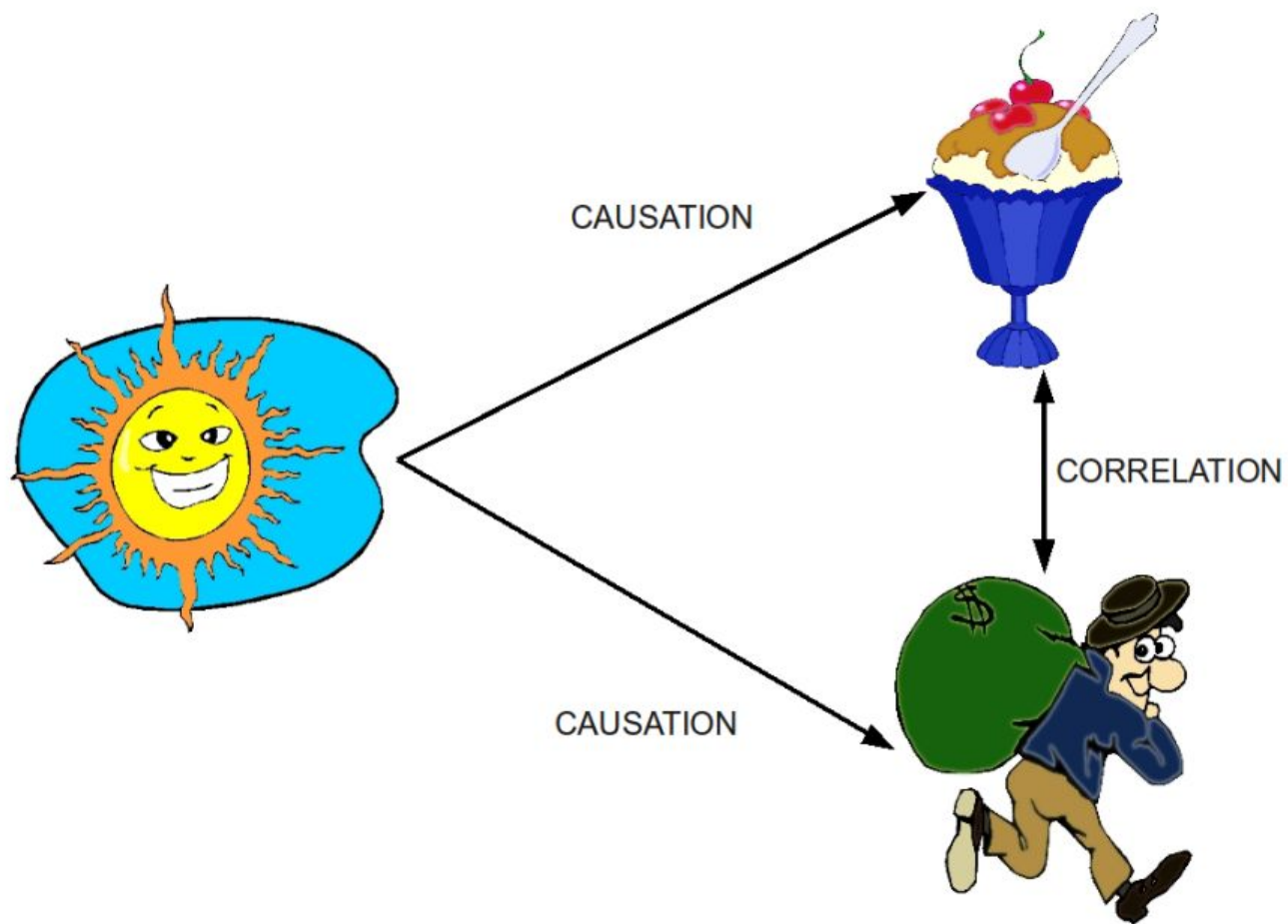
© MARK ANDERSON

WWW.ANDERSTOONS.COM

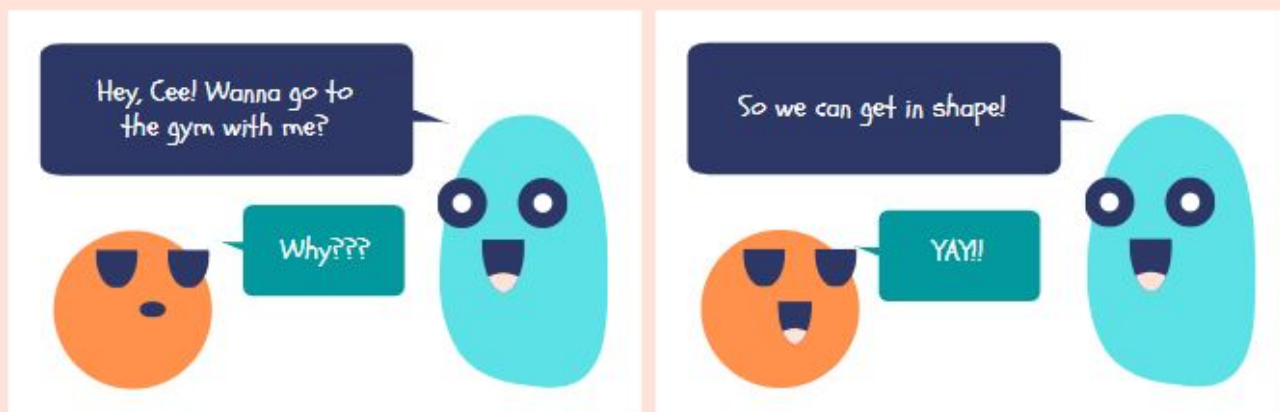


"It's important to remember that correlation does not imply causation. Besides, we all know it was Brian."

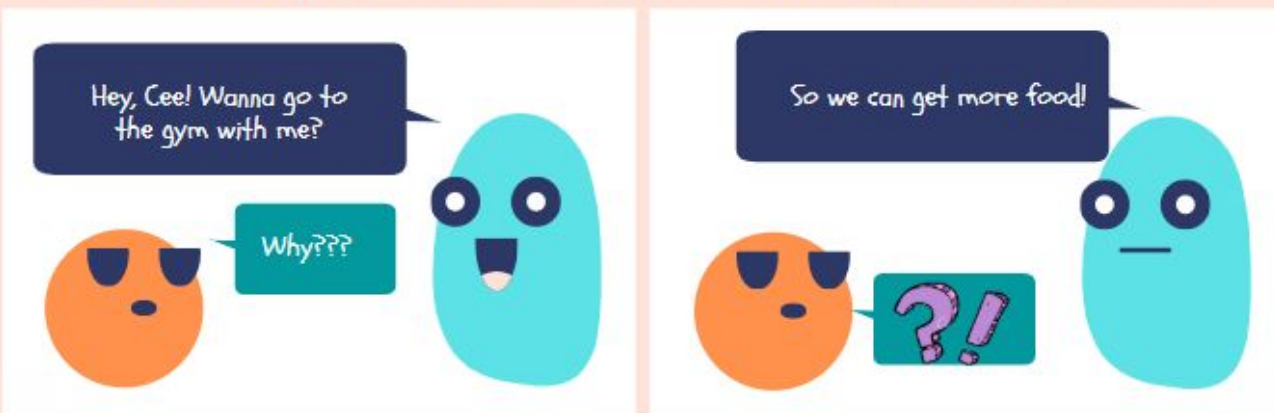




CORRELATION WITH CAUSATION



CORRELATION WITHOUT CAUSATION



Pandas `get_dummies` (One-Hot Encoding)

Sample:

Pandas `get_dummies` (One-Hot Encoding)

City		Houston	Rome	Madrid	London
Houston	one-hot encoding →	1	0	0	0
Rome		0	1	0	0
Madrid		0	0	1	0
London		0	0	0	1

Minimum Mean Absolute Error (MAE)

The diagram illustrates the Minimum Mean Absolute Error (MAE) formula with the following components and annotations:

- MAE**: The Minimum Mean Absolute Error metric.
- $=$** : The equals sign.
- $\frac{1}{n}$** : A blue box containing the fraction 1 over n. A blue line points to it from the text "Divide by the total number of data points".
- Σ** : The summation symbol. A black line points to it from the text "Sum of".
- $|$** : The absolute value bars.
- y** : A green box containing the variable y. A green line points to it from the text "Actual output value".
- $-$** : The minus sign.
- \hat{y}** : An orange box containing the predicted value y-hat. A yellow line points to it from the text "Predicted output value".
- $|$** : The absolute value bars.
- The absolute value of the residual**: A black bracket underneath the entire expression $|y - \hat{y}|$ points to it from this text.

$$MAE = \frac{1}{n} \sum |y - \hat{y}|$$

$$Error = Y_t - \hat{Y}_t$$

Mean squared error

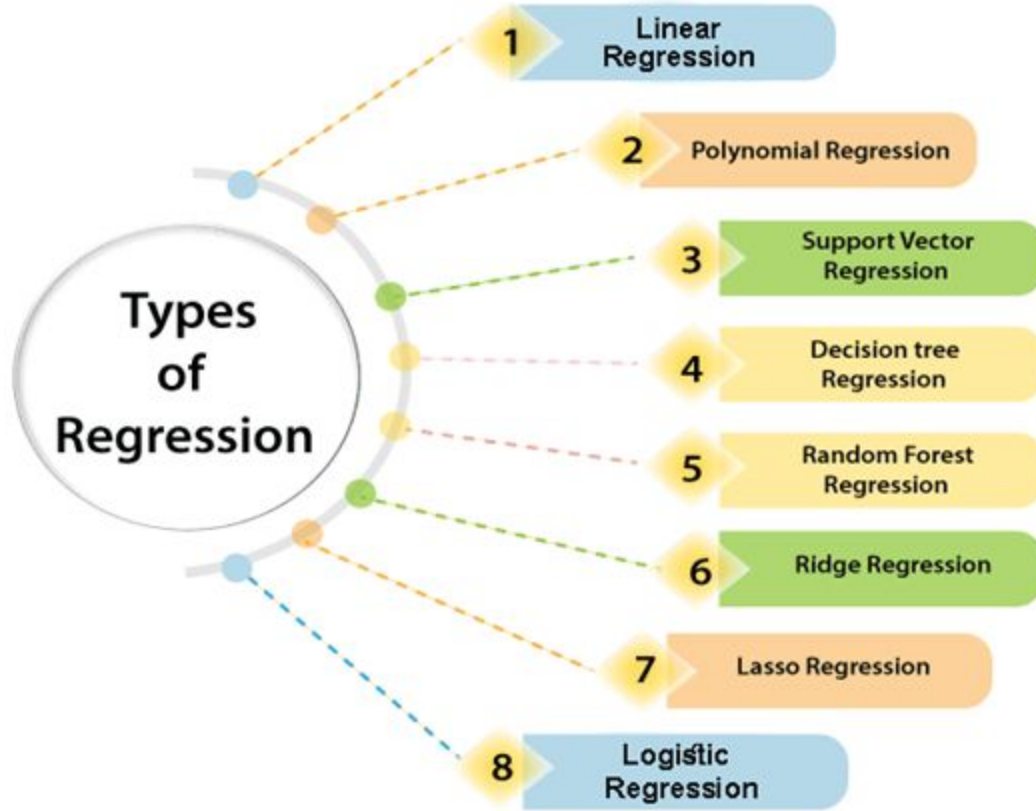
$$MSE = \frac{1}{n} \sum_{t=1}^n e_t^2$$

Root mean squared error

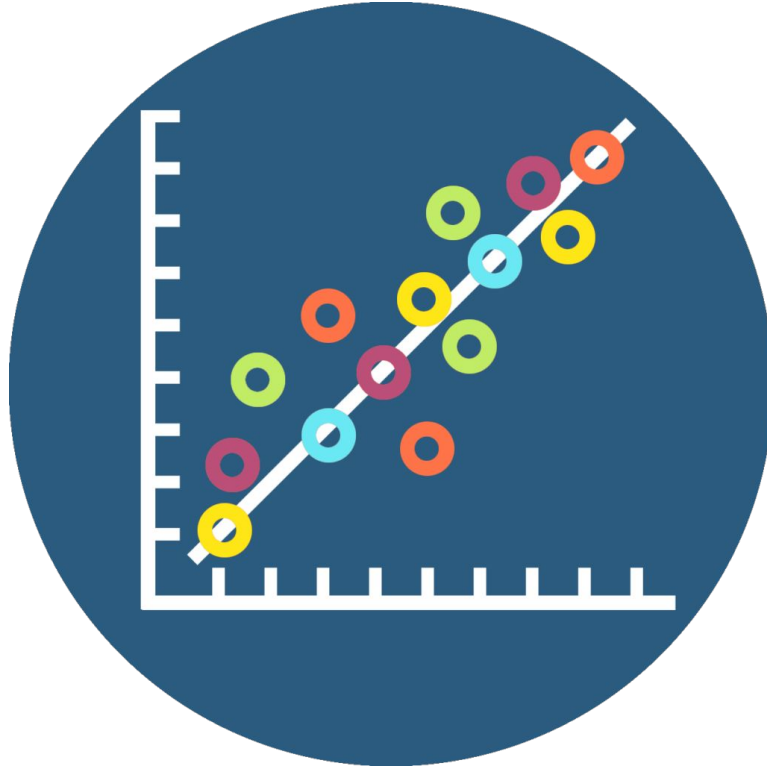
$$RMSE = \sqrt{\frac{1}{n} \sum_{t=1}^n e_t^2}$$

Mean absolute error

$$MAE = \frac{1}{n} \sum_{t=1}^n |e_t|$$



Let's Start Coding: Regression Model



✓ Import necessary libraries:

✓
1s

```
[2] 1 import pandas as pd
      2 import numpy as np
      3
      4 import matplotlib.pyplot as plt
      5 import seaborn as sns
      6
      7 from sklearn.model_selection import train_test_split
      8 from sklearn.linear_model import LinearRegression
      9 from sklearn.metrics import mean_absolute_error, mean_squared_error
     10 import statsmodels.api as sm
```

✓ Load the dataset:

✓
0s

```
[3] 1 df = pd.read_csv('house_prices_dataset.csv')
```


✓ EDA (Exploratory Data Analysis)

✓
0s

```
[4] 1 df.head()
```

	Id	MSSubClass	MSZoning	LotFrontage	LotArea	Street	Alley	LotShape	LandContour
0	1	60	RL	65.0	8450	Pave	NaN	Reg	Lvl
1	2	20	RL	80.0	9600	Pave	NaN	Reg	Lvl
2	3	60	RL	68.0	11250	Pave	NaN	IR1	Lvl
3	4	70	RL	60.0	9550	Pave	NaN	IR1	Lvl
4	5	60	RL	84.0	14260	Pave	NaN	IR1	Lvl

5 rows x 81 columns



0s

```
[5] 1 df.shape
```

```
(2919, 81)
```



0s

```
[6] 1 df.columns
```

```
Index(['Id', 'MSSubClass', 'MSZoning', 'LotFrontage', 'LotArea', 'Street',  
      'Alley', 'LotShape', 'LandContour', 'Utilities', 'LotConfig',  
      'LandSlope', 'Neighborhood', 'Condition1', 'Condition2', 'BldgType',  
      'HouseStyle', 'OverallQual', 'OverallCond', 'YearBuilt', 'YearRemodAdd',  
      'RoofStyle', 'RoofMatl', 'Exterior1st', 'Exterior2nd', 'MasVnrType',  
      'MasVnrArea', 'ExterQual', 'ExterCond', 'Foundation', 'BsmtQual',  
      'BsmtCond', 'BsmtExposure', 'BsmtFinType1', 'BsmtFinSF1',  
      'BsmtFinType2', 'BsmtFinSF2', 'BsmtUnfSF', 'TotalBsmtSF', 'Heating',  
      'HeatingQC', 'CentralAir', 'Electrical', '1stFlrSF', '2ndFlrSF',  
      'LowQualFinSF', 'GrLivArea', 'BsmtFullBath', 'BsmtHalfBath', 'FullBath',  
      'HalfBath', 'BedroomAbvGr', 'KitchenAbvGr', 'KitchenQual',  
      'TotRmsAbvGrd', 'Functional', 'Fireplaces', 'FireplaceQu', 'GarageType',  
      'GarageYrBlt', 'GarageFinish', 'GarageCars', 'GarageArea', 'GarageQual',  
      'GarageCond', 'PavedDrive', 'WoodDeckSF', 'OpenPorchSF',  
      'EnclosedPorch', '3SsnPorch', 'ScreenPorch', 'PoolArea', 'PoolQC',  
      'Fence', 'MiscFeature', 'MiscVal', 'MoSold', 'YrSold', 'SaleType',  
      'SaleCondition', 'SalePrice'],  
      dtype='object')
```

Data Fields

- **SalePrice:** The property's sale price in dollars. This is the target variable that you're trying to predict.
- **MSSubClass:** The building class
- **MSZoning:** The general zoning classification
- **LotFrontage:** Linear feet of street connected to the property
- **LotArea:** Lot size in square feet
- **Street:** Type of road access
- **Alley:** Type of alley access
- **LotShape:** General shape of the property
- **LandContour:** Flatness of the property
- **Utilities:** Type of utilities available
- **LotConfig:** Lot configuration
- **LandSlope:** Slope of the property
- **Neighborhood:** Physical locations within Ames city limits
- **Condition1:** Proximity to the main road or railroad
- **Condition2:** Proximity to the main road or railroad (if a second is present)
- **BldgType:** Type of dwelling
- **HouseStyle:** Style of dwelling

✓ Handle missing values:

✓
0s



```
1 # Check for missing values
2
3 missing_values = df.isnull().sum()
4 missing_values
```



```
Id          0
MSSubClass  0
MSZoning    4
LotFrontage 486
LotArea      0
...
MoSold      0
YrSold      0
SaleType     1
SaleCondition 0
SalePrice   1459
Length: 81, dtype: int64
```

✓
0s



```
1 missing_values = missing_values[missing_values > 0]  
2 missing_values
```



MSZoning	4
LotFrontage	486
Alley	2721
Utilities	2
Exterior1st	1
Exterior2nd	1
MasVnrType	24
MasVnrArea	23
BsmtQual	81
BsmtCond	82
BsmtExposure	82
BsmtFinType1	79
BsmtFinSF1	1
BsmtFinType2	80

✓
0s

```
[9] 1 # Replace missing values in numeric columns with mean and in categorical columns with mode
    2
    3 numeric_cols = df.select_dtypes(include=np.number).columns.tolist()
    4 categorical_cols = df.select_dtypes(exclude=np.number).columns.tolist()
    5
    6 for col in numeric_cols:
    7     df[col].fillna(df[col].mean(), inplace=True)
    8
    9 for col in categorical_cols:
   10     df[col].fillna(df[col].mode()[0], inplace=True)
```

Goals Scored Over the Last 7 Games

1 3 4 6 6 7 8

mean
average **5**

mode
most common **6**

median
middle **6**

range
largest - smallest **7**

1 3 4 6 6 7 8

$$\text{mean} = \underset{\text{Total Sum}}{35} \div \underset{\text{Total Numbers}}{7} = \mathbf{5}$$

The mean is 5 goals per game.

1 3 4 **6 6** 7 8

$$\text{mode} = \mathbf{6}$$

The mode is 6 goals.

Look for numbers that repeat!

Mean

The **mean** is the average.

1. Add up all the values to find a total
2. Divide the total by the number of values you added together.

$$2+2+3+5+5+7+11 = 35$$

(There are 7 values)

$$35 \div 7 = 5$$

→ The **mean** is 5

Median

The **median** is the middle value.

1. Write the values in numerical order
2. The median is the middle value. If there are 2 values in the middle, find the average of these 2 numbers.

2, 2, 3, 5, 7, 11

→ The **median** is 5

Mode

The **mode** is the most frequent value.

1. Count how many of each value appears.
2. The mode is the value that appears the most often.

Note: You can have more than one mode.

2, 2, 3, 5, 5, 7, 11

→ The **modes** are 2 and 5

Range

The **range** is the difference between the lowest and highest value

1. Subtract the lowest value from the highest value in the data set

2, 2, 3, 5, 5, 7, 11

$$11 - 2 = 9$$

→ The **range** is 9

1, 3, 3, 3, 5, 6, 6, 9, 9, 9

There are two modes

3 9

✓
0s

```
[10] 1 # Reassessing for any remaining missing values  
      2  
      3 missing_values = df.isnull().sum()  
      4 missing_values
```

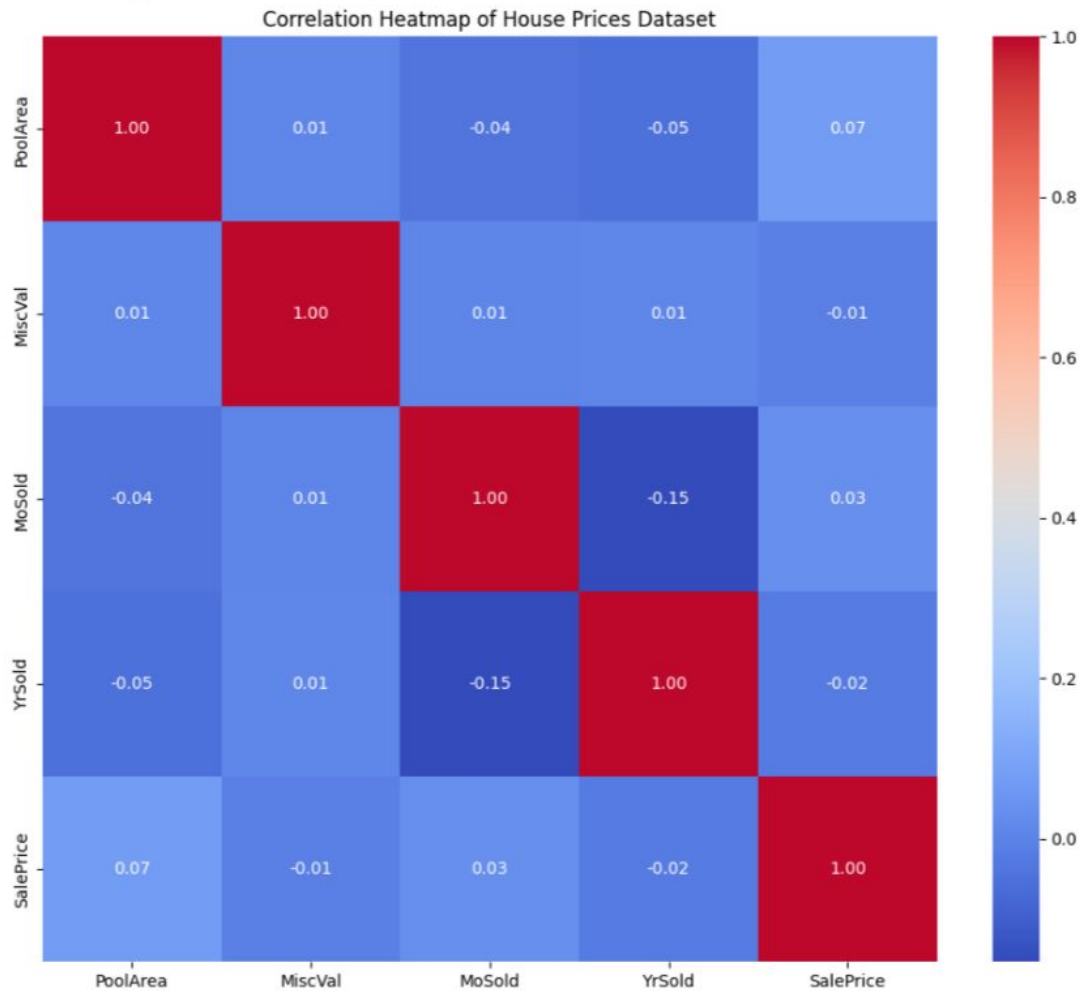
```
Id          0  
MSSubClass  0  
MSZoning    0  
LotFrontage 0  
LotArea     0  
           ..  
MoSold      0  
YrSold      0  
SaleType    0  
SaleCondition 0  
SalePrice   0  
Length: 81, dtype: int64
```

✓ Correlation

- A statistic that measures the degree to which two securities move in relation to each other.



```
1 # Calculate the correlations
2 correlation_matrix = df.iloc[:, -10:].corr()
3
4 # Set the figure size
5 plt.figure(figsize=(12, 10))
6
7 # Create a heatmap
8 sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt='.2f')
9 plt.title('Correlation Heatmap of House Prices Dataset')
10 plt.show()
```

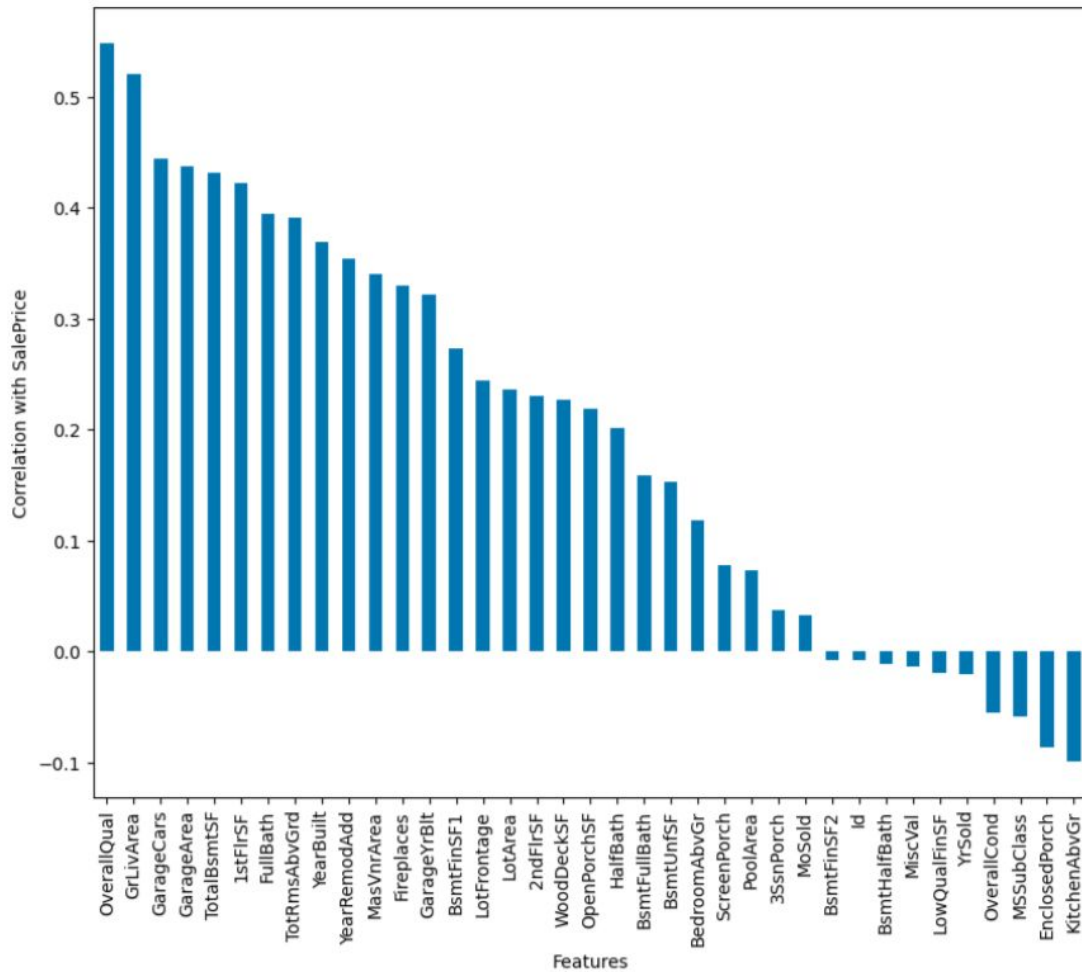


✓
0s



```
1 # Calculate correlations of features with the target 'SalePrice'
2 correlation_with_target = df.corr()['SalePrice'].sort_values(ascending=False)
3
4 # Plotting correlations with the target
5 plt.figure(figsize=(10, 8))
6 correlation_with_target.drop('SalePrice').plot(kind='bar')
7 plt.title('Correlation of Features with SalePrice')
8 plt.xlabel('Features')
9 plt.ylabel('Correlation with SalePrice')
10 plt.show()
```

Correlation of Features with SalePrice





```
1 # Calculate correlations of features with the target 'SalePrice'
2 correlation_with_target = df.corr()['SalePrice'].abs().sort_values(ascending=False)
3
4 # Select the top five features correlated with SalePrice
5 top_five_features = correlation_with_target[1:6] # Excluding SalePrice itself
6
7 # Display the top five correlated features
8 top_five_features
```

```
OverallQual    0.548617
GrLivArea      0.520311
GarageCars     0.444406
GarageArea     0.437654
TotalBsmtSF    0.431912
Name: SalePrice, dtype: float64
```

✓ Dummy Code (One-Hot Encoding)

✓
0s

```
[14] 1 df['SaleCondition'].unique()
```

```
array(['Normal', 'Abnorml', 'Partial', 'AdjLand', 'Alloca', 'Family'],  
      dtype=object)
```

✓
0s

```
[15] 1 # Assuming 'data' is your DataFrame containing both numerical and categorical columns  
    2 df = pd.get_dummies(df)
```

✓
0s



```
1 df.columns[-20:]
```

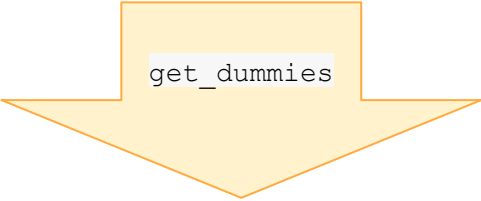


```
Index(['Fence_MnWw', 'MiscFeature_Gar2', 'MiscFeature_Othr',  
      'MiscFeature_Shed', 'MiscFeature_TenC', 'SaleType_COD', 'SaleType_CWD',  
      'SaleType_Con', 'SaleType_ConLD', 'SaleType_ConLI', 'SaleType_ConLw',  
      'SaleType_New', 'SaleType_Oth', 'SaleType_WD', 'SaleCondition_Abnorml',  
      'SaleCondition_AdjLand', 'SaleCondition_Alloca', 'SaleCondition_Family',  
      'SaleCondition_Normal', 'SaleCondition_Partial'],  
      dtype='object')
```


✓
0s

```
[14] 1 df.shape # before using get_dummies
```

```
(2919, 81)
```



get_dummies

✓
0s

```
[18] 1 df.shape # after using get_dummies
```

```
(2919, 290)
```

✓ Feature selection based on correlation:

Replace 'Feature1', 'Feature2', 'Feature3', 'Feature4', 'Feature5' with the names of the top 5 features you've identified from the correlation plot.

✓
0s

```
[14] 1 selected_features = ['Feature1', 'Feature2', 'Feature3', 'Feature4', 'Feature5']
```

✓
0s

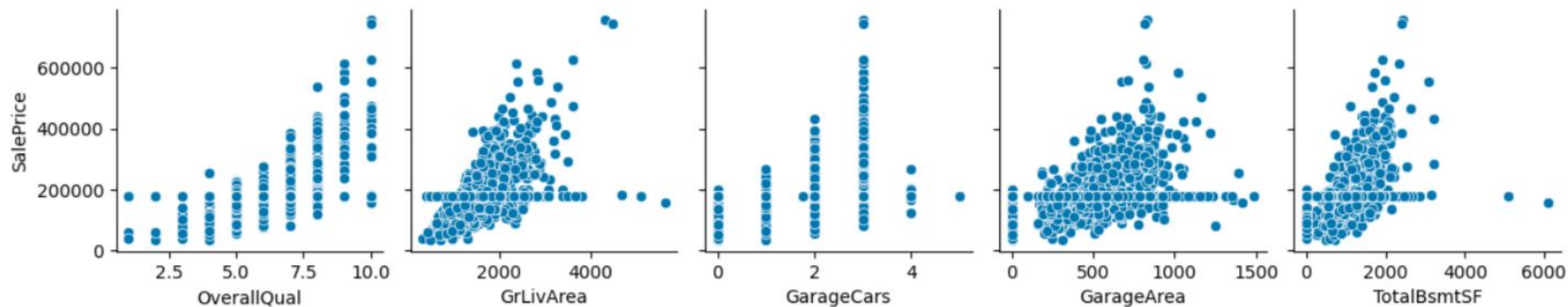
```
[15] 1 selected_features = ['OverallQual', 'GrLivArea', 'GarageCars', 'GarageArea', 'TotalBsmtSF']
```

✓
0s

```
[17] 1 X = df[selected_features]  
    2 y = df['SalePrice']
```

✓ 1s `1 sns.pairplot(data=df, y_vars=['SalePrice'], x_vars=selected_features)`

`<seaborn.axisgrid.PairGrid at 0x7f4640474250>`



- ✓ Split the dataset into training and testing sets:

[illegible]

✓ Create and fit a Linear Regression model:

✓
0s [19] 1 model = LinearRegression()

✓
0s [20] 1 model.fit(X_train, y_train)

▼ LinearRegression
LinearRegression()

✓
0s [21] 1 y_pred = model.predict(X_test)

✓ Calculate evaluation metrics: (RMSE, MAE, MSE)

✓
0s



```
1 mae = mean_absolute_error(y_test, y_pred)
2 mse = mean_squared_error(y_test, y_pred)
3 rmse = np.sqrt(mse)
4
5 # Print the calculated metrics
6 print(f"Mean Absolute Error (MAE): {mae}")
7 print(f"Mean Squared Error (MSE): {mse}")
8 print(f"Root Mean Squared Error (RMSE): {rmse}")
```



```
Mean Absolute Error (MAE): 31831.12330200309
Mean Squared Error (MSE): 1852270153.820561
Root Mean Squared Error (RMSE): 43038.00824643911
```

✓ Perform statistical analysis using p-values and obtain regression model summary:

✓
0s



```
1 X_train = sm.add_constant(X_train)
2 model_sm = sm.OLS(y_train, X_train).fit()
3 p_values = model_sm.pvalues
4 print("P-values:")
5 print(p_values)
6 print("Regression Model Summary:")
7 print(model_sm.summary())
```

```
P-values:
const      1.219841e-18
OverallQual 2.418788e-30
GrLivArea  1.007386e-38
GarageCars 2.729036e-01
GarageArea 4.382539e-02
TotalBsmtSF 1.102933e-10
dtype: float64
Regression Model Summary:
```

OLS Regression Results						
Dep. Variable:	SalePrice	R-squared:	0.400			
Model:	OLS	Adj. R-squared:	0.398			
Method:	Least Squares	F-statistic:	310.0			
Date:	Fri, 05 Jan 2024	Prob (F-statistic):	7.89e-255			
Time:	06:17:57	Log-Likelihood:	-28303.			
No. Observations:	2335	AIC:	5.662e+04			
Df Residuals:	2329	BIC:	5.665e+04			
Df Model:	5					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	3.672e+04	4131.783	8.888	0.000	2.86e+04	4.48e+04
OverallQual	1.092e+04	940.561	11.612	0.000	9077.150	1.28e+04
GrLivArea	30.2332	2.281	13.257	0.000	25.761	34.705
GarageCars	2996.5169	2732.404	1.097	0.273	-2361.590	8354.714
GarageArea	19.3333	9.586	2.017	0.044	0.536	38.131
TotalBsmtSF	16.9178	2.610	6.482	0.000	11.799	22.036
Omnibus:	1075.456	Durbin-Watson:	1.991			
Prob(Omnibus):	0.000	Jarque-Bera (JB):	16166.517			
Skew:	1.780	Prob(JB):	0.00			
Kurtosis:	15.389	Cond. No.	8.99e+03			

P-values:
const 1.219841e-18
OverallQual 2.418788e-30
GrLivArea 1.007386e-38
GarageCars 2.729036e-01
GarageArea 4.382539e-02
TotalBsmtSF 1.102933e-10

dtype: float64

Regression Model Summary:

OLS Regression Results

```
=====
Dep. Variable:      SalePrice      R-squared:      0.400
Model:              OLS           Adj. R-squared:    0.398
Method:             Least Squares  F-statistic:    310.0
Date:               Fri, 05 Jan 2024 Prob (F-statistic): 7.89e-255
Time:               06:17:57      Log-Likelihood: -28303.
No. Observations:   2335          AIC:             5.662e+04
Df Residuals:       2329          BIC:             5.665e+04
Df Model:           5
Covariance Type:    nonrobust
=====
```

	coef	std err	t	P> t	[0.025	0.975]
const	3.672e+04	4131.783	8.888	0.000	2.86e+04	4.48e+04
OverallQual	1.092e+04	940.561	11.612	0.000	9077.150	1.28e+04
GrLivArea	30.2332	2.281	13.257	0.000	25.761	34.705
GarageCars	2996.5169	2732.404	1.097	0.273	-2361.680	8354.714
GarageArea	19.3333	9.586	2.017	0.044	0.536	38.131
TotalBsmtSF	16.9178	2.610	6.482	0.000	11.799	22.036

```
=====
Omnibus:            1075.456      Durbin-Watson:      1.991
Prob(Omnibus):      0.000        Jarque-Bera (JB):    16166.517
Skew:               1.780        Prob(JB):           0.00
Kurtosis:           15.389       Cond. No.            8.99e+03
=====
```

- Explain the regression equation

Week 7: Assignment

Apply **regression modelling** techniques to **predict house prices** accurately.

The primary objective of this assignment is to create regression models using the House Prices dataset to predict property sale prices accurately.

