

## Correlation

Correlation in machine learning refers to the statistical relationship between two or more variables. It measures how changes in one variable are associated with changes in another variable. Understanding correlation is crucial in various machine learning tasks, including feature selection, data preprocessing, and model evaluation.

Here are different types of correlation coefficients commonly used in machine learning:

1. **Pearson Correlation Coefficient:** This measures the linear relationship between two continuous variables. It ranges from -1 to 1, where -1 indicates a perfect negative correlation, 1 indicates a perfect positive correlation, and 0 indicates no correlation.
2. **Spearman's Rank Correlation Coefficient:** Unlike Pearson correlation, Spearman's correlation measures the strength and direction of association between two variables, but it does not assume that the variables are normally distributed. Instead, it considers the ranks of the data points. It is suitable for both continuous and ordinal variables.
3. **Kendall's Tau:** Similar to Spearman's correlation, Kendall's Tau also measures the association between two variables based on the ranks of the data. It is particularly useful when dealing with small sample sizes.

Correlation is often used in feature selection to identify redundant or highly correlated features, which can lead to overfitting and increased computational complexity. In such cases, one of the correlated features may be removed to improve model performance and interpretability.

Moreover, correlation analysis helps in understanding the relationship between input variables and the target variable. Features with high correlation to the target variable are likely to be more predictive and can improve model accuracy.

However, it's essential to note that correlation does not imply causation. Even if two variables are highly correlated, it does not necessarily mean that changes in one variable cause changes in the other. Causation requires additional evidence and careful experimentation.

In summary, correlation analysis is a fundamental tool in machine learning for understanding relationships between variables, identifying redundant features, and improving model performance.

## Dataset:

**we create the dataset with python code:**

```
import numpy as np
import pandas as pd

# Set random seed for reproducibility
np.random.seed(42)

# Number of data points
num_entries = 300

# Generate apartment size in meters squared
apartment_size = np.random.randint(50, 150, size=num_entries)

# Generate apartment location
cities = ['Tbilisi', 'Batumi', 'Kutaisi', 'Rustavi', 'Gori']
location = np.random.choice(cities, size=num_entries)

# Generate apartment prices based on size and location
# Tbilisi prices are higher, and there's a strong correlation between
# size and price
price = 1000 * apartment_size + 5000 * (location == 'Tbilisi') +
np.random.normal(0, 10000, size=num_entries)

# Create DataFrame
data = pd.DataFrame({'Apartment Size (m²)': apartment_size,
                    'Location': location,
                    'Price (GEL)': price})

# Save to CSV
data.to_csv('apartment_data.csv', index=False)

print("Dataset saved to apartment_data.csv")
```

**Some rows from the resulting dataset:**

Apartment Size (m <sup>2</sup> )	Price (GEL)	Location
101	96666.9659309 744	Kutaisi
142	144038.227275 27466	Rustavi
64	74591.0920261 4527	Tbilisi
121	129725.415914 30479	Rustavi
110	121543.543384 08924	Kutaisi
70	70226.1826942 4679	Tbilisi
132	125357.725566 95578	Rustavi
136	149601.113662 48172	Rustavi
124	103629.240650 2786	Kutaisi
124	144700.322805 42672	Tbilisi
137	133976.757579 23693	Kutaisi
149	153954.371710 09704	Tbilisi
73	66176.5739887 9847	Gori
52	51518.0451446 41874	Batumi
71	83711.3822598 0915	Batumi

## Code demonstration:

```
import pandas as pd

# Read the CSV file
data = pd.read_csv('apartment_data.csv')

# Convert location column to one-hot encoding
location_dummies = pd.get_dummies(data['Location'])

# Concatenate one-hot encoded location columns with the original data
data_encoded = pd.concat([data, location_dummies], axis=1)

# Drop the original location column
data_encoded.drop(columns=['Location'], inplace=True)

# Compute the correlation matrix
correlation_matrix = data_encoded.corr()

# Print correlation matrix
print("Correlation Matrix:")
print(correlation_matrix)
```

Correlation Matrix:

	Apartment Size (m²)	Price (GEL)	Batumi	Gori	Kutaisi	Rustavi	Tbilisi
Apartment Size (m²)	1.000000	0.942732	-0.017101	-0.055617	0.011068	0.005104	0.057761
Price (GEL)	0.942732	1.000000	-0.052138	-0.090927	0.033855	-0.015366	0.127548
Batumi	-0.017101	-0.052138	1.000000	-0.260378	-0.226285	-0.257790	-0.255198
Gori	-0.055617	-0.090927	-0.260378	1.000000	-0.235678	-0.268491	-0.265792
Kutaisi	0.011068	0.033855	-0.226285	-0.235678	1.000000	-0.233336	-0.230990
Rustavi	0.005104	-0.015366	-0.257790	-0.268491	-0.233336	1.000000	-0.263150
Tbilisi	0.057761	0.127548	-0.255198	-0.265792	-0.230990	-0.263150	1.000000