# Correlation

Correlation serves as a statistical tool to gauge the strength and direction of the linear relationship between two variables. It quantifies the tendency of two variables to change together. The correlation coefficient (often represented as 'r') spans from -1 to 1, with interpretations as follows:

$$r = 1$$

r=1 denotes a perfect positive correlation, implying that as one variable increases, the other variable also increases proportionally.

$$r = -1$$

r=−1 signifies a perfect negative correlation, indicating that as one variable increases, the other variable decreases proportionally.

$$r = 0$$

r=0 suggests no linear correlation between the variables.

It's crucial to note that correlation doesn't necessarily imply causation; rather, it merely measures the degree of association between variables. In machine learning, correlation analysis finds applications in feature selection, understanding variable relationships, detecting multicollinearity, and uncovering patterns or anomalies in data.

## Cybersecurity Application: Detecting Coordinated Attacks

Correlation finds significant utility in cybersecurity for detecting coordinated attacks or spotting patterns within network traffic, system logs, or other security-related datasets. For instance, in a scenario where we possess data regarding various network events (e.g., failed login attempts, port scans, malware detections) across multiple systems or endpoints, we can compute the correlation between these events to unveil potential coordinated attacks or compromised systems.

Sample Data:

Let's consider a dataset containing the following attributes for numerous network events across diverse systems:

System ID

Event Type (e.g., failed login, port scan, malware detection)

Timestamp

Below is a snippet to generate a small random sample dataset using Python:

```python
import numpy as np

import pandas as pd

# Sample data - 20 events across 5 systems

data = {

    'system_id': np.random.randint(1, 6, size=20),

    'event_type': np.random.randint(1, 4, size=20),  # 1: failed login, 2: port scan, 3: malware

    'timestamp': pd.date_range(start='2023-05-01', end='2023-05-05', periods=20)

}

df = pd.DataFrame(data)
```

Python Code for Correlation:

Here's the Python code snippet to compute correlations:

```python
import pandas as pd

import numpy as np

from scipy.stats import pearsonr


# Calculate correlation between event types

corr_matrix = df['event_type'].corr(df['event_type'])

print("Correlation Matrix:\n", corr_matrix)

# Calculate correlation between event types and system IDs

for event_type in df['event_type'].unique():
```

```python
    subset = df[df['event_type'] == event_type]

    system_ids = subset['system_id'].values

    timestamps = subset['timestamp'].values

    r, p = pearsonr(system_ids, timestamps)

    print(f"Correlation between event type {event_type} and system IDs: {r:.2f}")
```

This code snippet computes the correlation between different event types to identify potential coordinated attacks. Additionally, it calculates the correlation between event types and system IDs to detect if certain systems are more susceptible to specific event types.