# Correlation

Correlation is a statistical measure that describes the strength and direction of the linear relationship between two variables. It quantifies the extent to which two variables tend to change together. The correlation coefficient, usually denoted as r, ranges from -1 to 1, where:

$r = 1$ indicates a perfect positive correlation (as one variable increases, the other variable also increases)
$r = -1$ indicates a perfect negative correlation (as one variable increases, the other variable decreases)
$r = 0$ indicates no linear correlation between the variables

Correlation does not necessarily imply causation; it only measures the degree of association between variables.
In machine learning, correlation analysis is often used for feature selection, understanding the relationships between variables, detecting multicollinearity (high correlation between predictors), and identifying potential patterns or anomalies in data.

## Cybersecurity Application: Detecting Coordinated Attacks

In cybersecurity, correlation can be used to detect coordinated attacks or identify patterns in network traffic, system logs, or other security-related data sources.
For instance, if we have data on various network events (e.g., failed login attempts, port scans, malware detections) across multiple systems or endpoints, we can calculate the correlation between these events to identify potential coordinated attacks or compromised systems.

Sample Data:
Assume we have a dataset containing the following features for various network events across multiple systems:

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

We can generate a small random sample dataset:

```
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:

```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():
    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 calculates the correlation between different event types (e.g., failed logins, port scans, and malware detections) to identify potential coordinated attacks. Additionally, it computes the correlation between event types and system IDs to detect if certain systems are more prone to specific types of events.