# SNS COLLEGE OF ENGINEERING

Kurumbapalayam (Po), Coimbatore – 641 107

**An Autonomous Institution**

Accredited by NBA – AICTE and Accredited by NAAC – UGC with 'A' Grade
Approved by AICTE, New Delhi & Affiliated to Anna University, Chennai

**DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING-IoT Including CS & BCT**

COURSE NAME :19SB701  PATTERN RECOGNITION TECHNIQUES IN CYBER CRIME

IV YEAR / VII SEMESTER

Unit IV- **MALWARE ANALYSIS AND NETWORK TRAFFIC ANALYSIS**

Topic  :Feature Engineering – Detection with Data and Algorithms

31-08-2024

Feature Engineering – Detection with Data and Algorithms   / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE

1

Feature engineering is a crucial step in the data science pipeline that involves transforming raw data into a format that is better suited for machine learning models.

Effective feature engineering can significantly enhance the performance of algorithms used for anomaly detection, classification, and other predictive tasks.

It involves creating, selecting, and refining features that capture important patterns and relationships in the data.

**Feature Engineering – Detection with Data and Algorithms   / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE**

**Key Aspects of Feature Engineering for Detection:**

**1. Understanding the Data:**

    **1. Exploration:** Start by exploring the dataset to understand its structure, types of variables, distributions, and potential relationships.

    **2. Domain Knowledge:** Use domain-specific knowledge to guide the creation of features that are relevant to the detection task.

31-08-2024

Feature Engineering – Detection with Data and Algorithms   / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE

3

**Feature Creation:**

•**Statistical Features:** Compute statistical measures such as mean, variance, skewness, and kurtosis for each feature. These can help capture underlying patterns in the data.

- Example: In time-series data, features like moving averages or rolling statistics can be useful.

•**Temporal Features:** For time-series data, create features like lag values, differences, or trends.

- Example: In anomaly detection in sensor data, include features like the difference between consecutive measurements.

•**Aggregated Features:** Aggregate data over time or different segments to capture overall patterns.

- Example: Sum, average, or count of events over specific intervals.

•**Interaction Features:** Combine existing features to capture interactions or relationships.

- Example: In network traffic data, combine source and destination IP addresses to create a feature representing communication patterns.

**Feature Transformation:**

•**Normalization/Scaling:** Scale features to a common range to ensure that all features contribute equally to the model.

- Example: Standardize features to have zero mean and unit variance, or normalize to a [0,1] range.

•**Encoding Categorical Variables:** Convert categorical variables into numerical format using techniques like one-hot encoding or label encoding.

- Example: Encode categorical features such as user roles or product categories.

•**Dimensionality Reduction:** Use techniques like Principal Component Analysis (PCA) or t-Distributed Stochastic Neighbor Embedding (t-SNE) to reduce the number of features while preserving important information.

- Example: Reduce high-dimensional text data into a smaller set of principal components.

31-08-2024

Feature Engineering – Detection with Data and Algorithms   / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE

5

**Feature Selection:**

•**Filter Methods:** Use statistical tests to select features that have strong relationships with the target variable.

- Example: Use chi-square tests for categorical features or correlation coefficients for numerical features.

•**Wrapper Methods:** Evaluate the performance of different subsets of features by training models and selecting the best-performing subset.

- Example: Recursive Feature Elimination (RFE) for feature selection.

•**Embedded Methods:** Perform feature selection as part of the model training process.

- Example: Lasso (L1 regularization) in regression models, which penalizes less important features.

**Feature Engineering for Anomaly Detection:**

•**Threshold-Based Features:** Create features that capture deviations from known thresholds or expected ranges.

- Example: In fraud detection, create features that indicate whether transaction amounts exceed typical thresholds.

•**Seasonal Decomposition:** Decompose time-series data into seasonal, trend, and residual components to better capture anomalies.

- Example: Detect anomalies in sales data by analyzing deviations from seasonal patterns.

•**Distance-Based Features:** Compute distances between data points to identify anomalies.

- Example: In clustering-based anomaly detection, compute distances from cluster centroids.

**Examples of Feature Engineering in Detection:**

**1.Fraud Detection in Financial Transactions:**

    **1.Raw Data:** Transaction amount, time, location, merchant type.

    **2.Engineered Features:** Transaction frequency, average transaction amount per day, deviation from average spending, time since last transaction, merchant category.

**2.Network Intrusion Detection:**

    **1.Raw Data:** Source IP, destination IP, port numbers, protocol types, packet sizes.

    **2.Engineered Features:** Number of connections per IP, connection duration, number of unique ports used, traffic volume per IP.

Feature Engineering – Detection with Data and Algorithms / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE

## 3.Predictive Maintenance:

1. **Raw Data:** Sensor readings (temperature, pressure, vibration), time stamps.
2. **Engineered Features:** Moving averages of sensor readings, rate of change in readings, frequency of readings exceeding thresholds, time since last maintenance.

## Tools and Techniques:

- **Pandas:** For data manipulation and feature creation in Python.
- **Scikit-Learn:** Provides tools for feature selection and transformation, including standardization and encoding.
- **Featuretools:** An open-source Python library for automated feature engineering.
- **TensorFlow and PyTorch:** For deep learning-based feature engineering, especially when dealing with complex data types like images or sequences.

31-08-2024

Feature Engineering – Detection with Data and Algorithms  / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE

9

**MCQ**

1.What is the primary purpose of feature engineering in anomaly detection?

A) To reduce the size of the dataset.

B) To transform raw data into a format that better captures patterns and improves model performance.

C) To increase the complexity of the data for better visualization.

D) To create a large number of irrelevant features.

**Answer: B**

2.Which of the following is an example of a feature transformation technique?

A) Normalization

B) Feature Selection

C) Feature Creation

D) Data Cleaning

**Answer: A**

31-08-2024

Feature Engineering – Detection with Data and Algorithms   / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE

11

3. In the context of anomaly detection, what does the term "threshold-based feature" refer to?

A) A feature that captures how often a data point exceeds a predefined threshold.

B) A feature that calculates the mean value of a dataset.

C) A feature that represents the average distance between data points.

D) A feature that encodes categorical variables into numerical values.

**Answer: A**

31-08-2024

Feature Engineering – Detection with Data and Algorithms   / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE

12

Any Query????

Thank you……

Feature Engineering – Detection with Data and Algorithms   / 23ITB201-DATA STRUCTURES & ALGORITHMS /Mr.R.Kamalakkannan/CSE-IOT/SNSCE