Ahmet Emre Usta                                                         04.12.2023
2200765036

## AIN427 Midterm

## Question 1: Data Preprocessing

**a)**      In the realm of data science and machine learning, the effective preprocessing of data is a critical step that significantly influences the success of model outcomes. Real-world data often have many imperfections, including errors, inconsistencies, and missing values. Data cleaning, a vital component of preprocessing, is employed to address these issues. It involves techniques like imputation for missing values, outlier detection and handling, and correcting inaccuracies. Without this cleansing process, models may yield biased or invalid results due to flawed input data. This includes data transformation techniques, such as normalization and one-hot encoding, to standardize and reformat data for better algorithm compatibility, and data reduction methods, such as PCA, binning, and sampling, to manage large datasets efficiently without losing critical information. Addressing imbalanced datasets, particularly in classification problems, is also essential, employing techniques like SMOTE and under sampling to prevent model bias. Feature engineering, driven by domain knowledge, plays a vital role in creating or modifying features to capture essential patterns in data, aiding in more accurate predictions. Overall, the success of machine learning models heavily relies on these preprocessing steps, ensuring high-quality, relevant, and well-structured input data.

**b)**

**Normalization:** It scales numerical features to a uniform range, often between 0 and 1, which helps in stabilizing the training process and improving convergence in machine learning algorithms.

**Example:** If a dataset contains a feature like age ranging from 18 to 90 and another feature like income ranging from $20,000 to $100,000, normalization will scale these features into a 0-1 range so that no single feature dominates the model due to its scale.

**Handling Missing Values:** Addresses gaps in the dataset to avoid errors during analysis. Missing data can be due to various reasons, like errors in data collection or non-response.

**Example:** If a dataset with patient records has missing values in the blood pressure column, you might fill these gaps using the average or median blood pressure of the dataset or use more sophisticated imputation methods like predictive modeling.

**Encoding Categorical Variables:** Machine learning models usually require numerical input, so categorical data must be converted into a numerical format.

**Example:** For a dataset with a categorical feature like color (with values like 'Red', 'Blue', 'Green'), one-hot encoding would create three new binary features ('Color_Red', 'Color_Blue', 'Color_Green'), each representing the presence or absence of a color in each data point.

**Question 2: Association Rule Mining**

**a)**      Association rule mining, a key technique in data analysis, involves a two-step process that unveils hidden patterns within the datasets. The first step is to identify all frequent item sets that meet a specified minimum support threshold, which helps recognize the most common item combinations. The second step builds upon this by generating strong association rules from these item sets based on minimum confidence thresholds. These rules reveal significant relationships and dependencies between the items. Beyond its traditional use in market basket analysis, association rule mining has diverse applications across various fields. Bioinformatics aids in the identification of gene patterns, thereby contributing to genetic research and medical advancements. E-commerce plays a pivotal role in generating product recommendations and enhancing customer experience and business strategies. Additionally, this technique is increasingly employed in the analysis of sensor data in IoT applications, where it helps interpret vast streams of sensor data, leading to improved decision-making and system efficiency.

**b)**      The ECLAT algorithm, which is a variation of association rule mining, offers improved efficiency through its unique approach. By employing a Depth-First Search methodology, it stands apart from the traditional Apriori algorithm. In ECLAT, exploration of the dataset is conducted depth-first, focusing on extending an itemset only with items that are frequently associated with it. This selective extension results in a notably faster performance, as it avoids the exhaustive pair-wise comparison common in Apriori. Complementing this is the ECLAT's use of a vertical data

format. This format is particularly efficient for counting the support of item sets, a key step in association rule mining. By organizing data such that each item points to its transaction IDs, ECLAT can swiftly calculate itemset support through the intersection of these ID sets. This method is significantly more efficient than repetitive database scanning required in other approaches, further enhancing the speed and performance of the ECLAT algorithm.

**c)**     The Direct Hashing and Pruning (DHP) technique streamlines the process of identifying frequent item sets within large datasets. By employing a hash-based counting mechanism, DHP efficiently tallies pairs of items, thereby circumventing the need to retain and repeatedly access the entire dataset, a task that can be resource-intensive. Concurrently, the algorithm applies a dynamic pruning strategy. This strategy is grounded in the principle that all non-empty subsets of a frequent item set must be frequent. By eliminating itemset candidates that do not meet this criterion, the algorithm effectively reduces the search space, focuses on the most promising candidates, and enhances the overall efficiency. This synergistic use of hashing and pruning not only expedites the mining process but also minimizes memory usage, making it an adept technique for handling extensive datasets in association rule mining.

**d)**     https://www.kaggle.com/datasets/heeraldedhia/groceries-dataset

| # Member_number | A Date | A itemDescription |
|---|---|---|
| ID of customer | Date of purchase | Description of product purchased |
| 1000  —  5000 | **728** unique values | whole milk 6% / other vegetables 5% / Other (34365) 89% |
| 1808 | 21-07-2015 | tropical fruit |
| 2552 | 05-01-2015 | whole milk |
| 2300 | 19-09-2015 | pip fruit |
| 1187 | 12-12-2015 | other vegetables |
| 3037 | 01-02-2015 | whole milk |
| 4941 | 14-02-2015 | rolls/buns |
| 4501 | 08-05-2015 | other vegetables |
| 3803 | 23-12-2015 | pot plants |
| 2762 | 20-03-2015 | whole milk |

**Complementary Products:** Customers who buy certain items, such as coffee, might also tend to buy related items, such as sugar or cream.

*Rule Example:* Coffee -> Sugar

**Frequent Itemset:**  Identifying sets of items that are commonly purchased together.

*Rule Example:* Bread, Milk -> Eggs suggests that when customers buy bread and milk, they often buy eggs.

**Seasonal Purchases:** Discovering patterns in purchases related to seasons or holidays, such as increased sales of certain products during Christmas or the summer.

*Rule Example:* Christmas Season -> Mince Pies, Mulled Wine

**Promotional Strategies:** Analyzing the effect of promotions and discounts on the sales of other non-discounted products.

*Rule Example:* Discounted Chips -> Regular Priced Soft Drinks

**Cross-category Purchases:** Understanding cross-category purchase behavior, such as customers buying hygiene products when buying baby food.

*Rule Example:* Baby Food -> Diapers

**Time-Based Patterns:** Consider time-based patterns, such as the tendency to buy certain items on weekends or during specific hours of the day.

*Rule Example:* Saturday -> Snack Items

**Customer Segmentation:** Determine product affinities within customer segments based on factors such as the frequency of visits or total spending.

*Rule Example:* High Frequency Member -> Premium Products