



Department of Computer Science

CS571 – Data Preparation and Analysis

Instructor: Oleksandr Narykov

## **CHICAGO WEATHER PREDICTION USING MACHINE LEARNING TECHNIQUE**

Submitted by

Devendra singh portey, A20538603

Sivakrishna Golla, A20563027

Haritha chennuru venugopal reddy, A20531012

Satya mani srujan dommeti, A20594429

Sneha poojitha gade, A20558333

# Contents

<b>1 Introduction</b>	<b>4</b>
1.1 Abstract	4
1.1 Problem Statement	5
<b>2 Data Analysis</b>	<b>5</b>
2.1 Exploratory Data Analysis	5
2.1.1 File Description	5
2.1.2 Data Visualization	8
2.1.3 Insights from Visualization	11
<b>3 Unsupervised Learning Techniques</b>	<b>13</b>
3.1 K-Means Clustering	13
3.2 Hierarchical Clustering	16
3.3 DBSCAN Clustering	17
3.4 Anomaly detection with Isolation Forest	19
3.5 Anomaly detection with One-class SVM	21
3.6 Association rule with Apriori Algorithm	24
3.7 Feature Transformation with Auto Encoders	24
3.8 Feature Transformation with Kernel PCA	25
3.9 Feature Transformation with Kernel PCA	26
<b>4 Cross Validation Strategies</b>	<b>27</b>
4.1 Stability-based cross-validation for dimensionality reduction	27
4.2 Clustering Stability Validation	28
4.3 Anomaly Stability with Isolation Forest	28
4.4 Cross-Validation of Association Rule Learning	29
<b>5 Identify pitfalls and Performance Measure</b>	<b>29</b>
<b>6 Performance Enhancement Techniques</b>	<b>30</b>
6.1 Experiment 1	31
6.2 Experiment 2	32
6.3 Experiment 3	33

## List of Figures

Figures	Page No
Fig 1.1 Correlation between the features	8
Fig 1.2 Principal Component Analysis	9
Fig 1.3 UMAP (Uniform Manifold Approximation and Projection)	10
Fig 1.4 t-SNE (t-Distributed Stochastic Neighbor Embedding)	11
Fig 1.5 Elbow Method for optimal k	14
Fig 1.6 K-Means Clustering	15
Fig 1.7 Hierarchical Clustering	17
Fig 1.8 DBSCAN Clustering	19
Fig 1.9 Advanced Dimensionality Reduction with UMAP	21
Fig 1.10 Anomaly detection with Isolation Forest	22
Fig 1.11 Anomaly detection with one-class SVM	23
Fig 1.12 Feature Transformation with Auto Encoders	25
Fig 1.13 Feature Transformation with Kernel PCA	26
Fig 1.14 Stability-based cross-validation for dimensionality reduction	27

# 1 Introduction

## 1.1 Abstract

Weather prediction has been a fundamental pursuit for centuries, given its profound impact on human activities and the environment. Accurate weather forecasts are essential for various sectors, including agriculture, transportation, energy, disaster management, and public health. The ability to predict weather conditions not only helps in planning day-to-day activities but also plays a critical role in mitigating the adverse effects of extreme weather events like storms, floods, and heatwaves. In recent years, advancements in data science, artificial intelligence, and machine learning have revolutionized weather forecasting by leveraging large-scale historical data and powerful computational tools.

Traditional weather prediction methods rely heavily on physical and mathematical models of atmospheric dynamics, which, while effective, often struggle to handle the complexity and non-linearity of weather systems. Machine learning techniques offer a complementary approach, focusing on data-driven models to uncover hidden patterns and relationships among weather variables. By analyzing historical weather data, these models can generate forecasts with improved accuracy, often in shorter timeframes than conventional methods.

The importance of weather prediction extends beyond immediate applications. In agriculture, accurate forecasts enable farmers to plan planting, irrigation, and harvesting schedules, minimizing crop losses and maximizing yield. In transportation, weather predictions improve safety and efficiency by guiding airline schedules, maritime navigation, and road traffic management. The energy sector also benefits greatly, as renewable energy sources like solar and wind depend heavily on weather conditions. Moreover, early warnings of severe weather events can save lives and reduce economic losses by allowing timely evacuation and preparedness measures.

This project seeks to address the growing need for accurate weather forecasting by analyzing an extensive dataset comprising hourly weather data from 2021 to 2024. By applying advanced machine learning techniques, the aim is to identify relationships between variables such as temperature, humidity, precipitation, wind speed, and pressure, and use these insights to predict future weather conditions effectively. This approach not only aims to

enhance predictive accuracy but also to provide actionable insights that can benefit a wide range of stakeholders in weather-sensitive domains.

## 1.2 Problem Statement

In this project, we aim to predict the weather conditions in Chicago by analyzing historical weather data and identifying patterns that can help forecast key parameters such as temperature, precipitation, humidity, and wind speed. The goal is to develop a machine learning model capable of predicting weather conditions accurately, providing valuable insights into future weather patterns. Accurate weather forecasting is crucial for various industries, including transportation, agriculture, and urban planning, to make informed decisions based on future climate scenarios.

We will be working with a weather dataset containing hourly weather observations for Chicago, including features like temperature, humidity, wind speed, pressure, and precipitation. The model will use this data to predict weather conditions for future time periods. The steps involved in this project are as follows:

1. **Explore the Dataset:** Analyze the dataset to identify correlations between features and the target variable (e.g., temperature, precipitation). Assess the impact of each feature on the target and evaluate independence assumptions to ensure valid modeling.
2. **Visualization:** Apply various visualization strategies such as correlation plots, dimensionality reduction techniques (PCA, UMAP, t-SNE) to identify underlying patterns and relationships between variables. Insights gained from these visualizations will guide feature selection and model development.
3. **Unsupervised Learning:** Explore the data using unsupervised learning techniques to identify clusters and patterns that may not be immediately obvious in the data.
4. **Cross-Validation Strategy:** Choose an appropriate cross-validation strategy to ensure the model generalizes well to unseen data, minimizing overfitting and underfitting.
5. **Initial Model Training:** Train a simple machine learning model to predict weather conditions and evaluate its performance. Use a validation set for hyperparameter

tuning and/or early stopping to optimize the model. Analyze its performance using cross-validation to identify potential pitfalls and limitations.

6. **Model Improvement:** Propose and implement strategies to improve the model's performance, such as feature selection, regularization, and increasing model complexity. Conduct at least two more experiments to test the impact of these changes.

The challenge in this project lies in effectively handling the temporal and spatial complexity of weather data to make accurate short-term forecasts. By leveraging machine learning techniques and advanced visualization methods, the goal is to improve the accuracy and reliability of weather predictions for Chicago, offering actionable insights that can help various industries plan for future weather events.

## 2 Data Analysis

In this study we followed the below steps sequentially to design an model for predicting the performance of players.

### 2.1 Exploratory Data Analysis

#### 2.1.1 File Description

The dataset contains historical weather observations for Chicago and includes key features related to weather conditions, such as temperature, humidity, wind speed, pressure, and precipitation. The data is structured in a tabular format, with each row representing an hourly weather observation. Dataset contains hourly weather observations for Chicago with the following features:

1. **Date/Time:** Timestamp of the observation.
2. **Temperature (°C):** Air temperature in degrees Celsius.
3. **Dew Point (°C):** Temperature at which air becomes saturated with moisture.
4. **Humidity (%):** Relative humidity percentage.
5. **Pressure (hPa):** Atmospheric pressure in hectopascals.
6. **Wind Speed (km/h):** Wind speed in kilometers per hour.
7. **Wind Direction (°):** Wind direction in degrees (0° = north).
8. **Precipitation (mm):** Amount of precipitation in millimeters.
9. **Visibility (km):** Distance at which objects can be seen in kilometers.
10. **Cloud Cover (%):** Percentage of sky covered by clouds.
11. **Condition:** Weather condition (e.g., clear, cloudy, rain, snow).

The observations cover multiple weather parameters, which can be used to analyze correlations and predict future weather conditions. Missing or invalid data points may be handled through imputation techniques or removed based on the analysis requirements. The target variable could be any of the weather parameters, such as temperature or precipitation,

depending on the modelling goals. This dataset is valuable for forecasting weather conditions and identifying patterns that can influence decision-making in sectors like transportation, agriculture, and urban planning.

2.1.2 Data Visualization

Fig 1.1 below shows the Correlation between the features visualization helps to identify relationships between different weather features in the dataset. Strong correlations (either positive or negative) between features can provide insights into how they interact, while weak or no correlations may suggest that features are independent of one another. It can also aid in identifying potential multicollinearity issues when building predictive models.

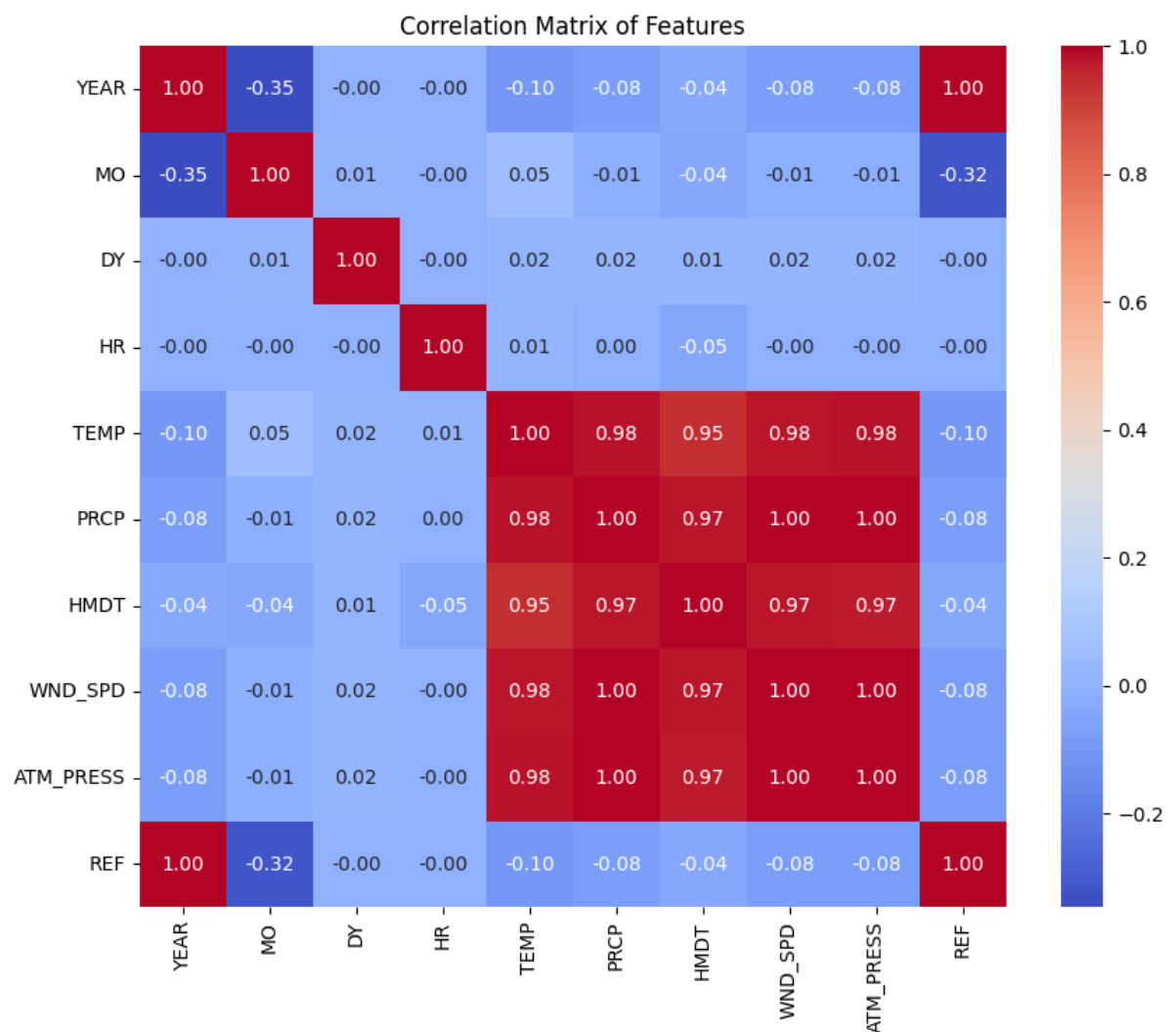


Fig 1.1 Correlation between the features



Fig 1.2 below shows the Principal Component Analysis(PCA) Visualization to reduce the dataset's dimensionality to two principal components. It selects numeric features, applies PCA, and transforms the data into two components (PC1 and PC2). A scatter plot is then generated to visualize the data in the 2D space of the principal components. This helps identify patterns or clusters in the data and provides insights into the most important features driving the variance in the dataset.

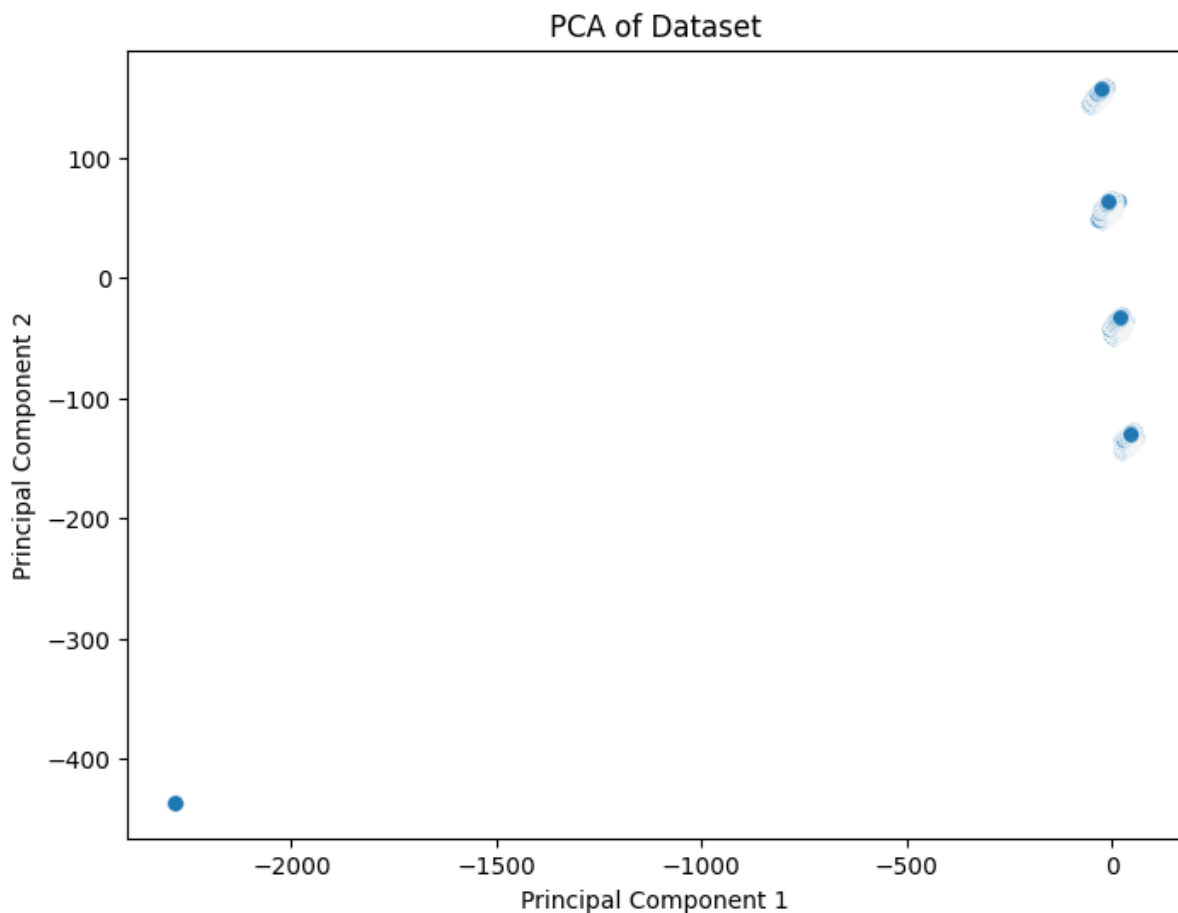


Fig 1.2 Principal Component Analysis

Fig 1.3 below shows the UMAP (Uniform Manifold Approximation and Projection). The code Used applies UMAP (Uniform Manifold Approximation and Projection) to reduce the dataset's dimensionality to two components. It selects the numeric features from the dataset, applies UMAP, and transforms the data into two dimensions (UMAP1 and UMAP2). A scatter plot is then created to visualize the transformed data in a 2D space. This visualization helps to explore

the relationships and structure of the dataset, identifying clusters or patterns that may not be apparent in the higher-dimensional space.

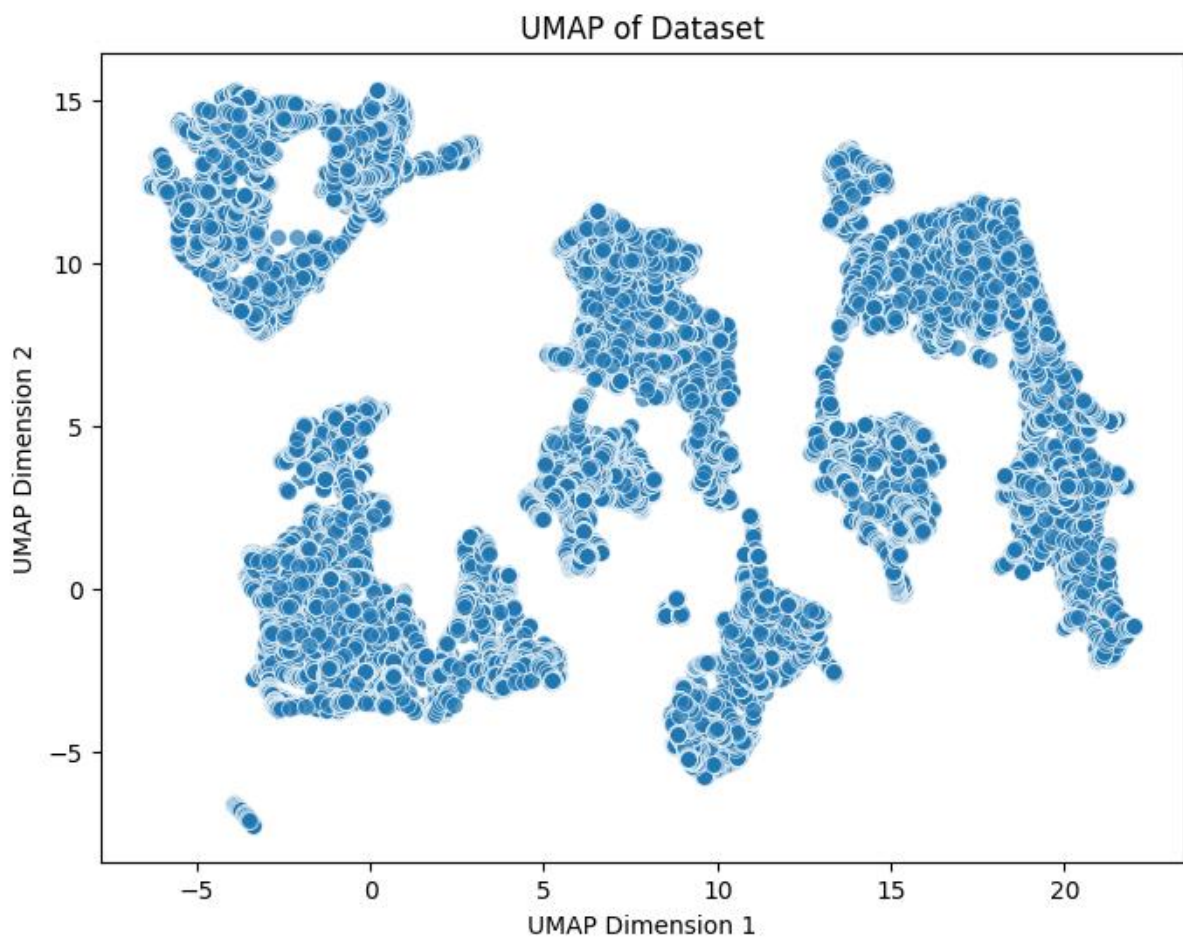


Fig 1.3 UMAP (Uniform Manifold Approximation and Projection)

Fig 1.4 shows the t-SNE (t-Distributed Stochastic Neighbor Embedding). The code applies t-SNE (t-Distributed Stochastic Neighbor Embedding) to reduce the dataset's dimensionality to two components. t-SNE is a technique often used for visualizing high-dimensional data by minimizing the divergence between probability distributions of points in high-dimensional and lower-dimensional spaces.

In this case, the numeric data is transformed into two dimensions (t-SNE1 and t-SNE2). A scatter plot is then created to show these two components, helping to visualize clusters or patterns in the data. This type of visualization is useful for detecting underlying structures or groups within the data that may not be obvious in the higher-dimensional feature space.

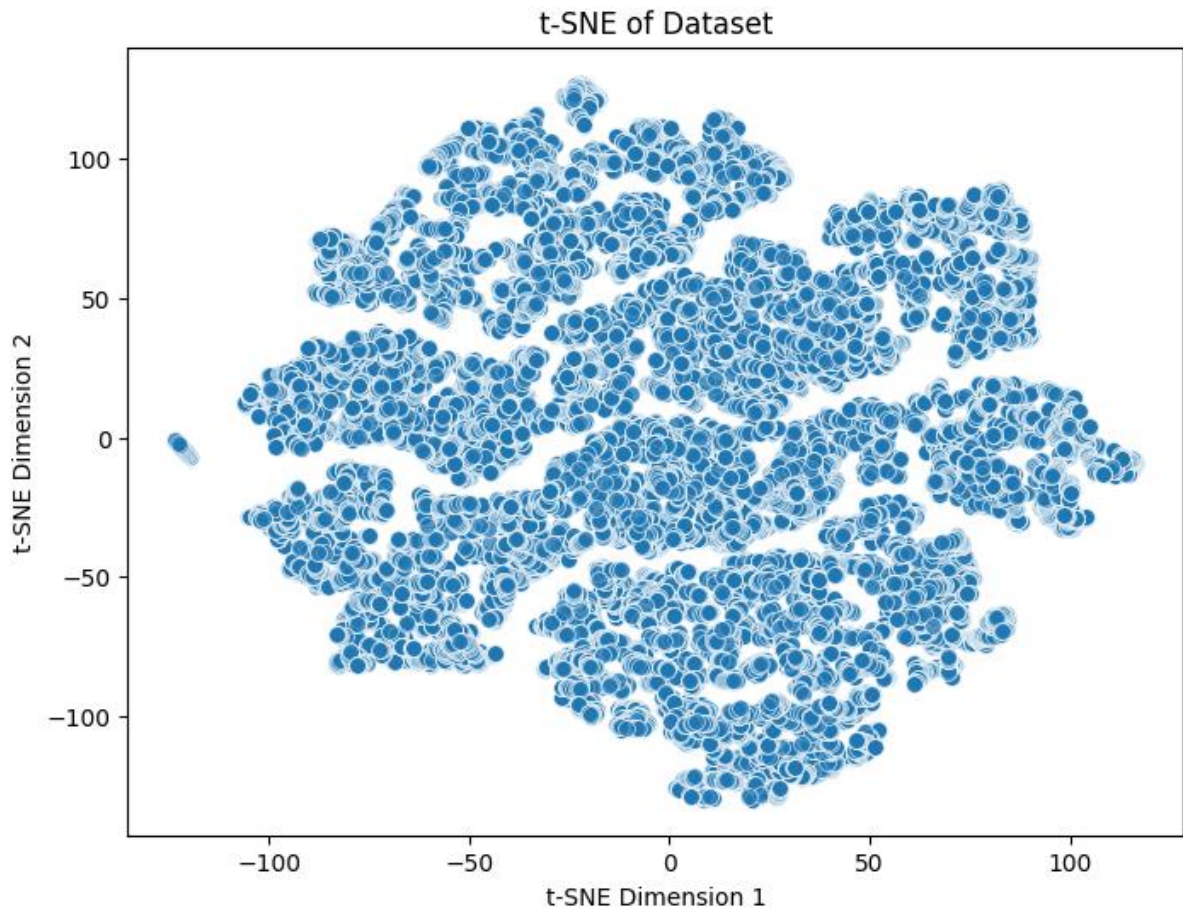


Fig 1.4 t-SNE (t-Distributed Stochastic Neighbor Embedding)

### 2.1.3 Insights from the Visualizations

The insights gained from visualizations like PCA, UMAP, and t-SNE can provide a better understanding of the dataset's structure and help identify patterns or relationships among features. Here's an explanation of potential insights from these techniques based on the dataset:

#### PCA (Principal Component Analysis)

**Dimensionality Reduction:** PCA reduces the data to two dimensions (PC1 and PC2) while preserving as much variance (information) as possible.

**Clustering or Patterns:** If distinct clusters or groupings appear in the PCA plot, it suggests that certain features in the dataset are highly correlated and may form natural groupings. For example, if the data points form tight clusters, it could indicate specific weather conditions

(e.g., clear, cloudy, rainy) are related to similar patterns of other features (like temperature and humidity).

Variance Explained: PCA also helps understand how much variance is explained by each component. If the first few components explain a high percentage of the variance, this indicates that the data's structure is well captured with a few components.

### **UMAP (Uniform Manifold Approximation and Projection):**

Non-linear Relationships: UMAP is often better at preserving non-linear structures than PCA. If the UMAP plot shows well-separated clusters, it suggests that there are distinct patterns or categories in the data that are not just linear in nature.

Visualizing Clusters: UMAP can highlight more subtle relationships. For example, it might reveal that high temperatures and humidity levels are often associated with specific weather conditions (rain, snow) that may not be immediately obvious in higher dimensions.

Local vs Global Structure: UMAP provides a balance between global and local structure, making it useful for exploring complex relationships that can exist in weather data.

### **t-SNE (t-Distributed Stochastic Neighbor Embedding)**

Local Structure Preservation: t-SNE focuses on preserving local structure, which is useful for visualizing how individual observations relate to each other. If similar weather observations (e.g., temperature, humidity) cluster together, it suggests that these features share a strong relationship.

Cluster Formation: Just like PCA and UMAP, t-SNE can reveal clusters, but it's especially good at showing tightly packed groups that may represent similar weather conditions. For example, t-SNE might reveal a cluster where temperature and humidity levels correlate closely with rain or snow conditions.

High Density Areas: Regions with dense points could indicate weather conditions that are most common in the dataset, such as typical daily weather patterns or common environmental conditions.

## General Insights:

**Feature Relationships:** The visualization techniques help reveal how features like temperature, humidity, wind speed, and precipitation interact. For example, PCA or t-SNE may show that high humidity is closely related to rain or snow conditions, or that low temperatures tend to correlate with higher precipitation.

**Anomalies and Outliers:** Clusters in UMAP or t-SNE plots might also help identify outliers or anomalies. For instance, if one or two points are far removed from other points, it might indicate unusual weather conditions that are rare in the dataset.

**Modeling Insights:** These visualizations help refine feature selection for modeling. If features form distinct groups or clusters, it suggests which features are more informative for predicting weather conditions, making them valuable for further model development and training.

By using these dimensionality reduction and visualization techniques, you can explore hidden patterns, relationships, and clusters in the data, which will aid in building more accurate and interpretable weather prediction models.

## 3 Unsupervised Learning Techniques

Let's explore the different unsupervised learning algorithms to analyze the data for predicting the performance.

### 3.1 K-Means Clustering

The Elbow Method for Optimal k plot shows the inertia (sum of squared distances from each point to its assigned cluster center) for different values of k, the number of clusters, in K-Means clustering. Here's how to interpret the plot:

- **Inertia Decrease:** The inertia starts high for smaller values of k and decreases rapidly as k increases, indicating that the points are more tightly grouped within clusters.
- **Elbow Point:** The "elbow" is located where the inertia starts to decrease at a slower rate. In this plot, it appears around  $k=3$ , suggesting that 3 clusters are optimal. After this point, increasing k yields diminishing returns in terms of reducing inertia.

This plot helps to determine the best  $k$  for clustering by balancing model complexity and performance. The optimal  $k$  is where the inertia reduction starts to slow down, as seen in this case with 3 clusters.

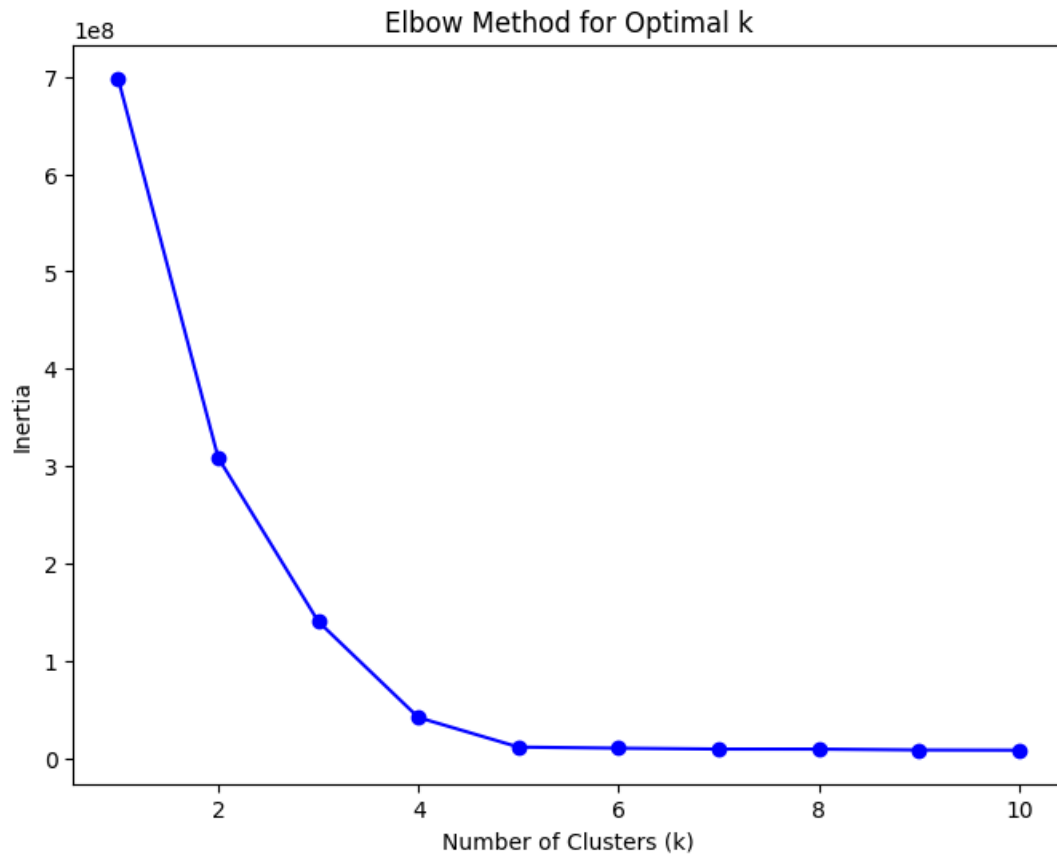


Fig 1.5 Elbow Method for optimal k

### K-Means Clustering Plot:

The second plot shows the results of K-Means clustering applied to the dataset using PCA (Principal Component Analysis) for dimensionality reduction. This plot visualizes how the data points are grouped into 3 clusters (based on the optimal  $k$  from the elbow method) in the 2D PCA space.

- Clusters: The data points are coloured according to the cluster they belong to. The clusters are well-separated, with each cluster corresponding to distinct groups of data in the lower-dimensional space.
- PCA: The PCA components show the major variance directions of the data in a reduced 2D space, helping in visualizing the cluster structure.

The combination of the elbow method and the clustering plot demonstrates how K-Means identifies distinct groups in the dataset, and using dimensionality reduction techniques like PCA helps in visualizing these groupings in a simpler form.

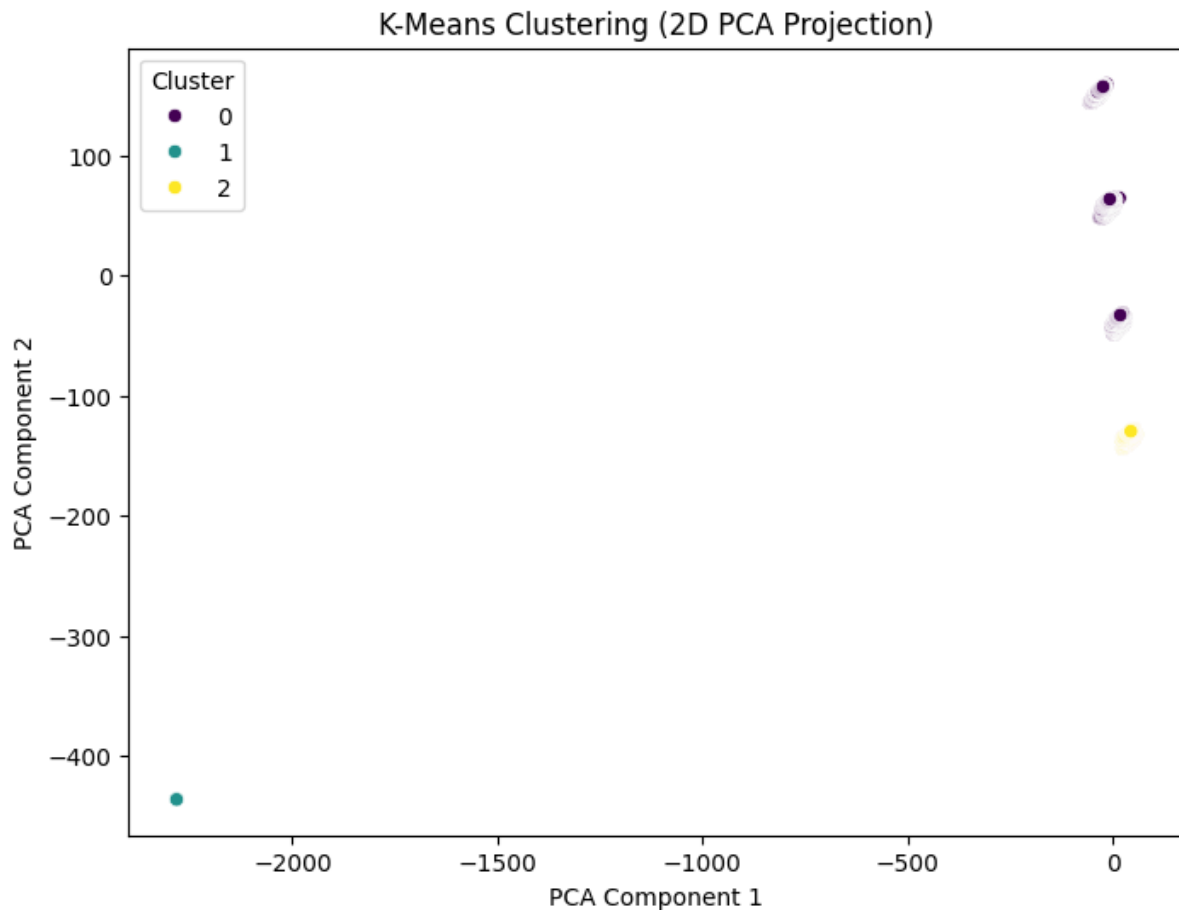


Fig 1.6 K-Means Clustering

## 3.2 Hierarchical Clustering

The Hierarchical Clustering Dendrogram provides a tree-like structure that illustrates how data points are grouped together based on their similarity.

- **Euclidean Distance:** The vertical axis represents the Euclidean distance, which indicates how far apart the clusters are in terms of their feature values. The higher the value on the axis, the more distinct the clusters are from each other.
- **Tree Structure:** The horizontal axis represents the data points, and the lines connecting clusters show how data points are grouped as the distance between them decreases.

At the bottom of the dendrogram, each data point starts as its own cluster, and as we move upwards, the clusters merge.

- **Ward's Method:** The method used here is Ward's linkage, which minimizes the variance within clusters. The clusters are merged based on the smallest increase in total within-cluster variance.
- **Color Coding:** The plot also visualizes different groups in different colors. For instance, the green and red branches of the tree indicate separate clusters formed at a certain threshold of Euclidean distance.

### **Insights:**

- **Cluster Identification:** By cutting the dendrogram at a specific height, we can determine the number of clusters that best represent the data. The height at which clusters merge shows how similar the data points are within those clusters.
- **Cluster Relationships:** The branching patterns suggest how closely related certain data points or groups are. For example, if two branches join at a low distance, those data points are very similar to each other.

This dendrogram can guide the choice of the optimal number of clusters by visually identifying the appropriate cut-off point for cluster formation.



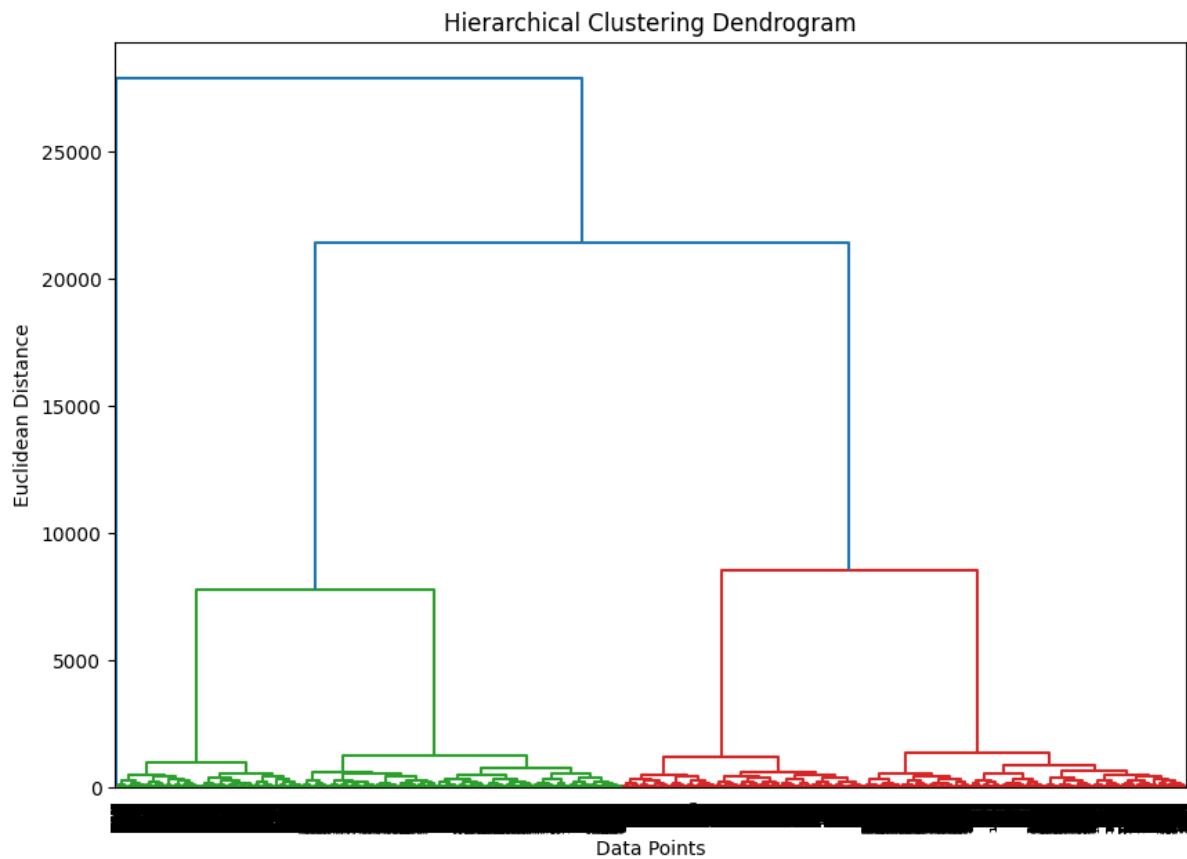


Fig 1.7 Hierarchical Clustering

### 3.3 DBSCAN Clustering

In the DBSCAN (Density-Based Spatial Clustering of Applications with Noise) clustering approach:

- Eps (epsilon): This parameter defines the radius of the neighborhood around a data point. In this case, it's set to 0.5. Points within this radius of each other are considered neighbors.
- Min\_samples: This parameter defines the minimum number of points required to form a dense region. Here, it is set to 5, meaning that at least 5 points within the eps distance are needed to form a cluster.

**Key Elements:**

- **Cluster Assignment:** The `fit_predict` method assigns each data point to a cluster. Data points labeled -1 are considered noise, meaning they do not belong to any cluster.
- **Visualization:** The results are projected onto 2D using PCA for better visualization. The points are colored based on their assigned cluster from DBSCAN.

**Insights:**

- **Cluster Identification:** DBSCAN is useful for identifying clusters of varying shapes and sizes. It does not require the number of clusters to be predefined. It also efficiently handles noise in the data by assigning a cluster label of -1 to outliers.
- **Density-based Clustering:** Unlike K-means, which assumes spherical clusters, DBSCAN can detect arbitrarily shaped clusters based on density.
- **Noise Detection:** The noise points, marked by -1, can be observed in the scatter plot, showing that DBSCAN has successfully identified points that do not belong to any dense region.

This visualization helps in understanding the structure of the data and can be a powerful tool when dealing with irregular data distributions.

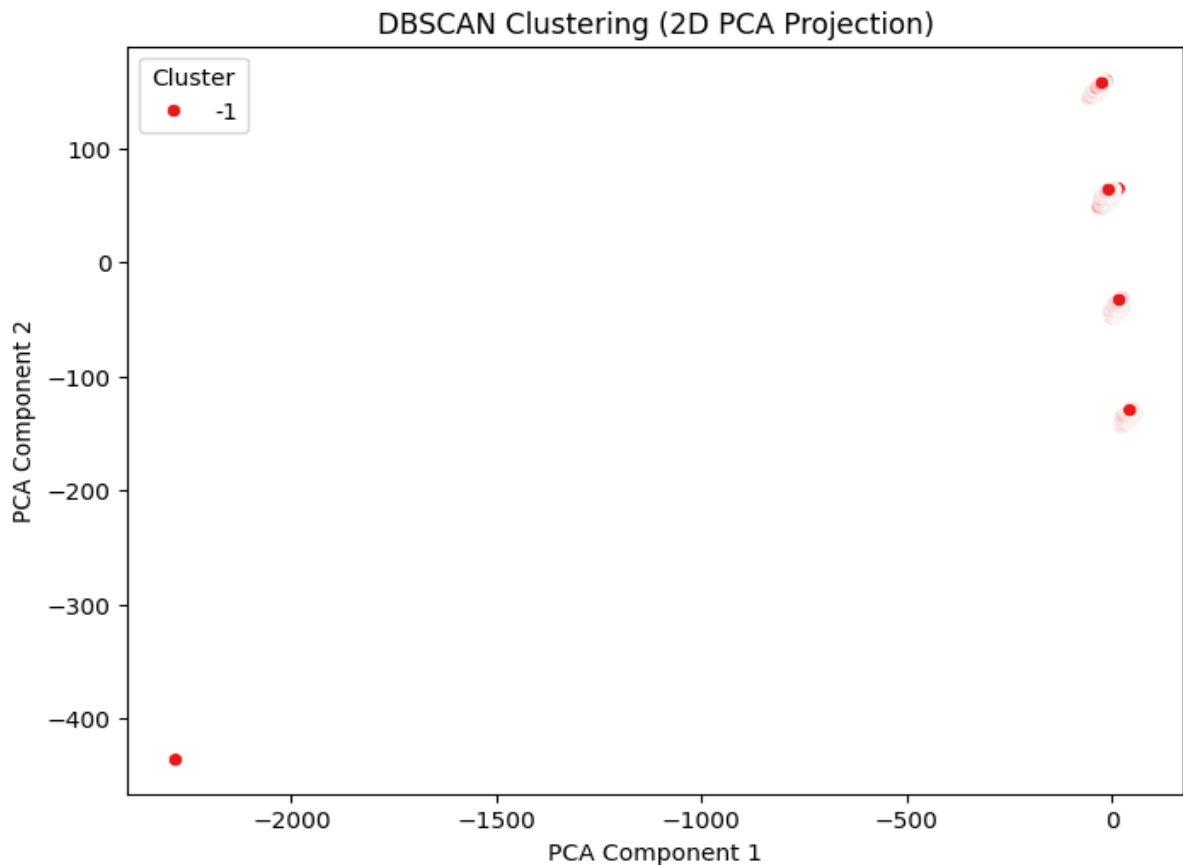


Fig 1.8 DBSCAN Clustering

### 3.4 Advanced Dimensionality Reduction with UMAP

In this UMAP (Uniform Manifold Approximation and Projection) dimensionality reduction approach:

- UMAP is a non-linear dimensionality reduction technique that preserves both local and global structure in the data while reducing it to a lower-dimensional space, making it ideal for visualization in 2D or 3D.
- In this code, the dataset is projected into 2 dimensions (`n_components=2`), while the original dataset's numeric features are selected for processing.
- The K-Means clusters obtained earlier are used to color the points in the UMAP plot. This helps us to visualize how the clusters are distributed in the reduced space.

**Insights:**

1. **Cluster Separation:** UMAP helps in visually assessing how well the K-Means clusters are separated in the lower-dimensional space. If the clusters are well-separated in the plot, it indicates that K-Means clustering has successfully grouped similar data points together.
2. **Visualization of High-dimensional Data:** UMAP is particularly effective in maintaining the local structure of the data while reducing the dimensions. This makes it useful for visualizing high-dimensional datasets and identifying patterns that might not be apparent in the original feature space.
3. **Cluster Consistency:** By coloring the data points according to the K-Means clusters, we can assess the consistency of the clustering algorithm. If the clusters are distinct in the UMAP plot, it suggests that the algorithm performed well. If the points from different clusters overlap, it might indicate the need for refinement in the clustering process (perhaps a different number of clusters or a different algorithm).

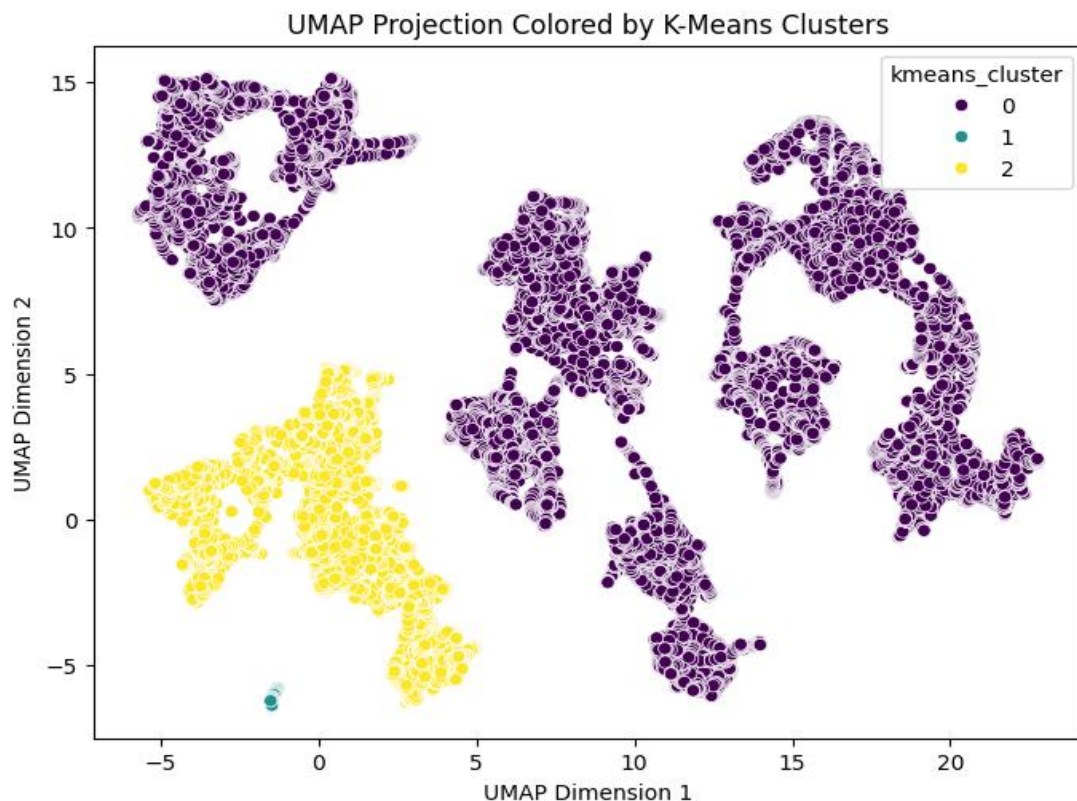


Fig 1.9 Advanced Dimensionality Reduction with UMAP

### 3.5 Anomaly detection with Isolation Forest

Anomaly detection identifies unusual data points or patterns that deviate significantly from the norm. Isolation Forest is an unsupervised anomaly detection algorithm that isolates anomalies efficiently by building random trees. It works on the principle that anomalies, being distinct and rare, are easier to isolate with fewer splits. The algorithm assigns an anomaly score based on how quickly a data point is isolated, making it a scalable and effective method for detecting outliers in large, high-dimensional datasets.

This approach focuses on isolating anomalies by detecting data points that deviate significantly from the majority. After loading the dataset, numeric features are extracted for analysis. The Isolation Forest model is trained with a contamination parameter set to 5%, indicating the expected proportion of anomalies in the data. The model assigns a label of -1 to outliers and 1 to normal points, and these labels are added to the dataset under a new column called `Outlier_IF`.

To visualize the results, a scatter plot is created, using the index of each data point on the x-axis and a feature (e.g., TEMP) on the y-axis. Normal points are displayed in blue, while anomalies are shown in red, effectively highlighting the deviations in temperature data. This visualization provides an intuitive way to interpret the presence and distribution of anomalies within the dataset.

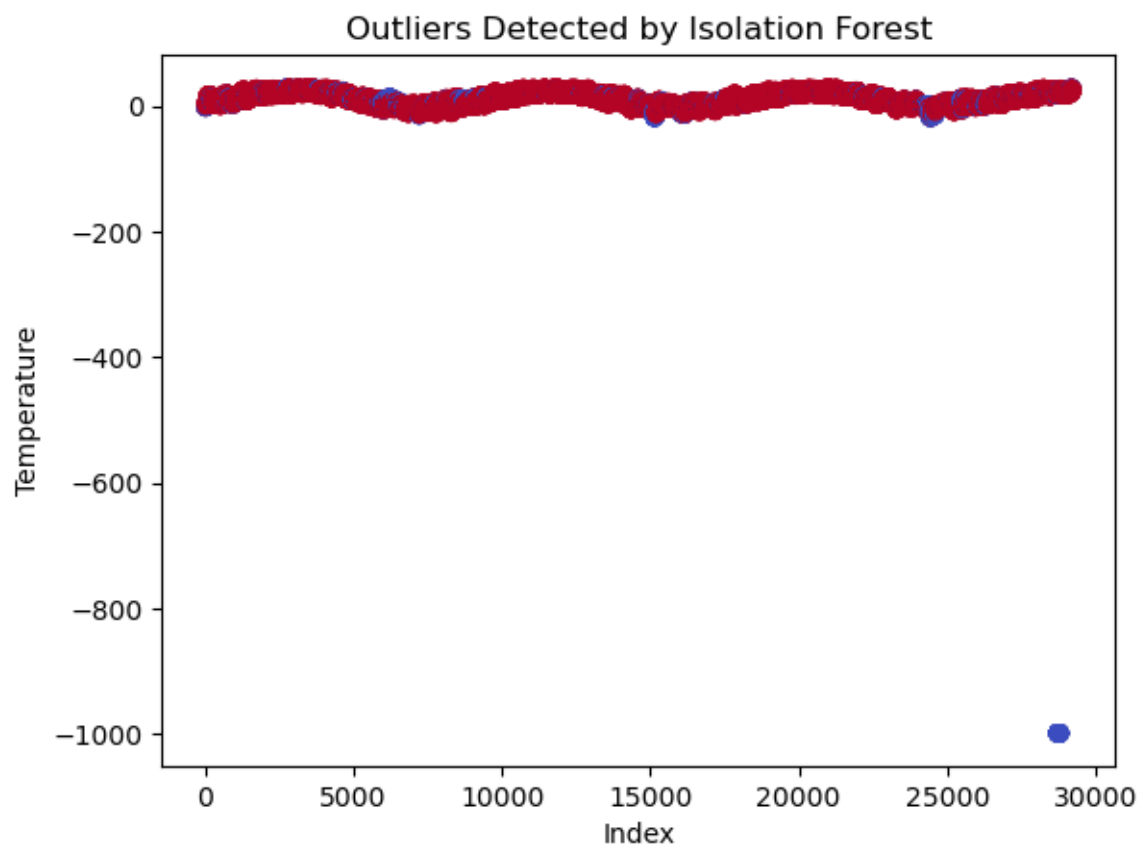


Fig 1.10 Anomaly detection with Isolation Forest

### 3.6 Anomaly detection with One-class SVM

**One-Class SVM** is an unsupervised learning algorithm used for anomaly detection. It learns the distribution of normal data and identifies outliers by classifying data points that fall outside this learned boundary as anomalies. The model uses a kernel function (like RBF) to create a decision boundary and can detect outliers in high-dimensional and complex data. Key parameters include  $\nu$ , which controls the proportion of outliers, and  $\gamma$ , which

affects the shape of the boundary. This approach is commonly used in applications such as fraud detection and defect analysis

The plot generated using the One-Class SVM model highlights outliers within the Temperature data. Here are key insights:

- **Outlier Detection:** The red points indicate data instances that deviate significantly from the majority of the data, suggesting unusual or anomalous temperature readings.
- **Distribution Insight:** Most of the data points (shown in blue) are considered typical temperatures that conform to the learned model. The presence of outliers implies that there are specific instances with temperatures far from the norm, possibly indicative of errors, rare events, or special cases.
- **Application:** This type of anomaly detection can be useful in identifying data quality issues, monitoring systems for abnormal behaviour, or flagging potential incidents in time-series data.



Fig 1.11 Anomaly detection with one-class SVM

### 3.7 Association rule with Apriori Algorithm

**Association rule learning** is a technique used to identify relationships between items in large datasets. It is commonly used in market basket analysis to discover patterns, like "customers who buy bread often also buy butter."

**FP-Growth** is an efficient algorithm that finds **frequent itemsets** by constructing an FP-tree, which helps identify items that appear together frequently without scanning the database repeatedly.

**Association rules** are generated from these frequent itemsets, evaluated based on **support** (how often items appear together) and **confidence** (how likely the rule is true). This method is useful for making predictions, recommendations, and uncovering hidden patterns.

This code uses the FP-Growth algorithm to analyze weather data by categorizing temperature, precipitation, and humidity into "Low," "Medium," and "High" bins. It creates transactions from these categories, finds frequent itemsets with a minimum support of 10%, and generates association rules with at least 70% confidence.

#### Key Insights:

- **Frequent Patterns:** Combinations of weather conditions that appear frequently together.
- **Association Rules:** Insights like "If temperature is 'Medium' and precipitation is 'Low', then humidity is likely 'High' with a confidence of 70%."

### 3.8 Feature Transformation with Auto Encoders

An **autoencoder** is a type of neural network used for unsupervised learning that compresses input data into a lower-dimensional representation (encoding) and reconstructs it to match the original input. This approach is beneficial for **dimensionality reduction**, **feature extraction**, and **anomaly detection**.



The code provided demonstrates the use of an autoencoder model to learn a compressed representation of input data. The encoder part of the model reduces the dimensionality of the data to a 2D space. This compressed data can be visualized in a scatter plot, where the axes represent the encoded dimensions. By plotting the encoded data, we can gain insights into patterns, clusters, or potential outliers within the dataset, helping to understand the relationships and distribution of features after reduction. This visualization can be particularly useful for identifying structures within the data that may not be apparent in higher dimensions

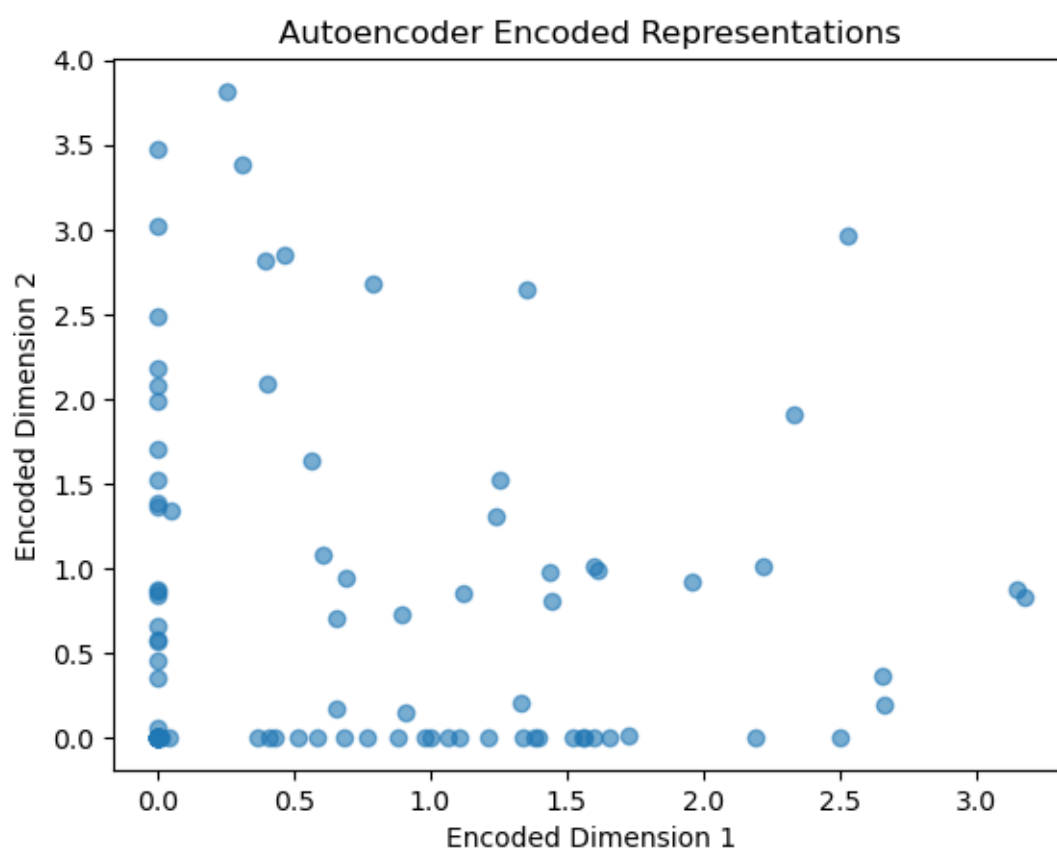


Fig 1.12 Feature Transformation with Auto Encoders

### 3.9 Feature Transformation with Kernel PCA

Kernel PCA (Principal Component Analysis) is an extension of standard PCA that uses kernel methods to handle non-linear relationships in data. While traditional PCA is limited to linear dimensionality reduction, Kernel PCA can project data into higher-dimensional space using a kernel function (such as RBF), enabling the capture of non-linear structures.

In the provided code, Kernel PCA is applied to scaled data, reducing it to two dimensions for visualization. The gamma parameter determines the influence of each training sample, with smaller values giving broader influence and larger values creating more localized influence. The results are then plotted, showcasing how data points are distributed in the new feature space. This helps reveal patterns or clusters in data that linear methods might miss, providing a deeper understanding of complex data structures.

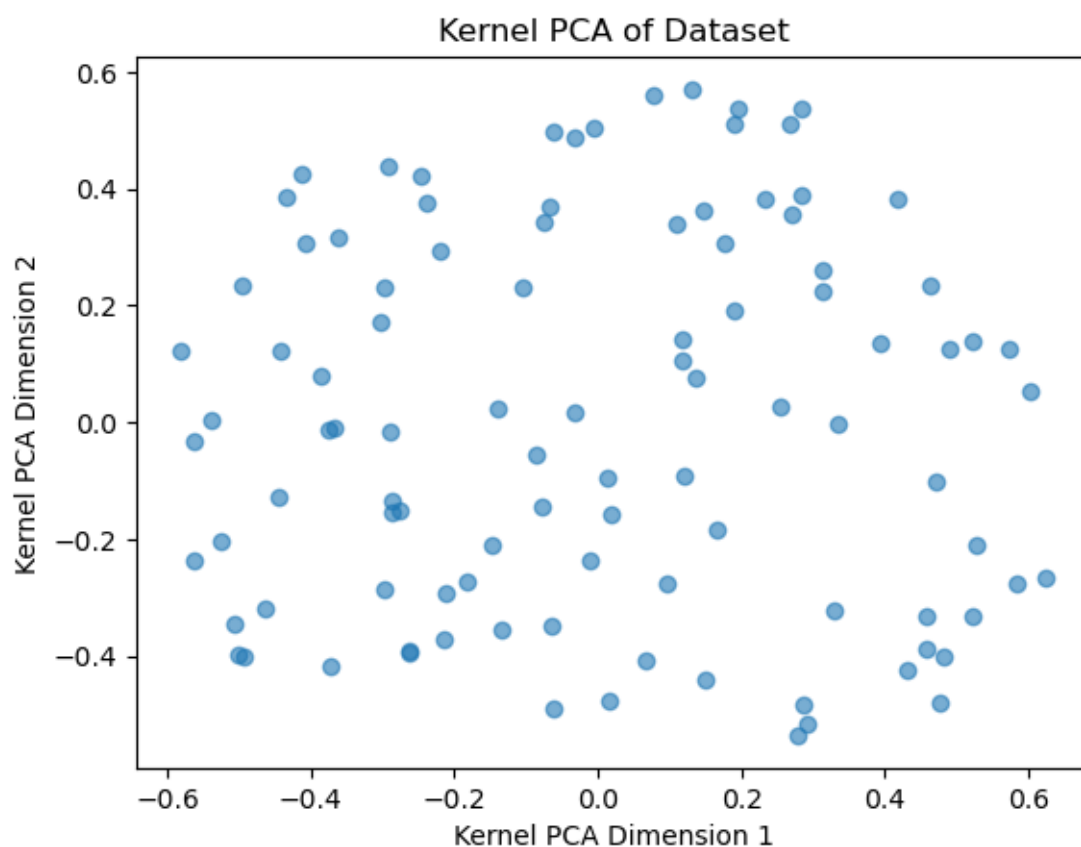


Fig 1.13 Feature Transformation with Kernel PCA

## 4 Cross Validation Strategies

Cross-validation is a technique used to assess the performance and generalizability of a machine learning model. It helps to ensure that the model performs well not just on a training set but also on unseen data. This technique involves splitting the data into multiple subsets or "folds" and iteratively training the model on some folds while testing it on the remaining fold(s). Below are some common cross-validation strategies.

### 4.1 Stability-based cross-validation for dimensionality reduction

Stability-based cross-validation for dimensionality reduction is a method used to evaluate how robust a technique like Principal Component Analysis (PCA) is when applied to different subsets of data. The idea is to split the dataset into two or more parts and apply PCA to each part separately. By comparing the results, you can determine whether the reduced dimensions are consistent across different subsets of the data.

In this approach:

- Data Splitting: The dataset is divided into two halves or subsets.
- PCA Application: PCA is performed on each subset to extract principal components.
- Visualization and Comparison: The principal components from each subset are visualized to assess consistency.

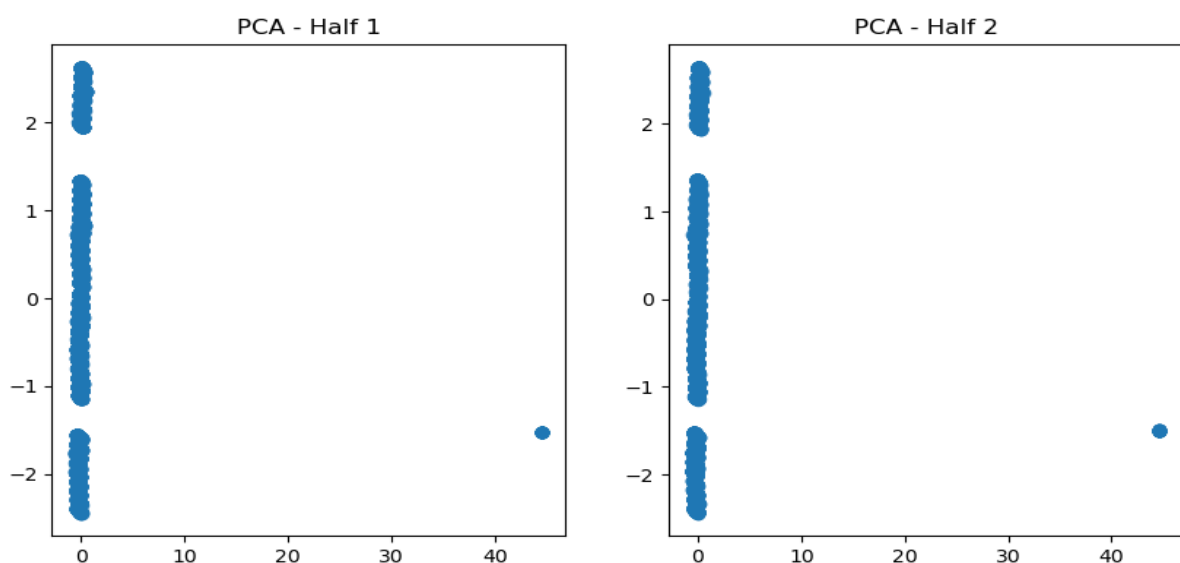


Fig 1.14 Stability-based cross-validation for dimensionality reduction

## 4.2 Clustering Stability Validation

Clustering stability validation is a technique used to assess the reliability and robustness of clustering results. It involves running clustering algorithms on different subsets or random splits of a dataset multiple times and comparing the results to ensure consistency. The goal is to determine whether the algorithm identifies similar groupings across different data partitions. Common metrics used for this evaluation include the **Adjusted Rand Index (ARI)**, which measures the agreement between two clustering results, and **Silhouette Score**, which assesses how well-separated the clusters are. Stable clustering models show high consistency and are more trustworthy for real-world applications, as they reliably group similar data points together.

The `cluster_stability_kmeans` function evaluates the stability of K-Means clustering by splitting the data into training and testing sets multiple times and calculating the Adjusted Rand Index (ARI) for each split. This index measures the agreement between cluster assignments for different splits, with scores ranging from -1 (no agreement) to 1 (perfect agreement).

The function calculates ARI for each split and reports the mean ARI as an overall stability metric. High mean ARI values indicate stable and consistent clustering, while lower values suggest that clustering results are sensitive to data variations. This approach helps assess the reliability of the K-Means algorithm for the given dataset.

## 4.3 Anomaly Stability with Isolation Forest

Anomaly stability with Isolation Forest checks how consistent the model's anomaly detection results are across different data splits. The data is divided into training and test sets multiple times, with the model trained on the training set and tested on the test set. If the model consistently identifies the same observations as anomalies, it indicates stability, confirming the model's reliability for detecting outliers in varying data.

The `anomaly_stability_isolation_forest` function tests the stability of the Isolation Forest algorithm for anomaly detection. It splits the dataset into training and test sets, trains the model on the training data, and predicts anomalies on the test data. The function then prints

how many anomalies were detected in the test set, helping to assess the model's reliability in detecting anomalies when applied to different data samples

The function evaluates the consistency of the Isolation Forest model's anomaly detection by testing it on different data subsets. This helps determine if the model can reliably detect anomalies in new data. For more thorough validation, techniques like cross-validation with multiple splits or using ARI scores can be added to better assess model stability.

#### **4.4 Cross-Validation of Association Rule Learning**

Cross-validation of association rule learning involves partitioning the dataset into multiple subsets to train and test the model's ability to generate meaningful rules consistently. This process helps evaluate the robustness and generalizability of the rules identified by ensuring that they hold across different data samples. Unlike traditional cross-validation used in predictive models, association rule cross-validation focuses on validating the frequency and confidence of the rules across subsets, which can help in identifying rules that are truly significant and not specific to a single subset of data. This approach provides a more reliable assessment of rule stability and performance when applied to unseen data.

The code uses the Apriori algorithm to discover association rules in a dataset and validates them by comparing rules from a training set and a test set. The data is preprocessed by binning numeric columns, converting them to one-hot encoded binary format, and splitting into training and test subsets. The Apriori algorithm generates frequent itemsets, from which rules are created based on the lift metric. The rules from both the training and test sets are displayed to check for consistency. This helps assess the stability and generalizability of the rules, ensuring they are reliable when applied to new data.

## 5 Identify Pitfalls and Performance Measure

Identifying pitfalls and performance measures in data analysis and machine learning is crucial for ensuring that models and methodologies provide reliable and accurate insights.

### Pitfalls:

- **Overfitting:** Model fits training data too closely, impacting new data performance.
- **Underfitting:** Model is too simplistic, missing key patterns.
- **Data Issues:** Poor quality or biased data can skew results.
- **Feature Problems:** Irrelevant or excessive features can degrade model accuracy.
- **Evaluation Bias:** Using inappropriate metrics can misrepresent model quality.
- **Data Leakage:** Including future information during training can inflate results.

### Performance Measures:

- **Accuracy:** General success rate; may not be reliable for imbalanced data.
- **Precision, Recall, F1-Score:** Evaluate classification models, especially for imbalanced classes.
- **AUC-ROC:** Measures how well the model separates classes.
- **MAE/MSE:** Metrics for regression tasks that assess average error.
- **Adjusted Rand Index (ARI):** Used for clustering validation.
- **Cross-Validation:** Splits data into training and validation sets to ensure generalizability.

In our Project,

Overfitting: If cross-validation scores are much higher than test accuracy.

Underfitting: If cross-validation scores are low, it indicates underfitting.

Hyperparameter tuning: The performance depends on the proper choice of hyperparameters.

Early Stopping: Although not necessary for Logistic Regression, it's useful in deep learning models.

## 6 Performance Enhancement Techniques

### 6.1 Experiment 1

#### Feature Selection and Regularization with Logistic Regression (L1 + L2)

Feature selection and regularization are techniques used in machine learning to enhance model performance and reduce overfitting. Feature selection helps choose the most relevant features to build a simpler and more interpretable model. Regularization adds a penalty to the loss function to discourage overly complex models.

In Logistic Regression:

- L1 Regularization (Lasso) can shrink some feature coefficients to zero, effectively selecting a subset of features.
- L2 Regularization (Ridge) reduces the impact of less important features but keeps them in the model.

Combining these techniques leads to models that generalize better and are more efficient.

The code snippet applies feature selection and evaluates logistic regression models with L1 (Lasso) and L2 (Ridge) regularization.

- L1 Regularization (Lasso) encourages sparsity by shrinking some feature coefficients to zero, aiding in feature selection.
- L2 Regularization (Ridge) penalizes large coefficients but keeps all features, helping to avoid overfitting.
- Feature Selection with RFE identifies the top 3 most important features to improve model efficiency and performance.

The code uses cross-validation to compare the accuracy of models trained with L1 and L2 regularization on the selected features, showing which approach performs better on the dataset.

## 6.2 Experiment 2

### Increasing Model Complexity with Random Forest

Random Forest is an ensemble learning method that creates multiple decision trees and combines their predictions to improve accuracy and reduce overfitting. It is effective for handling complex data and high-dimensional spaces due to its random sampling of data and features. While it is robust and can handle noise, it may become computationally intensive with a large number of trees. Hyperparameter tuning, such as adjusting the number of trees (`n_estimators`) or the depth of trees (`max_depth`), helps optimize performance.

The code snippet demonstrates using Random Forest with hyperparameter tuning to improve model performance. It employs `GridSearchCV` to find the best combination of hyperparameters (`n_estimators`, `max_depth`, `min_samples_split`) by evaluating multiple configurations using cross-validation.

#### Key Steps:

1. Model Creation: A `RandomForestClassifier` is initialized.
2. Hyperparameter Tuning: `GridSearchCV` tests various parameter combinations.
3. Model Fitting: The best parameters are chosen and fitted to the training data.
4. Evaluation: Cross-validation scores are calculated to assess the model's accuracy.

#### Output:

- Best hyperparameters for the Random Forest model.
- Cross-validation scores and the average accuracy, indicating the model's robustness.



## 6.3 Experiment 3

### Ensemble Method - Stacking

Stacking combines the predictions of multiple base models (Logistic Regression, Random Forest, and SVC) and uses a final estimator (here, another Logistic Regression) to make the final prediction based on the base models' outputs.

#### Key Points:

- **Base Estimators:** Models like Logistic Regression, Random Forest, and Support Vector Classifier (SVC) are used as base learners.
- **Final Estimator:** A Logistic Regression model combines the outputs of the base models to make the final prediction.
- **Cross-Validation:** The model is evaluated with 5-fold cross-validation to estimate its performance. The mean score gives an overall idea of how well the stacking ensemble generalizes to new data.

#### Benefits:

- Stacking can leverage the strengths of different models, potentially leading to better performance compared to using any single model.
- The final estimator can adapt to combine base models effectively, enhancing predictive accuracy.