# COMPARATIVE ANALYSIS OF ML ALGORITHM IN DROUGHT PREDICTION

*Report submitted to*

*Haldia Institute of Technology*

*Haldia for the award of the degree*

*Of*

**Bachelor of Technology in Computer Science & Engineering**

*By*

**Kritik Tiwary(10300120079)**

**Mayank Kumar(10300120086)**

**Priyanshu(10300120110)**

**Rahul Raman(10300120114)**

**DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING**
**HALDIA INSTITUTE OF TECHNOLOGY, HALDIA**

**DECLARATION**

We certify that

   a.  The work contained in this report is original and has been done by me/us under the guidance of my/our supervisor(s).

   b.  The work has not been submitted to any other Institute for any degree or diploma.

   c.  I/We have followed the guidelines provided by the Institute in preparing the report.

   d.  I/We have conformed to the norms and guidelines given in the Ethical Code of Conduct of the Institute. e. Whenever I/we have used materials (data, theoretical analysis, figures, and text) from other sources, I / we have given due credit to them by citing them in the text of the report and giving their details in the references.

Signature of the Students

# CERTIFICATE

This is to certify that the Dissertation Report entitled, "**Comparative Analysis of ML Algorithms in Drought Prediction**" submitted **by Mr. Kritik Tiwary , Mr. Mayank Kumar, Mr. Priyanshu , Mr. Rahul Raman** to Haldia Institute of Technology, Haldia, India, is a record of bonafide Project work carried out by him/her under my/our supervision and guidance and is worthy of consideration for the award of the degree of Bachelor of Technology in Computer Science and Engineering of the Institute.

.....................................
Project Mentor

…………………...
HOD

…………………..
Project Coordinator

III

# ACKNOWLEDGEMENT

We would like to express our sincere gratitude to all those who played a pivotal role in the successful completion of our final year project, "Comparative Analysis of ML Algorithms in Drought Prediction."

First and foremost, our deepest appreciation goes to Dr. Arindam Giri, our dedicated mentor, whose unwavering guidance and expertise have been instrumental in shaping our project. Dr. Arindam Giri's insightful feedback and encouragement propelled us forward, and we are truly grateful for the wealth of knowledge he shared with us throughout this journey.

Our heartfelt thanks also extend to our fellow group members for their collaborative spirit and tireless efforts. The synergy within our team significantly contributed to the depth and quality of our research, making the project a fulfilling and rewarding experience.

We are indebted to the faculty members who provided valuable insights and constructive criticism, aiding us in refining our methodology and expanding our understanding of machine learning algorithms in the context of drought prediction

Lastly, we would like to acknowledge the support of our friends and family whose encouragement and understanding sustained us during the challenges of this project.

This project has been a journey of growth, learning, and collaboration, and each person mentioned above has played an integral role in its success. Thank you for being part of this enriching experience.

# ABSTRACT

Drought is a major global challenge with significant economic and environmental consequences. Early and accurate prediction of drought is crucial for proactive mitigation and management strategies. Machine learning (ML) algorithms have emerged as promising tools for drought prediction due to their ability to handle complex spatiotemporal relationships. This project investigates the performance of various ML algorithms in predicting drought occurrence and severity.

This final year project undertakes a holistic exploration of drought prediction efficacy through the lens of machine learning (ML) algorithms. The project unfolds in sequential phases. Rigorous data wrangling techniques are applied to ensure data integrity and completeness, setting the foundation for subsequent analyses.

The project advances through an extensive phase of exploratory data analysis (EDA), unraveling intricate patterns, trends, and anomalies within the collected datasets. This process not only enhances our understanding of the data but also informs subsequent feature selection and extraction strategies.

In the feature-rich landscape, the project delves into the comparative analysis of prominent ML algorithms, including SVM, Decision Trees and KNN.

The outcomes of this project promise to contribute to the evolving field of drought prediction, offering insights into the nuanced performance of ML algorithms in diverse environmental contexts.This project will contribute valuable knowledge to the domain of drought forecasting and promote the application of ML in mitigating the devastating effects of drought.

# CONTENTS

# INDEX OF IMAGES

# **Chapter-1 Introduction**

Droughts remain to be one of the most dangerous hazards, having a serious and large-scale impact on environment, society and economy. Recent events like the Summer 2018 drought in huge parts of Central Europe led to severe forest fires and crop failures. The damage was estimated to several hundred millions euros solely in Germany (Federal Ministry of Food and Agriculture, 2018). Moreover the effect of global warming leads to major changes in the earth's climate system, having a direct influence on the frequency and severity of extreme events like droughts (Spinoni et al., 2016). An increase in frequency of drought occurrence is a major threat for current and future generations, and comprehensive knowledge on the phenomenon of drought is needed in order to take action early and to prevent humanitarian catastrophes. This goes in conjunction with drought prediction. A tool for drought prediction would enable to mitigate the dangers connected to drought occurrences, such that it would advise stakeholders to store the maximal possible amount of water in the endangered regions. This would help to mitigate the water shortage when the drought arrives. Measures for demand reduction could like that be introduced earlier, and in better adjusted extent; this would help to reduce the economic and societal damage. To mitigate the effects of droughts the information on the their onset is of crucial importance. This can be derived from a drought index. A variety of drought indices exist, which are typically defined according to statistical and physical measures. These are mostly taking into account atmospheric and soil variables. Among the most popular ones are the Standardized Precipitation Index (SPI), Standardized Precipitation Evaporation Index (SPEI), Soil Moisture Percentile (SMP), and Palmer Drought Severity Index (PDSI). Standardized Precipitation Index (SPI) is adopted as the standard meteorological index by World Meteorological Organization (2012). It is a measure of meteorological drought based on the probability of occurrence of certain

precipitation amounts in the area of interest (Sheffield and Wood, 2011). Studies on drought prediction by Belayneh et al. (2016) and Bonaccorso et al. (2015) use SPI as a prediction variable for the forecast. Forecasting of any physical phenomenon can either be done by a physical, conceptual or data-driven model. The latter ones are widely used due to their rapid development times and the flexibility in input parameters. McGovern et al. (2017) argues that AI-methods have a high potential for prediction of extremes due to the ability of machine learning methods to learn from past data, to handle large amounts of input variables, to integrate physical understanding into the models and to discover additional knowledge from the data. A review on seasonal drought prediction given by Hao et al. (2018) identifies two typical predictor groups of variables: large-scale climate indices that reflect the atmosphere-ocean circulation patterns and local climate variables. The first ones are known to correlate with the precipitation patterns in special regions and therefore are naturally correlated with the occurrence of drought. The teleconnection indices important for European precipitation include North Atlantic Oscillation (NAO), Scandinavian Oscillation (SCA), East Atlantic/Western Russia Oscillation (EA/WR), East Atlantic Oscillation (EA) and Atlantic Multidecadal Oscillation (AMO) (Hao et al., 2018). As shown by Folland et al. (2009) a positive NAO index in summer is associated with dry and warm conditions in the north-west of Europe, whereas southern Europe and the Mediterranean experience cooler and wetter conditions. More information on the influence of the NAO, SCA, EA, and EAWR on the European climate can be found in Folland et al. (2009), Bueh and Nakamura (2007), Mikhailova and Yurovsky (2016), Lim (2015), Barnston and Livezey (1987) and Sheffield et al. (2009). A positive phase of AMO is associated with humid conditions over Great Britain and parts of Scandinavia and with dry conditions in the Mediterranean (Sheffield and Wood, 2011, p. 26); the negative phase is associated with a reversed pattern: dry conditions in Great Britain and wet conditions in the Mediterranean. A study by Sheffield et al. (2009) showed a

correlation between the amount of droughts and AMO of 62% with a significance at the 90% level. A recent study by Bonaccorso et al. (2015) uses NAO for prediction of probability of drought occurrence for Sicily. The local climate variables like precipitation, temperature, soil moisture were also used as inputs to reflect the conditions at the time the prediction occurs. Belayneh et al. (2016) and Bonaccorso et al. (2015) used SPI for the past months as input variable to the algorithm. A study by Morid et al. (2007) used precipitation as an input parameter. This paper examines the possibilities of meteorological drought prediction with the lead time of one month applying artificial neural networks (ANN) for two domains with different climate: one with Mediterranean (Lisbon), one with continental climate (Munich) (Ceglar et al., 2019). Both sites experienced an increase of drought frequency when comparing 2015 and 1950 and are projected to keep rising under RCP4.5 as well as RCP 8.5. (Spinoni et al., 2017). Observational data offers only a limited field for drought investigation as it can be seen from the following approximation. Systematical weather observations started in 1781 by Societas Meteorologica Palatina (Kington, 1980). In this study SPI1<-1 is used as a threshold for drought occurrence. It corresponds to the 15% driest months (John Keyantash, 2018) and can be estimated by a total amount of 430 events up to the year 2020 (Eq. 1).

$$(2020 - 1781) \: yr \cdot 12 \: months/yr \cdot 15\% = 430 \: events \qquad (1)$$

Compared to that CRCM5-LE offers a total amount of roughly 4500 events when using the first 50 years from the climate simulation data (1955-2005) (see Eq. 2).

$$50 \: yr/member \cdot 50 \: members \cdot 12 \: months/yr \cdot 15\% = 4500 \: events \qquad (2)$$

This is a difference of an order of magnitude. The more data is available the better the predictions that can be derived by a drought predicting machine learning model and the more can be learned about drought formation. According to von Trentini et al. (2020) precipitation

in summer and winter derived from the European gridded data set (E-OBS) does fall to a high percentage into the range produced by CRCM5-LE for the historic period. Therefore, the CRCM5-LE proves applicable to this study and its larger amount of extreme events can be used as input to the machine learning algorithms. In this study a variety of ANNs were trained. Best performing models were investigated to using explainable AI methods to understand the results. While no comparable study exists for the Munich domain, Santos et al. (2014) performed a drought prediction based on SPI6 for Portugal for the months April, May and June using the following input variables: sea surface temperatures (JFM), NAO (DJFM) and cumulative precipitation (NDJFM for SP I6April, DJFM for SP I6May, JFM for SP I6June). Best results were achieved for the prediction of SPI6 for April with a correlation coefficient of 0.98. SPI6 for May and June referred to a correlation coefficient of 0.78 and 0.77 respectively.

# Chapter-2 Machine Learning

This study investigates drought predictability applying the technique of supervised machine learning for this purpose. Machine learning is a promising tool for the analysis of complex and data-rich phenomena as droughts (McGovern et al., 2017). The python package Keras, a high-level neural network package, was used for the design of the machine learning models (Chollet et al., 2015), as it allows to design neural networks in an easy way by adding layers. Three crucial elements are needed to perform drought prediction by supervised machine learning: input data, a target variable to be predicted and a computation pipeline, which includes the machine learning algorithm. The data from the years 1957 - 1999 was used as training data, the years 2000-2005 were used for the testing purpose. A small fraction of the training data was used for the validation of the machine learning algorithms. The target variable chosen for the prediction of droughts was SPI. Two classes for the prediction were identified in the following way: $SPI1 < -1$ was defined as an event and was initialized with 1, $SPI1 > -1$ was initialized with 0 and corresponded to a non-drought event. The lead time of one month was chosen for the prediction After the feature selection 27 variables originating directly from the CRCM5-LE dataset were used as input, each of them in a timeseries of 12 months. In addition to those the teleconnection indices NAO, SCA, EA, EA/WR, AMO and AMO10 were used as input. For this analysis we used a supervised machine learning algorithm, an Artificial Neural Network (ANN). ANNs are algorithms whose design is inspired by the architecture of the human brain with its neurons (Russell and Norvig, 2009).; they both consists of connected nodes. A link between the node i and the node j serves to propagate the activation $a_i$ from i to j. To each connection a numeric weight $w_{i,j}$ is assigned. The output of the node is computed by:

$$a_i = g(in_j) = g\left(\sum_{i=0}^{n} w_{i,j} a_i\right) \qquad (4)$$

(Russell and Norvig, 2009, p. 728). The activation function defines the output of the node. In order to have stable learners with confident predictions a function with a soft threshold is recommended (Russell and Norvig, 2009). In this study the following three activation functions were used: Sigmoid, Rectified Linear Unit (ReLU), Exponential Linear Unit (ELU). Sigmoid activation is especially useful for the output layer (Russell and Norvig, 2009), while ReLU and ELU both have the property of allowing very fast optimization (Maas, 2013) . Sigmoid function, also called logistic function, is defined in the following way:

$$Logistic(x) = \frac{1}{1+e^{-x}} \qquad (5)$$

(Russell and Norvig, 2009). This function has an output between 0 and 1. This can be interpreted as a probability of belonging to the class 1. One of the main disadvantages of the sigmoid activation function is the vanishing gradient problem: at higher, almost saturated layers with values of 1 or -1, the gradients become nearly 0 resulting in a slow optimization convergence (Russell and Norvig, 2009, p. 726). ReLU refers to Rectified Linear Unit and shows better performance when dealing with the vanishing gradient problem (Maas, 2013). ReLU is defined in the following way:

$f(x) = max(0,x)$ $\qquad (6)$

ELU refers to the Exponential Linear Unit and was introduced by Clevert et al. (2016). Clevert et al. (2016) claim that in experiments the ELU activation led to faster learning and

significantly better generalization performance than ReLU and sigmoid activation. The
function is defined as:

$$f(x) = \begin{cases} x \text{ if } x > 0 \\ \alpha(exp(x) - 1) \text{ if } x \leq 0 \end{cases} \qquad (7)$$

α controls the value to which an ELU saturates for negative inputs. Per default the value is set
to 1 such that the function saturates at -1. Two kinds of layers were used in this study: Dense
and Dropout. Dense refers to a regular fully connected neural network layer. Dropout refers
to a layer which is randomly setting a fraction of inputs to zero at each update. This technique
is used to prevent overfitting and therefore improving the performance of the algorithm
(Chollet et al., 2015). The first part of the study concentrated on the methodological search
for the best performing algorithms. A pipeline to search for the best performing architecture,
value for L2 regularization and loss function was built up. The model performance was
evaluated using Accuracy and F1-score (Sasaki, 2007). The latter one is especially useful
when training on datasets with an imbalanced class distribution, as it is in the case of our
dataset. Accuracy is defined in the following way:

$$Accuracy = \frac{Number\ of\ right\ predictions}{Total\ number\ of\ samples} \qquad (8)$$

F1-score is a harmonic measure between precision and recall. Precision is the amount of true positives with respect to the
amount of positively classified data. Recall is the amount of true positives with respect to the total number of positives in the
data. F1-score is defined in the following way:

$$F1 - score = 2\frac{Precision \cdot Recall}{Precision + Recall} \qquad (9)$$

Due to the class imbalance within the dataset we require that the accuracy on each class is at
least 50%. In that case given the distribution of the test dataset of 1803 non-drought events to

387 droughts for Lisbon and 1848 non-drought events to 352 drought events for Munich a marginal F1-score of 0.26 for Lisbon and 0.24 for Munich is given. The second part of the study analyzed the best performing algorithms (one for Lisbon domain, one for Munich domain) by applying explainable AI methods. SHAP (SHapley Additive exPlanations) is a state of the art method for interpretation of machine learning models, which was inspired by game theory (Lundberg and Lee, 2017). It estimates for each input feature an average marginal contribution to the prediction of the result and therefore allows a comparison of the contributions among different features. In addition to that the difference in predictability among the seasons is calculated and compared to gain a better understanding on the influence of seasonal weather patterns.

# Chapter-3 Dataset

The US drought monitor is a measure of drought across the US manually created by experts using a wide range of data.

This datasets' aim is to help investigate if droughts could be predicted using only meteorological data, potentially leading to generalization of US predictions to other areas of the world.
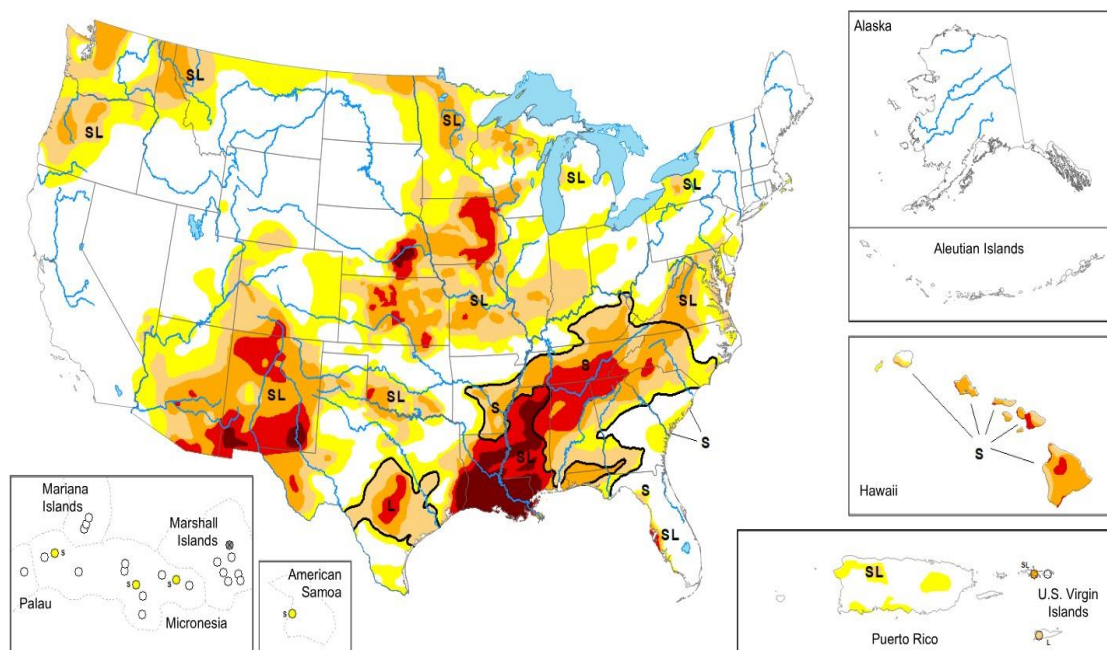


Fig. 3.1- Map of 28<sup>th</sup> November, 2023 taken from the official website of US Drought Monitor

**Content**

This is a classification dataset over six levels of drought, which is no drought (None in the dataset), and five drought levels shown below.

Each entry is a drought level at a specific point in time in a specific US county, accompanied by the last **90 days of 18 meteorological indicators** shown in the bottom of this description.

| Category | Description | Possible Impacts |
|---|---|---|
| D0 | Abnormally Dry | Going into drought:<br>▪ short-term dryness slowing planting, growth of crops or pastures<br>Coming out of drought:<br>▪ some lingering water deficits<br>▪ pastures or crops not fully recovered |
| D1 | Moderate Drought | ▪ Some damage to crops, pastures<br>▪ Streams, reservoirs, or wells low, some water shortages developing or imminent<br>▪ Voluntary water-use restrictions requested |
| D2 | Severe Drought | ▪ Crop or pasture losses likely<br>▪ Water shortages common<br>▪ Water restrictions imposed |
| D3 | Extreme Drought | ▪ Major crop/pasture losses<br>▪ Widespread water shortages or restrictions |
| D4 | Exceptional Drought | ▪ Exceptional and widespread crop/pasture losses<br>▪ Shortages of water in reservoirs, streams, and wells creating water emergencies |

To avoid data leakage, the data has been split into the following subsets.

| Split | Year Range (inclusive) | Percentage (approximate) |
|-------|------------------------|--------------------------|
| Train | 2000-2009 | 47% |
| Validation | 2010-2011 | 10% |
| Test | 2012-2020 | 43% |

## Dataset Imbalance

The dataset is imbalanced, as can be seen in the following graph.



Fig. 3.2- Dataset Imbalance

**On NaN values**: The drought scores are available weekly while the meteorological data points are available daily. To make using previous drought scores for prediction easier (e.g. by interpolating), I merged them into one file and set the drought scores to NaN were not available.

## Acknowledgements

# Meteorological Indicators

| Indicator | Description |
| --- | --- |
| WS10M_MIN | Minimum Wind Speed at 10 Meters (m/s) |
| QV2M | Specific Humidity at 2 Meters (g/kg) |
| T2M_RANGE | Temperature Range at 2 Meters (C) |
| WS10M | Wind Speed at 10 Meters (m/s) |
| T2M | Temperature at 2 Meters (C) |
| WS50M_MIN | Minimum Wind Speed at 50 Meters (m/s) |
| T2M_MAX | Maximum Temperature at 2 Meters (C) |
| WS50M | Wind Speed at 50 Meters (m/s) |
| TS | Earth Skin Temperature (C) |

| Indicator | Description |
| --- | --- |
| WS50M_RANGE | Wind Speed Range at 50 Meters (m/s) |
| WS50M_MAX | Maximum Wind Speed at 50 Meters (m/s) |
| WS10M_MAX | Maximum Wind Speed at 10 Meters (m/s) |
| WS10M_RANGE | Wind Speed Range at 10 Meters (m/s) |
| PS | Surface Pressure (kPa) |
| T2MDEW | Dew/Frost Point at 2 Meters (C) |
| T2M_MIN | Minimum Temperature at 2 Meters (C) |
| T2MWET | Wet Bulb Temperature at 2 Meters (C) |
| PRECTOT | Precipitation (mm day-1) |

# Chapter-4 Data Process Lifecycle

The data process lifecycle is a systematic approach to handling and analyzing data, encompassing several key steps that ensure the extraction of meaningful insights and the development of robust models. Each stage plays a crucial role in the overall process, contributing to the success of data-driven decision-making. Here's an overview of the key steps in the data process lifecycle:



Fig. 4 - Dataset Process Lifecycle

1. **Data Collection:**

   - This is the initial phase where raw data is gathered from various sources. Sources can include databases, APIs, sensors, web scraping, or any other means depending on the nature of the project.

   - Data collection is critical as the quality and relevance of the data collected directly impact the outcomes of subsequent stages in the process.

2. **Data Wrangling:**

- Once the data is collected, it often needs to be cleaned and transformed to address issues such as missing values, outliers, or inconsistencies.

- Data wrangling involves tasks like handling missing data, converting data types, and dealing with duplicates. The goal is to prepare the data for further analysis by ensuring it is accurate and in a suitable format.

3. **Exploratory Data Analysis (EDA):**

- EDA is the process of visually and statistically exploring the data to gain insights into its underlying patterns, relationships, and distributions.

- Descriptive statistics, data visualization techniques, and other exploratory techniques are used to summarize and understand the main characteristics of the dataset.

4. **Feature Selection and Extraction:**

- In this step, features that are most relevant to the problem at hand are selected or engineered.

- Feature selection involves choosing a subset of the most important features, while feature extraction involves creating new features from the existing ones. The goal is to reduce dimensionality and improve the model's performance.

5. **Model Development:**

- This step involves choosing and building a suitable machine learning model based on the problem and the characteristics of the data.

- The model is trained using a portion of the dataset, and its parameters are adjusted to optimize performance. This step may involve the use of various algorithms depending on the nature of the task, such as regression, classification, or clustering.

6. **Model Evaluation:**

- After the model is developed, it needs to be evaluated to ensure its performance meets the desired criteria. This is done using a separate set of data that the model has not seen during training.

- Evaluation metrics depend on the type of problem; for example, accuracy, precision, recall, F1-score for classification problems, or mean squared error for regression problems.

- Iterative refinement of the model may be necessary based on the evaluation results.

Throughout the data process lifecycle, it's important to note that these steps are often iterative. Feedback from one stage may necessitate revisiting previous stages to make adjustments. This cyclic nature ensures a thorough and effective approach to handling and extracting value from data.

# Chapter-5 Initial Implementation

## 5.1 Importing Libraries

Python code snippet imports necessary libraries for data analysis, visualization, and machine learning. Below is a description of each import statement:

1. import os: This module provides a way to interact with the operating system, allowing the script to perform operations such as reading or writing to the file system.

2. import subprocess: The subprocess module allows the script to spawn new processes, connect to their input/output/error pipes, and obtain their return codes.

3. import pandas as pd: The code utilizes the pandas library, renaming it as 'pd'. Pandas is widely used for data manipulation and analysis, offering data structures like dataframes that simplify working with structured data.

4. import numpy as np: This imports the NumPy library, an essential package for numerical computing in Python. It provides support for large, multi-dimensional arrays and matrices, along with mathematical functions to operate on these data structures.

5. import matplotlib.pyplot as plt: This line imports the pyplot module from the Matplotlib library, which is commonly used for creating static, interactive, and animated visualizations in Python.

6. %matplotlib inline: This is a Jupyter Notebook magic command that ensures Matplotlib plots are displayed directly below the code cell in the notebook.

7. import warnings: The warnings module allows the script to control warning messages issued by other modules and libraries.

8. warnings.filterwarnings('ignore'): This line suppresses all warning messages, making the output cleaner by ignoring non-critical warnings.

9. import seaborn as sns: Seaborn is a statistical data visualization library based on Matplotlib. It provides an aesthetically pleasing and informative representation of data.

10. from sklearn.model_selection import train_test_split: This code imports the train_test_split function from scikit-learn, a popular machine learning library. It is used for splitting datasets into training and testing sets.

11. from sklearn.preprocessing import StandardScaler: The StandardScaler class from scikit-learn is imported for standardizing feature data by removing the mean and scaling to unit variance.

12. from sklearn.feature_selection import RFE: Recursive Feature Elimination (RFE) is a feature selection method, and this import brings it into the script. It is used to select the most important features for a machine learning model.

13. from sklearn.ensemble import RandomForestClassifier: The RandomForestClassifier class from scikit-learn is imported, enabling the use of a Random Forest algorithm for classification tasks.

## 5.2 Reading the input data

```
drought_df = pd.read_csv('../input/us-drought-meteorological-data/train_timeseries
/train_timeseries.csv')
drought_df.head()
```

| | fips | date | PRECTOT | PS | QV2M | T2M | T2MDEW | T2MWET | T2M_MAX | T2M_MIN | ... | TS | WS10M | WS10M_MAX | WS10M_MIN | WS10M_RANGE |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1001 | 2000-01-01 | 0.22 | 100.51 | 9.65 | 14.74 | 13.51 | 13.51 | 20.96 | 11.46 | ... | 14.65 | 2.20 | 2.94 | 1.49 | 1.46 |
| 1 | 1001 | 2000-01-02 | 0.20 | 100.55 | 10.42 | 16.69 | 14.71 | 14.71 | 22.80 | 12.61 | ... | 16.60 | 2.52 | 3.43 | 1.83 | 1.60 |
| 2 | 1001 | 2000-01-03 | 3.65 | 100.15 | 11.76 | 18.49 | 16.52 | 16.52 | 22.73 | 15.32 | ... | 18.41 | 4.03 | 5.33 | 2.66 | 2.67 |
| 3 | 1001 | 2000-01-04 | 15.95 | 100.29 | 6.42 | 11.40 | 6.09 | 6.10 | 18.09 | 2.16 | ... | 11.31 | 3.84 | 5.67 | 2.08 | 3.59 |
| 4 | 1001 | 2000-01-05 | 0.00 | 101.15 | 2.95 | 3.86 | -3.29 | -3.20 | 10.82 | -2.66 | ... | 2.65 | 1.60 | 2.50 | 0.52 | 1.98 |

5 rows × 21 columns

Fig. 5.1 – Input Data

## 5.3 Initial exploration and data cleaning (Data Wrangling)

Initial exploration and data cleaning, often referred to as data wrangling, are critical steps in the machine learning pipeline. These steps involve preparing the raw data for analysis and modeling by addressing issues such as missing values, outliers, inconsistencies, and other data quality issues. Here's a detailed explanation:

1. **Understanding the Data**:

   - Before diving into data cleaning, it's essential to understand the dataset thoroughly. This includes:

      - Understanding the meaning and context of each variable.

      - Identifying the types of variables (e.g., numerical, categorical).

      - Examining summary statistics, distributions, and visualizations to gain insights into the data's characteristics.

2. **Handling Missing Values**:

   - Missing values are a common issue in real-world datasets and can adversely affect the performance of machine learning models. Strategies for handling missing values include:

     - Imputation: Replace missing values with a substitute value, such as the mean, median, mode, or a value predicted by a model.

     - Deletion: Remove observations or variables with a high percentage of missing values if they cannot be reasonably imputed. However, this should be done cautiously to avoid losing valuable information.

     - Advanced Techniques: Use advanced imputation techniques such as k-nearest neighbors (KNN) imputation or multiple imputation to preserve the underlying structure of the data.

3. **Dealing with Outliers**:

   - Outliers can skew the distribution of data and affect the performance of machine learning models. Strategies for handling outliers include:

     - Identification: Use statistical methods or visualizations to detect outliers, such as Z-score, IQR, or box plots.

     - Treatment: Depending on the nature of the outliers and the specific requirements of the problem, outliers can be removed, transformed, or Winsorized (capped at a certain percentile).

     - Domain Knowledge: Consider the domain-specific context to determine whether outliers are genuine data points or errors that need to be addressed.

4. **Handling Categorical Variables**:

   - Categorical variables represent discrete categories or groups and may require encoding before they can be used in machine learning models. Common approaches include:

    - One-Hot Encoding: Convert categorical variables into binary vectors, where each category becomes a binary feature.

    - Label Encoding: Map categorical values to integer labels.

    - Target Encoding: Encode categorical variables based on the target variable's mean or frequency within each category.

    - Ordinal Encoding: Encode ordinal variables with a meaningful order, such as low, medium, high.

5. **Feature Engineering**:

   - Feature engineering involves creating new features or transforming existing ones to improve model performance. This may include:

    - Creating interaction terms or polynomial features.

    - Transforming numerical variables to better approximate a normal distribution.

    - Extracting information from date/time variables, such as day of the week or month.

    - Dimensionality reduction techniques like principal component analysis (PCA) or feature selection algorithms to reduce the number of features.

6. **Data Scaling and Normalization**:

   - Scaling and normalization are essential preprocessing steps for many machine learning algorithms, particularly those based on distance metrics or gradient descent optimization. Common techniques include:

     - Min-Max Scaling: Rescale features to a specific range, often between 0 and 1.

     - Standardization: Transform features to have a mean of 0 and a standard deviation of 1.

     - Robust Scaling: Scale features using median and interquartile range to mitigate the impact of outliers.

7. **Validation**:

   - After data cleaning and preprocessing, it's crucial to validate the quality of the processed data. This can involve:

     - Splitting the data into training, validation, and test sets.

     - Assessing the distribution and characteristics of the cleaned data.

     - Monitoring performance metrics during model training to ensure that data preprocessing has improved model performance.

In summary, initial exploration and data cleaning (data wrangling) are essential steps in the machine learning pipeline to ensure that the data is clean, consistent, and suitable for modeling. These steps involve understanding the data, handling missing values and outliers, encoding categorical variables, performing feature engineering, scaling and normalizing features, and validating the quality of the processed data before model training.

```
drought_df.info()
<class 'pandas.core.frame.DataFrame'>
```

```
RangeIndex: 19300680 entries, 0 to 19300679
Data columns (total 21 columns):
 #   Column       Dtype
---  ------       -----
 0   fips         int64
 1   date         object
 2   PRECTOT      float64
 3   PS           float64
 4   QV2M         float64
 5   T2M          float64
 6   T2MDEW       float64
 7   T2MWET       float64
 8   T2M_MAX      float64
 9   T2M_MIN      float64
 10  T2M_RANGE    float64
 11  TS           float64
 12  WS10M        float64
 13  WS10M_MAX    float64
 14  WS10M_MIN    float64
 15  WS10M_RANGE  float64
 16  WS50M        float64
 17  WS50M_MAX    float64
 18  WS50M_MIN    float64
 19  WS50M_RANGE  float64
 20  score        float64
dtypes: float64(19), int64(1), object(1)
memory usage: 3.0+ GB
```

`drought_df.isnull().sum()`

```
fips                    0
date                    0
PRECTOT                 0
PS                      0
QV2M                    0
T2M                     0
T2MDEW                  0
T2MWET                  0
T2M_MAX                 0
T2M_MIN                 0
T2M_RANGE               0
TS                      0
WS10M                   0
WS10M_MAX               0
WS10M_MIN               0
WS10M_RANGE             0
WS50M                   0
WS50M_MAX               0
WS50M_MIN               0
WS50M_RANGE             0
score            16543884
dtype: int64
```

## 5.4 Missing Value Treatment

Missing value treatment, also known as handling missing data, is an essential step in preparing data for machine learning models. Missing values are common in real-world datasets and can adversely affect the performance and accuracy of machine learning algorithms. Missing value treatment involves strategies for dealing with these missing values in a way that preserves the integrity and usefulness of the data. Here's an explanation of missing value treatment in machine learning:

1. Identifying Missing Values:

   - The first step in missing value treatment is to identify where missing values exist within the dataset. Missing values may be represented by various symbols such as NaN (Not a Number), NA (Not Available), NULL, or blanks.

2. Understanding the Causes of Missing Values:

   - It's essential to understand why data is missing, as the appropriate treatment strategy may depend on the underlying reasons. Missing values can occur due to various factors, including:

     - Data Entry Errors: Human errors during data collection or data entry processes.

     - Missing at Random (MAR): The probability of data being missing depends on other observed variables in the dataset.

     - Missing Completely at Random (MCAR): The probability of data being missing is unrelated to any other variables in the dataset.

     - Missing Not at Random (MNAR): The probability of data being missing depends on unobserved or latent variables.

3. Handling Missing Values:

   - There are several strategies for handling missing values in machine learning datasets:

   - Deletion: Delete rows or columns with missing values. This approach is straightforward but can lead to a loss of valuable information.

   - Row Deletion (Listwise Deletion): Remove entire rows with missing values.

   - Column Deletion: Remove entire columns (features) with a high percentage of missing values.

   - Imputation: Replace missing values with substitute values. Common imputation techniques include:

   - Mean/Median/Mode Imputation: Replace missing values with the mean, median, or mode of the feature.

   - Forward Fill/Backward Fill: Propagate the last observed value forward or the next observed value backward.

   - Linear Interpolation: Estimate missing values based on the values of neighboring data points.

   - K-Nearest Neighbors (KNN) Imputation: Predict missing values based on the values of similar instances using the KNN algorithm.

   - Multiple Imputation: Generate multiple imputed datasets, where missing values are imputed multiple times using statistical models, and combine the results.

   - Prediction Models: Use machine learning algorithms to predict missing values based on the values of other features in the dataset.

- Domain Knowledge: Utilize domain knowledge to infer or estimate missing values based on the context of the problem and the relationships between variables.

4. Choosing the Right Strategy:

- The choice of missing value treatment strategy depends on factors such as the extent of missingness, the underlying causes of missing values, the distribution of the data, and the requirements of the machine learning task. It may involve a combination of approaches or iterative experimentation to determine the most effective strategy.

5. Validation:

- After applying missing value treatment, it's crucial to validate the quality of the processed data. This can involve assessing the distribution and characteristics of the data, monitoring the impact on model performance, and comparing different treatment strategies to determine the most suitable approach.

In summary, missing value treatment in machine learning involves identifying missing values, understanding their causes, and applying appropriate strategies such as deletion, imputation, prediction models, or domain knowledge to handle them effectively. The choice of strategy depends on factors such as the extent and nature of missingness, the distribution of the data, and the requirements of the machine learning task. Validation is essential to ensure that the treatment process does not introduce bias or compromise the integrity of the data.

# Removing the null values in the target variable as the drought score is only available for once in 7 days.

```
drought_df = drought_df.dropna()
drought_df.isnull().sum()
fips          0
```

```
date          0
PRECTOT       0
PS            0
QV2M          0
T2M           0
T2MDEW        0
T2MWET        0
T2M_MAX       0
T2M_MIN       0
T2M_RANGE     0
TS            0
WS10M         0
WS10M_MAX     0
WS10M_MIN     0
WS10M_RANGE   0
WS50M         0
WS50M_MAX     0
WS50M_MIN     0
WS50M_RANGE   0
score         0
dtype: int64
```

## 5.5 Reformatting the data

```
drought_df.dtypes
fips            int64
date            object
PRECTOT         float64
PS              float64
QV2M            float64
T2M             float64
T2MDEW          float64
T2MWET          float64
T2M_MAX         float64
T2M_MIN         float64
T2M_RANGE       float64
TS              float64
WS10M           float64
WS10M_MAX       float64
WS10M_MIN       float64
WS10M_RANGE     float64
WS50M           float64
WS50M_MAX       float64
WS50M_MIN       float64
WS50M_RANGE     float64
score           float64
dtype: object
```

This code snippet augments a DataFrame named **drought_df**. It introduces new columns ('ye

ar', 'month', 'day') by extracting date components from an existing 'date' column. Additionally

, it refines the 'score' column by rounding its values to integers. The final line prints the data t

ypes of each column in the DataFrame after these adjustments.

```
fips            int64
date           object
PRECTOT       float64
PS            float64
QV2M          float64
T2M           float64
T2MDEW        float64
T2MWET        float64
T2M_MAX       float64
T2M_MIN       float64
T2M_RANGE     float64
TS            float64
WS10M         float64
WS10M_MAX     float64
WS10M_MIN     float64
WS10M_RANGE   float64
WS50M         float64
WS50M_MAX     float64
WS50M_MIN     float64
WS50M_RANGE   float64
score           int64
year            int64
month           int64
day             int64
dtype: object
```

```
drought_df['fips'].nunique()
3108
drought_df['score'].round().value_counts()

0    1652230
1     466944
2     295331
3     196802
4     106265
5      39224
Name: score, dtype: int64
```

## 5.6 Exploratory Data Analysis

## Univariate Analysis - Descriptive statistics

```
# Descriptive statistics
```
a statistical summary of a DataFrame called drought_df. It uses the describe() method to generate descriptive statistics for numeric columns and categorical columns separately. The display function is employed to show these summaries in a visually appealing format.

The subsequent lines calculate and print the skewness and kurtosis for each numeric column in the DataFrame. Skewness measures the asymmetry of the data distribution, indicating whether the data is skewed to the left or right. Kurtosis measures the tail heaviness of the distribution, signifying how much data is in the tails and how sharp or flat the peak is.

|       | fips        | PRECTOT      | PS           | QV2M         | T2M          | T2MDEW        | T2MWET        | T2M_MAX      | T2M_MIN       | T2M_RAI    |
|-------|-------------|--------------|--------------|--------------|--------------|---------------|---------------|--------------|---------------|------------|
| count | 2.756796e+06 | 2.756796e+06 | 2.756796e+06 | 2.756796e+06 | 2.756796e+06 | 2.756796e+06  | 2.756796e+06  | 2.756796e+06 | 2.756796e+06  | 2.75679    |
| mean  | 3.067038e+04 | 2.714566e+00 | 9.664736e+01 | 7.875770e+00 | 1.289923e+01 | 7.049350e+00  | 7.084938e+00  | 1.876711e+01 | 7.519916e+00  | 1.12472    |
| std   | 1.497911e+04 | 6.247590e+00 | 5.444698e+00 | 4.721459e+00 | 1.097040e+01 | 1.019765e+01  | 1.014364e+01  | 1.160295e+01 | 1.061818e+01  | 4.03802    |
| min   | 1.001000e+03 | 0.000000e+00 | 6.649000e+01 | 1.400000e-01 | -3.544000e+01 | -3.544000e+01 | -3.546000e+01 | -3.003000e+01 | -4.085000e+01 | 1.60000    |
| 25%   | 1.904450e+04 | 0.000000e+00 | 9.583000e+01 | 3.810000e+00 | 4.580000e+00 | -8.800000e-01 | -8.400000e-01 | 1.036000e+01 | -5.700000e-01 | 8.37000    |
| 50%   | 2.921200e+04 | 1.900000e-01 | 9.828000e+01 | 6.940000e+00 | 1.421000e+01 | 7.810000e+00  | 7.810000e+00  | 2.062000e+01 | 8.260000e+00  | 1.12000    |
| 75%   | 4.600750e+04 | 2.260000e+00 | 9.994000e+01 | 1.145000e+01 | 2.200000e+01 | 1.567000e+01  | 1.567000e+01  | 2.797000e+01 | 1.628000e+01  | 1.40800    |
| max   | 5.604300e+04 | 1.686900e+02 | 1.037600e+02 | 2.212000e+01 | 3.933000e+01 | 2.687000e+01  | 2.687000e+01  | 4.775000e+01 | 3.228000e+01  | 3.01700    |

8 rows × 23 columns

|        | date       |
|--------|------------|
| count  | 2756796    |
| unique | 887        |
| top    | 2000-01-04 |
| freq   | 3108       |

Fig. 5.2- Univariate Analysis

## Univariate Analysis - Distribution of continuous variables

Univariate analysis, specifically in the context of machine learning, involves the examination of individual features (variables) in a dataset one at a time. When focusing on the distribution of continuous variables, univariate analysis entails exploring the characteristics and patterns of each continuous variable independently. Here's a detailed explanation:

1. **Data Understanding**:

   - Before performing univariate analysis, it's essential to understand the dataset and the variables it contains. Continuous variables are those that can take an infinite number of values within a certain range, such as age, income, temperature, etc.

2. **Summary Statistics**:

   - The first step in univariate analysis is to calculate summary statistics for each continuous variable. These statistics provide a basic understanding of the distribution of the data and include measures such as mean, median, mode, range, variance, standard deviation, skewness, and kurtosis.

   - Mean and median indicate the central tendency of the data, while variance and standard deviation measure the spread or dispersion. Skewness and kurtosis provide insights into the symmetry and shape of the distribution.

3. **Visualization**:

   - Visualizing the distribution of continuous variables is crucial for gaining deeper insights and identifying patterns or anomalies. Common visualization techniques include:

     - Histograms: A histogram displays the frequency distribution of values in a continuous variable by dividing the data into bins and plotting the number of observations in each bin.

     - Density Plots: Density plots provide a smoothed representation of the distribution of data, similar to histograms but with a continuous curve.

     - Box Plots: Box plots illustrate the distribution of data using quartiles, median, and outliers, making it easy to identify the central tendency and spread of the data, as well as any outliers.

- Violin Plots: Violin plots combine the features of box plots and density plots, providing a summary of the distribution along with a kernel density estimation.

4. **Interpretation**:

  - Analyzing the summary statistics and visualizations obtained from univariate analysis helps in understanding the characteristics of each continuous variable:

    - Central Tendency: Determine whether the data is centered around a particular value (mean, median) or if there are multiple modes.

    - Spread: Assess the variability or dispersion of the data using measures like variance and standard deviation.

    - Skewness and Kurtosis: Examine the symmetry and shape of the distribution. Positive skewness indicates a right-skewed distribution, while negative skewness indicates a left-skewed distribution. Kurtosis measures the tails or peakedness of the distribution.

    - Outliers: Identify any extreme values that deviate significantly from the rest of the data, which may require further investigation or outlier treatment.

5. **Feature Engineering**:

  - Univariate analysis can guide feature engineering decisions by highlighting variables with meaningful patterns or distributions that are relevant to the machine learning task at hand. For example, transformations such as log transformation or normalization may be applied to address skewness or heteroscedasticity.

6. **Validation**:

- It's essential to validate the findings of univariate analysis through cross-validation or by assessing the impact of individual variables on the model's performance. This helps in ensuring that the insights obtained from univariate analysis are relevant and beneficial for the machine learning task.

In summary, univariate analysis of continuous variables in machine learning involves examining the distribution, summary statistics, and visualizations of individual features to gain insights into their characteristics and patterns. It serves as a fundamental step in data exploration and feature understanding, guiding subsequent data preprocessing and modeling efforts.

1. **measures_column_list**: A list containing the names of meteorological measures such as precipitation (PRECTOT), surface pressure (PS), specific humidity at 2 meters (QV2M), temperature at 2 meters (T2M), and various wind speed measures (WS10M, WS50M, etc.).

2. **drought_df_measures**: A new DataFrame created by selecting only the columns specified in **measures_column_list** from the original **drought_df**. This DataFrame now includes only the meteorological measures of interest.

3. **for col_name in measures_column_list:**: A loop iterates through each meteorological measure in the list.

4. **plt.figure():** Initiates a new figure for each measure.

5. **plt.hist(drought_df_measures[col_name], density=True):** Generates a histogram for the values of the current meteorological measure. The **density=True** parameter normalizes the histogram to represent a probability density.

6. **x_name = col_name:** Assigns the current measure's name to the x-axis label.

7. **plt.xlabel(x_name):** Sets the x-axis label using the measure's name.

8. **y_name = 'Density':** Assigns the y-axis label.

9. **plt.ylabel(y_name):** Sets the y-axis label.

10. **plt.title('Distribution of {x_name}'.format(x_name=x_name)):** Sets the title of the histogram, indicating the distribution of the current meteorological measure.
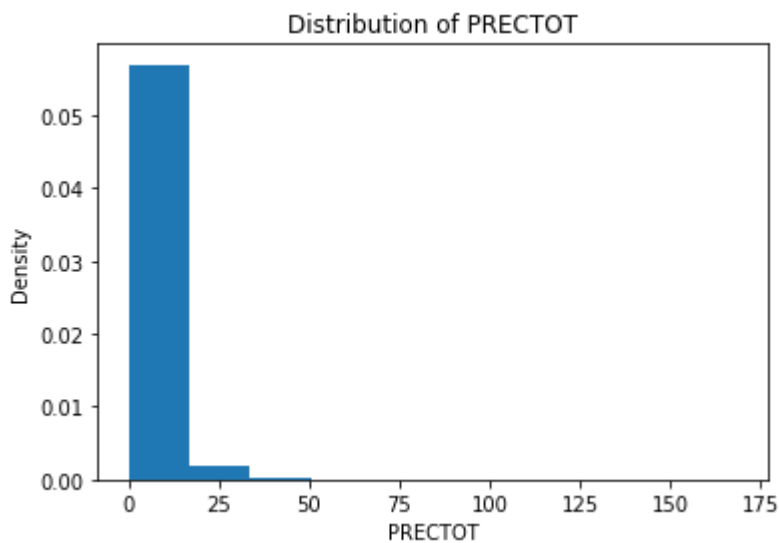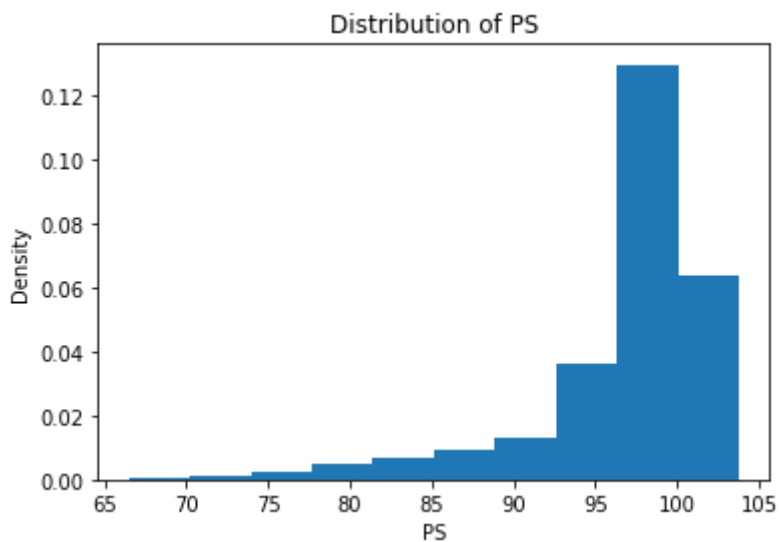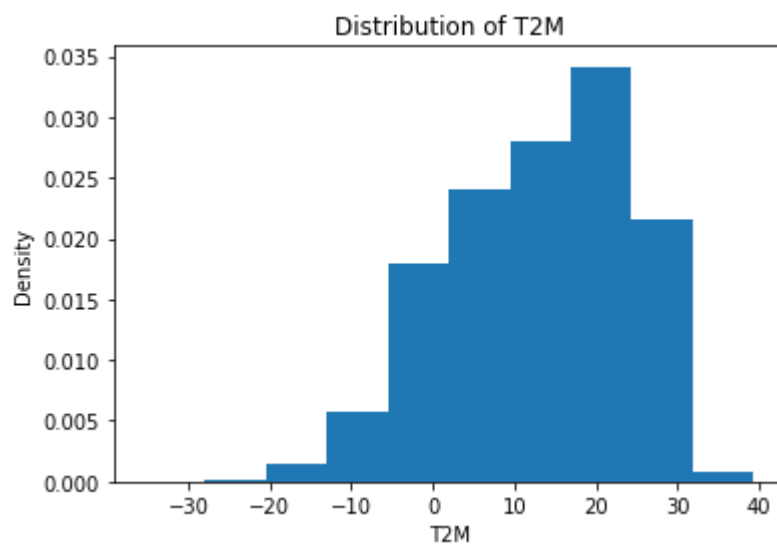


Fig. 5.3- Distribution of PRECTOT

Fig. 5.4-Distribution of PS, QV2M and T2M

Distribution of T2MWET
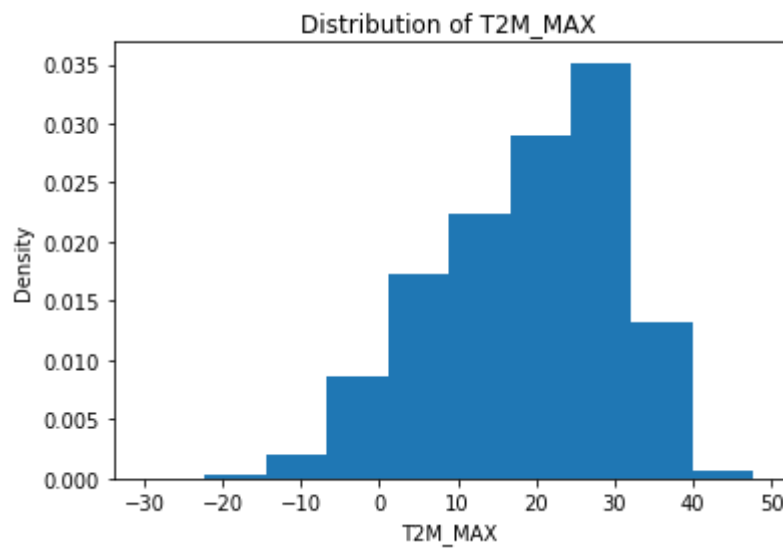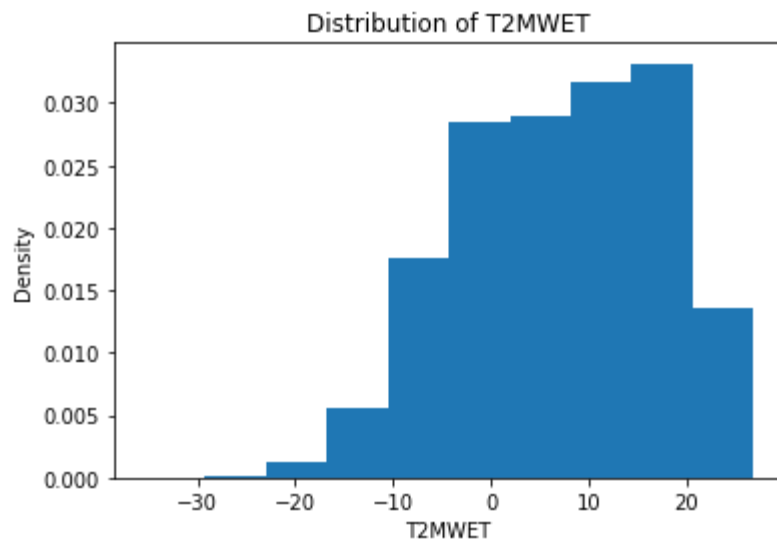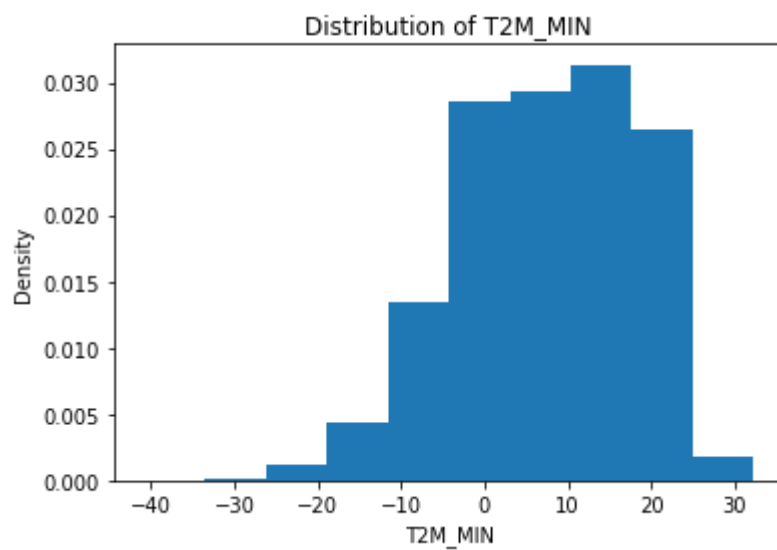

Distribution of T2M_MAX

Fig. 5.5-Distribution of T2MDEW, T2MWET, T2M_MAX


Distribution of T2M_MIN

Fig. 5.6-Distribution of T2M_MIN, T2M_RANGE, TS

Distribution of WS10M_MAX


Distribution of WS10M_MIN

Fig. 5.7-Distribution of WS10M, WS10M_MAX, WS10M_MIN


Distribution of WS10M_RANGE

Fig. 5.8-Distribution of WS10M_RANGE, WS50M, WS50M_MAX

Fig. 5.9-Distribution of WS50M_MIN, WS50M_RANGE

## 5.7 Outlier Treatment

Outliers are data points that significantly differ from the rest of the observations in a dataset. They can occur due to various reasons such as measurement errors, experimental errors, or rare events. In machine learning, outliers can adversely affect the performance and accuracy of models by skewing the training process and influencing the learned relationships between features and the target variable.

Outlier treatment refers to the process of identifying and handling these outliers in a dataset. Here's a detailed explanation of outlier treatment in machine learning:

1. **Identification of Outliers**:

   - Visual Inspection: One common method for identifying outliers is through visual inspection of the data using statistical plots such as scatter plots, box plots, or histograms. Data points that are significantly distant from the main cluster of points may be considered outliers.

- Statistical Methods: Statistical techniques like Z-score, modified Z-score, or interquartile range (IQR) can be employed to identify outliers. For example, data points with a Z-score greater than a certain threshold (usually 2 or 3) are considered outliers.

2. **Handling Outliers**:

- Removal: One approach to dealing with outliers is to simply remove them from the dataset. However, this should be done cautiously, as removing too many outliers may result in loss of valuable information and potentially bias the model.

- Transformation: Data transformation techniques such as log transformation, square root transformation, or Box-Cox transformation can sometimes be effective in mitigating the impact of outliers.

- Imputation: Outliers can be replaced with a more representative value such as the mean, median, or mode of the feature. This approach helps in retaining the overall structure of the data while minimizing the influence of outliers.

- Binning: Binning involves dividing the data into intervals or bins and replacing outliers with the bin boundaries. This method can be useful when dealing with continuous numerical data.

- Winsorization: Winsorization involves capping the extreme values of the data at a certain percentile (e.g., 1st and 99th percentiles). This approach helps in reducing the impact of outliers without completely removing them.

- Model-based Approaches: Some machine learning algorithms are robust to outliers or can handle them implicitly. Using such algorithms, like decision trees or random forests, may alleviate the need for explicit outlier treatment.

3. **Validation**:

   - After handling outliers, it's crucial to validate the effectiveness of the chosen approach. This can be done by evaluating the model's performance metrics before and after outlier treatment using techniques such as cross-validation or holdout validation.

4. **Domain Knowledge**:

   - Domain knowledge plays a vital role in outlier treatment. Understanding the context of the data and the underlying processes that generate outliers can help in making informed decisions about how to handle them effectively.

5. **Iterative Process**:

   - Outlier treatment is often an iterative process, where different techniques are applied, and the results are evaluated iteratively until satisfactory results are obtained. It may involve trying different approaches and comparing their impact on the model's performance.

In summary, outlier treatment in machine learning involves identifying and handling data points that deviate significantly from the rest of the dataset. It requires a combination of statistical techniques, domain knowledge, and careful validation to ensure that the outlier treatment does not adversely affect the model's performance.

## Identifying Outliers

Identifying outliers in the meteorological measures dataset (drought_df_measures). It utilizes a series of boxplots to visually represent the distribution of values for each meteorological measure and highlight potential outliers. Here's a breakdown of the code:

   1. plt.figure(figsize=(10,40)): Initiates a new figure with a specified size (10 inches in width and 40 inches in height) to accommodate multiple subplots.

2. for x in (range(1,19)):: Initiates a loop to iterate through the first 18 meteorological measures in the dataset.

3. plt.subplot(19,1,x): Creates a subplot for each meteorological measure, arranging them vertically in a single column.

4. sns.boxplot(x=drought_df_measures.columns[x-1], data=drought_df_measures): Generates a boxplot for the current meteorological measure, highlighting the distribution of values and revealing any potential outliers. The x-axis is labeled with the name of the meteorological measure.

5. x_name = drought_df_measures.columns[x-1]: Assigns the name of the current meteorological measure to the variable x_name.

6. plt.title(f'Distribution of {x_name}'): Sets the title of each subplot, indicating the distribution of the values for the current meteorological measure.

7. plt.tight_layout(): Adjusts the layout to prevent overlapping subplots and enhance overall visualization.

Fig. 5.10- Outliers - 1

Fig. 5.11- Outliers - 2

Fig. 5.12- Outliers - 3

The code calculates and prints the total number of rows in the DataFrame

**drought_df_measures** and then iterates through numeric columns, identifying and

printig the count of values that fall beyond the standard outlier limit (defined as three

times the standard deviation from the mean) for each column. This provides insight

into the presence of potential outliers in the dataset for each meteorological measure,

aiding in the identification and understanding of data points that may significantly

deviate from the norm.

OUTOUT:

```
Total rows =  2756796
Number of values beyond standard outlier limit in  PRECTOT
65933
Number of values beyond standard outlier limit in  PS
73197
Number of values beyond standard outlier limit in  QV2M
1
Number of values beyond standard outlier limit in  T2M
4531
Number of values beyond standard outlier limit in  T2MDEW
2023
Number of values beyond standard outlier limit in  T2MWET
1814
Number of values beyond standard outlier limit in  T2M_MAX
3384
Number of values beyond standard outlier limit in  T2M_MIN
```

```
6944
Number of values beyond standard outlier limit in  T2M_RANGE
3628
Number of values beyond standard outlier limit in  TS
4762
Number of values beyond standard outlier limit in  WS10M
29954
Number of values beyond standard outlier limit in  WS10M_MAX
23387
Number of values beyond standard outlier limit in  WS10M_MIN
39901
Number of values beyond standard outlier limit in  WS10M_RANGE
35979
Number of values beyond standard outlier limit in  WS50M
23090
Number of values beyond standard outlier limit in  WS50M_MAX
25985
Number of values beyond standard outlier limit in  WS50M_MIN
19569
Number of values beyond standard outlier limit in  WS50M_RANGE
33808
```

## Removing values beyond the standard outlier limit

This code aims to address outliers in the meteorological dataset represented by the

DataFrame **drought_df**. For each meteorological measure, the code filters out values that fall

beyond the standard outlier limit, defined as three times the standard deviation from the mean

for that specific measure. The snippet iterates through each column, applying these outlier

removal conditions to cleanse the data.

Here's a more detailed breakdown:

1. **Iteration Through Columns:** The code uses a series of statements for each

   meteorological measure, ensuring that outlier removal is performed individually for

   measures such as precipitation (**PRECTOT**), surface pressure (**PS**), specific humidity

   at 2 meters (**QV2M**), temperature at 2 meters (**T2M**), wind speed (**WS10M**,

   **WS50M**), and others.

2. **Outlier Removal Conditions:** For each measure, the code constructs conditions to exclude data points that deviate significantly from the mean. These conditions are based on the standard outlier limit, calculated as three times the standard deviation from the mean for each specific meteorological measure.

3. **DataFrame Update:** The DataFrame **drought_df** is successively updated with each measure-specific outlier removal operation, resulting in a refined dataset where extreme values are mitigated for each meteorological variable.

4. **Print Total Rows:** The snippet concludes by printing the total number of rows in the updated DataFrame, indicating the overall impact of the outlier removal process. In this case, the total number of rows is reported as 2,474,338 after the outlier removal.

In summary, this code is designed to enhance the quality and reliability of the meteorological dataset by systematically eliminating outliers for each measure, contributing to a more robust and representative dataset for subsequent analysis or modeling.

## Univariate Analysis - Distribution of categorical variables

This code focuses on analyzing the distribution of categorical variables within a DataFrame named **drought_df**. The snippet selects specific categorical columns ('score', 'year', 'month', 'day') and creates a new DataFrame (**drought_df_categorical**) containing only these columns. The distribution of each categorical variable is visualized using bar plots.

Here's a breakdown of the code:

1. **categorical_column_list:** A list containing the names of categorical columns, including 'score', 'year', 'month', and 'day'.

2. **drought_df_categorical:** A new DataFrame created by selecting only the specified categorical columns from the original DataFrame **drought_df**.

3. **plt.figure(figsize=(10,40)):** Initiates a new figure with a specific size (10 inches in width and 40 inches in height) to accommodate multiple subplots for each categorical variable.

4. **for col_name in categorical_column_list::** Initiates a loop to iterate through each categorical column.

5. **plt.figure():** Initiates a new subplot for each categorical variable.

6. **drought_df_categorical[col_name].value_counts().plot(kind='bar'):** Generates a bar plot illustrating the distribution of unique values for the current categorical variable. The height of each bar represents the frequency of each category.

7. **x_name = col_name:** Assigns the name of the current categorical variable to the x-axis label.

8. **y_name = 'Density':** Assigns the y-axis label as 'Density'.

9. **plt.xlabel(x_name):** Sets the x-axis label using the categorical variable's name.

10. **plt.ylabel(y_name):** Sets the y-axis label.

11. **plt.title('Distribution of {x_name}'.format(x_name=x_name)):** Sets the title of the subplot, indicating the distribution of values for the current categorical variable.

12. **plt.tight_layout():** Adjusts the layout to prevent overlapping subplots and improve overall visualization.

In summary, this code provides a detailed exploration of the distribution of categorical variables in the **drought_df** dataset through individual bar plots for each variable, aiding in the understanding of the distribution patterns within these specific data columns.



Fig. 5.13- Univariate Analysis 1

Fig. 5.14- Univariate Analysis 2

## Bivariate Analysis

A scatter plot to visualize the relationship between two meteorological variables, 'QV2M' (specific humidity at 2 meters) and 'T2M' (temperature at 2 meters), from the DataFrame drought_df. The color of each point in the scatter plot is determined by the corresponding 'score' value. The x-axis is labeled as 'QV2M', the y-axis is labeled as 'T2M', and the plot is titled 'Variation of T2M vs QV2M'. The resulting plot provides insights into the correlation

between specific humidity and temperature, with the color indicating the associated 'score'



values.

Fig. 5.15- Bivariate Analysis of T2M vs QV2M

This code snippet generates a scatter plot to visually explore the relationship between two meteorological variables, 'T2M' (temperature at 2 meters) and 'T2MDEW' (dew point temperature at 2 meters), within the DataFrame **drought_df**. Each point in the scatter plot is colored based on the corresponding 'score' value. The x-axis is labeled as 'T2M', the y-axis is labeled as 'T2MDEW', and the plot is titled 'Variation of T2MDEW vs T2M'. This plot provides insights into the correlation between temperature and dew point temperature, with the color indicating associated 'score'.

Fig. 5.16- Bivariate Analysis of T2MDEW vs T2M

This code snippet first filters the DataFrame **drought_df** to create a new DataFrame,

**temp_df**, containing only the rows where the 'score' is equal to 5. Subsequently, it generates

a scatter plot to illustrate the relationship between two meteorological variables, 'WS10M'

(wind speed at 10 meters) and 'WS50M' (wind speed at 50 meters), from the original

DataFrame. Each point in the scatter plot is colored according to the 'score' values from the

entire dataset. The x-axis is labeled as 'WS10M', the y-axis is labeled as 'WS50M', and the

plot is titled 'Variation of WS50M vs WS10M'. This visualization allows for an examination

of the correlation between wind speeds at different heights, specifically focusing on instances

where the 'score' is equal to 5.

Fig. 5.17- Bivariate Analysis of WS50M vs WS10M

## Extracting Dependent and Independent Variables

```
independent_variables = drought_df.drop('score', 1)
independent_variables = independent_variables.drop('fips', 1)
independent_variables = independent_variables.drop('date', 1)
independent_variables.head()

target = drought_df['score']
target.head()
```

This code seems to be part of a data preprocessing step in a machine learning or data analysis task. Let's break it down step by step:

1. Removing Columns:

   - The DataFrame `drought_df` likely contains data related to drought, and it seems that you're preparing the data for modeling.

   - `drought_df.drop('score', 1)` removes the column named 'score' from the DataFrame. The second argument, `1`, signifies that the operation is column-wise.

   - Similarly, the subsequent lines remove columns named 'fips' and 'date' from the DataFrame.

2. Assigning Independent Variables:

   - After removing the columns 'score', 'fips', and 'date', the resulting DataFrame is
assigned to the variable `independent_variables`. This DataFrame presumably contains the in
dependent variables or features that will be used for modeling.

3. Assigning Target Variable:

   - The column 'score' from the original DataFrame `drought_df` is assigned to the
variable `target`. This column likely represents the target variable, which the model aims to p
redict based on the independent variables.

4. Displaying Data:

   - `.head()` is a method in Pandas DataFrame which returns the first few rows (by
default, the first 5 rows) of the DataFrame. So, `independent_variables.head()` and `
target.head()` are just displaying the first few rows of the DataFrame containing independent
variables and the Series containing the target variable, respectively.

In summary, this code prepares the data for modeling by removing certain columns
from the original DataFrame (`drought_df`), separating the independent variables from
the target variable, and displaying the first few rows of both the independent variables
DataFrame and the target variable Series.

```
3     1
10    2
17    2
24    2
31    1
Name: score, dtype: int64
```

## Correlation between independent variables for Feature Selection

```
correlation_plot = drought_df_measures.corr()
correlation_plot.style.background_gradient(cmap = 'RdYlGn')
```

| | PRECTOT | PS | QV2M | T2M | T2MDEW | T2MWET | T2M_MAX | T2M_MIN | T2M_RANGE | TS | WS10M | WS10M_MAX | W |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PRECTOT | 1.000000 | 0.068775 | 0.245081 | 0.093258 | 0.231035 | 0.230975 | 0.026773 | 0.144929 | -0.304171 | 0.089598 | 0.049730 | 0.060981 | 0. |
| PS | 0.068775 | 1.000000 | 0.282412 | 0.164160 | 0.341234 | 0.341252 | 0.111979 | 0.208285 | -0.225935 | 0.163830 | -0.080747 | -0.135905 | 0. |
| QV2M | 0.245081 | 0.282412 | 1.000000 | 0.870242 | 0.959385 | 0.960434 | 0.804338 | 0.906144 | -0.071547 | 0.862559 | -0.225449 | -0.256452 | -0 |
| T2M | 0.093258 | 0.164160 | 0.870242 | 1.000000 | 0.913530 | 0.914218 | 0.983356 | 0.981629 | 0.244357 | 0.997515 | -0.207874 | -0.220192 | -0 |
| T2MDEW | 0.231035 | 0.341234 | 0.959385 | 0.913530 | 1.000000 | 0.999970 | 0.854716 | 0.939934 | -0.015643 | 0.905184 | -0.238299 | -0.268686 | -0 |
| T2MWET | 0.230975 | 0.341252 | 0.960434 | 0.914218 | 0.999970 | 1.000000 | 0.855401 | 0.940629 | -0.015500 | 0.905911 | -0.237971 | -0.268292 | -0 |
| T2M_MAX | 0.026773 | 0.111979 | 0.804338 | 0.983356 | 0.854716 | 0.855401 | 1.000000 | 0.937762 | 0.407534 | 0.980101 | -0.216764 | -0.221671 | -0 |
| T2M_MIN | 0.144929 | 0.208285 | 0.906144 | 0.981629 | 0.939934 | 0.940629 | 0.937762 | 1.000000 | 0.065037 | 0.979134 | -0.206382 | -0.225829 | -0 |
| T2M_RANGE | -0.304171 | -0.225935 | -0.071547 | 0.244357 | -0.015643 | -0.015500 | 0.407534 | 0.065037 | 1.000000 | 0.241564 | -0.080163 | -0.043127 | -0 |
| TS | 0.089598 | 0.163830 | 0.862559 | 0.997515 | 0.905184 | 0.905911 | 0.980101 | 0.979134 | 0.241564 | 1.000000 | -0.189823 | -0.202713 | -0 |
| WS10M | 0.049730 | -0.080747 | -0.225449 | -0.207874 | -0.238299 | -0.237971 | -0.216764 | -0.206382 | -0.080163 | -0.189823 | 1.000000 | 0.952217 | 0. |
| WS10M_MAX | 0.060981 | -0.135905 | -0.256452 | -0.220192 | -0.268686 | -0.268292 | -0.221671 | -0.225829 | -0.043127 | -0.202713 | 0.952217 | 1.000000 | 0. |
| WS10M_MIN | 0.023346 | 0.022932 | -0.108789 | -0.125407 | -0.115920 | -0.115882 | -0.141911 | -0.112878 | -0.110952 | -0.110273 | 0.833340 | 0.690087 | 1. |
| WS10M_RANGE | 0.065755 | -0.198332 | -0.269203 | -0.209030 | -0.280702 | -0.280199 | -0.199614 | -0.225256 | 0.018746 | -0.196015 | 0.702896 | 0.866026 | 0. |
| WS50M | 0.069057 | -0.043315 | -0.205971 | -0.193196 | -0.204238 | -0.204143 | -0.195727 | -0.197991 | -0.041778 | -0.180665 | 0.966275 | 0.910717 | 0. |
| WS50M_MAX | 0.079508 | -0.091821 | -0.249961 | -0.206444 | -0.245323 | -0.245147 | -0.196236 | -0.225744 | 0.029737 | -0.193347 | 0.908750 | 0.946710 | 0. |
| WS50M_MIN | 0.057816 | 0.036238 | -0.081554 | -0.112579 | -0.082416 | -0.082497 | -0.133234 | -0.096593 | -0.128844 | -0.102367 | 0.795424 | 0.660428 | 0. |
| WS50M_RANGE | 0.047477 | -0.154479 | -0.246203 | -0.159589 | -0.239335 | -0.239029 | -0.126331 | -0.200157 | 0.163320 | -0.152434 | 0.412412 | 0.592380 | -0 |

*Attributes QV2M, T2M, T2MDEW, T2MWET, T2M_MAX, T2M_MIN and TS have shown strong positive correlation*
*Similary WS10M, WS10M_MAX and WS10M_MIN have shown a strong positve correlation*
*Likewise, WS50M, WS50M_MAX and WS50M_MIN show strong positive correlation*
*However, from the scatter plots above, we see significant variance between the data points, despite the strong positive correlation. Hence we'll retain all these variables, and try other feature selection methods.*

Fig. 5.16 – Correlation Heatmap

## Splitting into train and test

This code snippet is related to the process of splitting a dataset into training and testing sets for a machine learning model. The **train_test_split** function from scikit-learn is used to divide the dataset into training features (**X_train**), testing features (**X_test**), training target values (**y_train**), and testing target values (**y_test**).

The split is performed with a test size of 20% (specified by **test_size=0.2**), and the

random state is set to 0 for reproducibility.

```
Train features shape (1979470, 21)
Train target shape (1979470,)
Test features shape (494868, 21)
Test target shape (494868,)
```

## Standardizing the data

```
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
X_train
```

1. Standardization (or Z-score normalization):

   - Standardization is a preprocessing technique used to rescale the features (independent vari

ables) of a dataset so that they have a mean of 0 and a standard deviation of 1.

   - The formula for standardization is the original feature value, $\mu$ is the mean of the fe

ature values, is the standard deviation of the feature values, and is the

standardized feature value

   1. This transformation centers the data around zero and scales it such that it has a

standard deviation of

2. Importance of Standardization in Machine Learning:

   - Equalizing the Scale:** In many datasets, features can have vastly different

scales. For example, one feature might range from 0 to 1000 while another ranges

from 0 to 1. Without standardization, features with larger scales can dominate the

learning process, leading to biased model results.

   - **Improving Model Performance:** Certain machine learning algorithms are sensitive to

the scale of features. For instance, algorithms based on distances or gradients (

e.g., k-nearest neighbors, support vector machines, gradient descent-based algorithms) can be heavily influenced by the scale of features. Standardization helps mitigate these effects, improving the stability and performance of these algorithms.

   - **Facilitating Interpretability:** Standardization can make the coefficients or weights of the model more interpretable. Since the features are on the same scale, the coefficients represent the change in the target variable corresponding to a one-standard-deviation change in the feature.

3. Steps in Standardization:

   - Instantiate StandardScaler: Create an instance of the `StandardScaler` class, which will be used to scale the features.

   -Fit Transformation on Training Data: Calculate the mean and standard deviation of each feature from the training data and use these statistics to transform the training data.-Apply Transformation on Test Data: Apply the same transformation learned from
 the training data to the test data. This ensures consistency in scaling between training and test datasets.

In summary, standardization is a crucial preprocessing step in machine learning that ensures features are on the same scale, leading to improved model performance, stability, and interpretability. It helps to mitigate issues arising from varying feature scales and ensures that models effectively learn from the data.

This code appears to be performing a standardization preprocessing step on the features (independent variables) of a machine learning dataset using `StandardScaler` from the `sklearn.preprocessing` module. Let's break it down:

- `StandardScaler` is a preprocessing technique used to standardize features by removing the mean and scaling them to unit variance. In other words, it transforms the data such that it has a mean of 0 and a standard deviation of 1. This transformation ensures that the features have the same scale, which can be important for certain machine learning algorithms, particularly those based on distances or gradient descent optimization.

2. **Instantiating StandardScaler**:
   - `sc = StandardScaler()` creates an instance of the `StandardScaler` class. This instance will be used to transform the data.

3.*Transforming Training Data:
   - `X_train = sc.fit_transform(X_train)` applies the standardization transformation to the training data (`X_train`).
- `fit_transform()` method calculates the mean and standard deviation from the training data (`X_train`) and then applies the transformation. This ensures that both the an and standard deviation used for the transformation are based only on the training data.

4. Transforming Test Data:
- `X_test = sc.transform(X_test)` applies the same transformation that was learned from the training data (`X_train`) to the test data (`X_test`).
- `transform()` method applies the transformation without re-calculating the mean and standard deviation, ensuring consistency with the scaling applied to the training data.

```
array([[ 2.39997504, -0.78416609,  0.74119512, ..., -1.42190019,
        -0.16603116,  1.17395993],
       [ 2.44655066,  0.37551891, -0.51918943, ...,  0.61688719,
```

```
      1.00794799,  0.26443762],
[-0.51617644, -0.98505641, -1.00721711, ..., -1.62577893,
  1.30144278, -0.19032353],
...,
[ 2.73117947, -1.01929907, -0.79927488, ...,  0.61688719,
  1.59493757,  1.28765022],
[ 0.22385845,  0.40519588,  1.29500045, ...,  1.2285234 ,
  0.7144532 , -1.55460699],
[-0.51617644,  0.75903678,  1.43292132, ..., -1.01414272,
 -0.16603116,  0.94657935]])
```

# Chapter 6

# Algorithm used in Drought Prediction

## 6.1 Random Forest Algorithm

Machine learning, a fascinating blend of computer science and statistics, has witnessed incredible progress, with one standout algorithm being the **Random Forest**. **Random forests or Random Decision Trees** is a collaborative team of **decision trees** that work together to provide a single output. Originating in 2001 through Leo Breiman, Random Forest has become a cornerstone for machine learning enthusiasts. In this article, we will explore the fundamentals and implementation of **Random Forest Algorithm**.

**What is the Random Forest Algorithm?**

Random Forest algorithm is a powerful tree learning technique in Machine Learning. It works by creating a number of Decision Trees during the training phase. Each tree is constructed using a random subset of the data set to measure a random subset of features in each partition. This randomness introduces variability among individual trees, reducing the risk of overfitting and improving overall prediction performance. In prediction, the algorithm aggregates the results of all trees, either by voting (for classification tasks) or by averaging (for regression tasks) This collaborative decision-making process, supported by multiple trees with their insights, provides an example stable and precise results. Random forests are widely used for classification and regression functions, which are known for their ability to handle complex data, reduce overfitting, and provide reliable forecasts in different environments.

Fig 6.1.1 Random forest algorithm

**What are Ensemble Learning models?**

Ensemble learning models work just like a group of diverse experts teaming up to make decisions – think of them as a bunch of friends with different strengths tackling a problem together. Picture it as a group of friends with different skills working on a project. Each friend excels in a particular area, and by combining their strengths, they create a more robust solution than any individual could achieve alone.

Similarly, in ensemble learning, different models, often of the same type or different types, team up to enhance predictive performance. It's all about leveraging the collective wisdom of the group to overcome individual limitations and make more informed decisions in various machine learning tasks. Some popular ensemble models include- XGBoost, AdaBoost, LightGBM, Random Forest, Bagging, Voting etc.

**What is Bagging and Boosting?**

Bagging is an ensemble learning model, where multiple week models are trained on different subsets of the training data. Each subset is sampled with replacement and prediction is made by averaging the prediction of the week models for regression problem and considering majority vote for classification problem.

Boosting trains multiple based models sequentially. In this method, each model tries to correct the errors made by the previous models. Each model is trained on a modified version of the dataset, the instances that were misclassified by the previous models are given more weight. The final prediction is made by weighted voting.

**How Does Random Forest Work?**

The random Forest algorithm works in several steps which are discussed below–>

**Ensemble of Decision Trees:** Random Forest leverages the power of ensemble learning by constructing an army of Decision Trees. These trees are like individual experts, each specializing in a particular aspect of the data. Importantly, they operate independently, minimizing the risk of the model being overly influenced by the nuances of a single tree.

**Random Feature Selection:** To ensure that each decision tree in the ensemble brings a unique perspective, Random Forest employs random feature selection. During the training of each tree, a random subset of features is chosen. This randomness ensures that each tree focuses on different aspects of the data, fostering a diverse set of predictors within the ensemble.

**Bootstrap Aggregating or Bagging:** The technique of bagging is a cornerstone of Random Forest's training strategy which involves creating multiple bootstrap samples from the original dataset, allowing instances to be sampled with replacement. This results in different

subsets of data for each decision tree, introducing variability in the training process and making the model more robust.

**Decision Making and Voting:** When it comes to making predictions, each decision tree in the Random Forest casts its vote. For <u>classification tasks</u>, the final prediction is determined by the <u>mode</u> (most frequent prediction) across all the trees. In <u>regression tasks</u>, the average of the individual tree predictions is taken. This internal voting mechanism ensures a balanced and collective decision-making process.

**Key Features of Random Forest**

Some of the Key Features of Random Forest are discussed below–>

**High Predictive Accuracy:** Imagine Random Forest as a team of decision-making wizards. Each wizard (decision tree) looks at a part of the problem, and together, they weave their insights into a powerful prediction tapestry. This teamwork often results in a more accurate model than what a single wizard could achieve.

**Resistance to Overfitting:** Random Forest is like a cool-headed mentor guiding its apprentices (decision trees). Instead of letting each apprentice memorize every detail of their training, it encourages a more well-rounded understanding. This approach helps prevent getting too caught up with the training data which makes the model less prone to overfitting.

**Large Datasets Handling:** Dealing with a mountain of data? Random Forest tackles it like a seasoned explorer with a team of helpers (decision trees). Each helper takes on a part of the dataset, ensuring that the expedition is not only thorough but also surprisingly quick.

**Variable Importance Assessment:** Think of Random Forest as a detective at a crime scene, figuring out which clues (features) matter the most. It assesses the importance of each clue in solving the case, helping you focus on the key elements that drive predictions.

**Built-in Cross-Validation:** Random Forest is like having a personal coach that keeps you in check. As it trains each decision tree, it also sets aside a secret group of cases (out-of-bag) for

testing. This built-in validation ensures your model doesn't just ace the training but also performs well on new challenges.

**Handling Missing Values:** Life is full of uncertainties, just like datasets with missing values. Random Forest is the friend who adapts to the situation, making predictions using the information available. It doesn't get flustered by missing pieces; instead, it focuses on what it can confidently tell us.

**Parallelization for Speed:** Random Forest is your time-saving buddy. Picture each decision tree as a worker tackling a piece of a puzzle simultaneously. This parallel approach taps into the power of modern tech, making the whole process faster and more efficient for handling large-scale projects.

## Feature Selection using RFE and Random Forest algorithm

The use of the Random Forest Classifier for feature selection using Recursive Feature Elimination (RFE). Here's a summary of the code:

1. Model Initialization: A Random Forest Classifier model is instantiated with 10 decision trees specified by the hyperparameter n_estimators=10.

2. RFE Configuration: RFE is applied to the model with the goal of selecting 15 features (n_features_to_select=15). This parameter is determined through trial and error.

3. Model Fitting: The RFE is fit to the training data (X_train, y_train), and the process identifies and ranks features based on their importance.

4. Print Statements: The snippet prints the number of selected features (fit.n_features_), a boolean array indicating the selected features (fit.support_), and the ranking of features (fit.ranking_).

5. Selected Features: The code identifies the selected features based on their support and prints their names.

Overall, this code implements feature selection using RFE with a Random Forest

Classifier, providing information on the number of selected features, the boolean array indicat

ing support, the ranking of features, and the names of the selected features. This process aids

in identifying the most relevant features for model training and

prediction.

```
Num Features: 15
Selected Features: [False  True  True  True  True False  True  True  True  Tru
e  True False
 False  True  True  True False  True  True False  True]
Feature Ranking: [4 1 1 1 1 2 1 1 1 1 1 5 6 1 1 1 3 1 1 7 1]
Index(['PS', 'QV2M', 'T2M', 'T2MDEW', 'T2M_MAX', 'T2M_MIN', 'T2M_RANGE', 'TS',
       'WS10M', 'WS10M_RANGE', 'WS50M', 'WS50M_MAX', 'WS50M_RANGE', 'year',
       'day'],
      dtype='object')
```

It involves preprocessing the independent variables for a machine learning model.

Here's a breakdown of the code:

1. **Dropping Columns:** Several columns are dropped from the

   **independent_variables** DataFrame. These columns include 'PRECTOT',

   'T2MWET', 'WS10M_MAX', 'WS10M_MIN', 'WS50M_MIN', and 'month'.

   This step is likely performed to exclude certain features from the model.

2. **Displaying Head:** The **head()** function is used to display the first few rows of the mo

   dified **independent_variables** DataFrame after dropping the specified

   columns.

3. **Train-Test Split:** The data is split into training and testing sets using the **train_test_s

   plit** function from scikit-learn. The training set (**X_train**, **y_train**) constitutes 80% of

   the data, while the testing set (**X_test**, **y_test**) comprises 20%.

   The random state is set to 0 for reproducibility.

4. **Print Statements:** The snippet prints the shapes of the training and testing sets, indicating the number of samples and features in each set.

5. **Standardization:** The **StandardScaler** is used to standardize the features in both the training and testing sets. This step ensures that the features have zero mean and unit variance, a common preprocessing step for many machine learning models.

In summary, this code performs data preprocessing by dropping specified columns, splitting the data into training and testing sets, and standardizing the features using the **StandardScaler**. These steps prepare the data for subsequent machine learning model training and evaluation.

```
Train features shape (1979470, 15)
Train target shape (1979470,)
Test features shape (494868, 15)
Test target shape (494868,)
```

## Fixing class imbalance:

Class imbalance is a common problem in machine learning where the distribution of classes in the dataset is uneven. This means that some classes have significantly more samples than others, which can lead to biased models that perform poorly on the minority class. Addressing class imbalance is crucial for building robust and fair models, especially in applications such as fraud detection, medical diagnosis, and spam filtering. Here are several strategies to handle class imbalance effectively:

**Data-Level Approaches**

1. **Resampling Techniques**

   - **Oversampling:** Increase the number of minority class samples by duplicating them or creating synthetic samples. Techniques include:

     - **Random Oversampling:** Randomly duplicate minority class examples.

     - **SMOTE (Synthetic Minority Over-sampling Technique):** Generate synthetic examples by interpolating between existing minority samples.

     - **ADASYN (Adaptive Synthetic Sampling):** Generate synthetic samples with a focus on difficult-to-learn examples.

   - **Undersampling:** Reduce the number of majority class samples. Techniques include:

     - **Random Undersampling:** Randomly remove majority class examples.

- **Tomek Links:** Remove examples that are very close to examples of the other class.

- **NearMiss:** Select majority class examples that are closest to minority class examples.

2. **Hybrid Methods**

- Combine oversampling and undersampling methods to balance the dataset effectively. For example, you can use SMOTE to oversample the minority class and Tomek Links to undersample the majority class.

**Algorithm-Level Approaches**

1. **Cost-Sensitive Learning**

- Modify the learning algorithm to incorporate the cost of misclassifying minority class examples. This can be done by assigning higher misclassification penalties to minority class errors.

- Many machine learning algorithms, such as decision trees, SVMs, and neural networks, can be adjusted to account for different misclassification costs.

2. **Class Weighting**

- Adjust the weights of classes in the learning algorithm to give more importance to the minority class. Most libraries (e.g., scikit-learn) allow you to set class weights directly.

3. **Ensemble Methods**

- Use ensemble techniques that combine multiple models to improve classification performance on imbalanced datasets. Methods include:

- **Balanced Random Forests:** Modify the standard random forest algorithm to balance the class distribution when sampling subsets of data.

- **EasyEnsemble:** An ensemble of models trained on different balanced subsets created by undersampling the majority class.

**Evaluation Metrics**

Standard metrics like accuracy can be misleading for imbalanced datasets. Instead, consider the following metrics:

1. **Precision, Recall, and F1-Score**

   - Precision: Proportion of true positive predictions among all positive predictions.

   - Recall: Proportion of true positive predictions among all actual positive instances.

   - F1-Score: Harmonic mean of precision and recall.

2. **Confusion Matrix**

   - Provides a detailed breakdown of true positives, false positives, true negatives, and false negatives.

3. **ROC-AUC (Receiver Operating Characteristic - Area Under the Curve)**

   - Measures the trade-off between true positive rate and false positive rate.

4. **PR-AUC (Precision-Recall Area Under the Curve)**

   - Especially useful when the positive class is the minority.

**Upsampling using SMOTE:**

SMOTE (Synthetic Minority Over-sampling Technique) is a popular method for addressing class imbalance in datasets by generating synthetic samples for the minority class. Unlike simple oversampling, which duplicates existing minority class examples, SMOTE creates new synthetic samples by interpolating between existing ones. This approach helps to enhance the generalization of the model by creating more diverse training samples, thereby mitigating the risk of overfitting.

**How SMOTE Works**

1. **Selection of Minority Class Instances:** SMOTE first selects samples from the minority class.

2. **K-nearest Neighbors:** For each selected minority class instance, SMOTE identifies its k-nearest neighbors (typically using Euclidean distance).

3. **Synthetic Sample Generation:** New synthetic instances are created by randomly selecting one of the k-nearest neighbors and generating a synthetic point along the line segment joining the minority instance and its neighbor. This is achieved using the formula:

*synthetic_sample=minority_instance+gap×(neighbor−minority_instance)*

*where the gap is a random number between 0 and 1.*

**Benefits of SMOTE**

- **Improved Model Performance:** By balancing the class distribution, models trained with SMOTE typically perform better on the minority class.

- **Reduced Overfitting:** SMOTE helps create diverse synthetic samples, which can reduce the risk of overfitting compared to simple oversampling.

- **Versatility:** SMOTE can be applied to a variety of classifiers and is not limited to specific types of algorithms.

**Limitations of SMOTE**

- **Overlapping Classes:** SMOTE can create synthetic examples that overlap with other classes, potentially introducing noise.

- **Higher Dimensionality:** In high-dimensional spaces, synthetic samples may not be as effective due to the curse of dimensionality.

- **Computational Cost:** SMOTE can be computationally expensive, especially for large datasets.

SMOTE is a powerful technique for handling class imbalance in machine learning. By generating synthetic samples for the minority class, it helps create a more balanced and representative dataset, improving the performance and robustness of machine learning models. However, it is essential to consider its limitations and apply appropriate evaluation metrics to ensure the effectiveness of the resulting model.

**Upsampling** is a technique used to address class imbalance in datasets by increasing the number of samples in the minority class. This can help prevent the model from being biased

toward the majority class and improve the performance on the minority class. There are several methods to perform upsampling, each with its own advantages and use cases.

Upsampling is a crucial technique for handling class imbalance in machine learning. Whether using simple methods like random oversampling or more sophisticated techniques like SMOTE and ADASYN, the goal is to create a balanced dataset that allows the model to perform well on both majority and minority classes. By carefully selecting and implementing the appropriate upsampling technique, and by using suitable evaluation metrics, you can build more robust and fair models for imbalanced datasets.

**Steps to Implement SMOTE for Upsampling**

1. **Understand the Dataset:**

   - **Assess Imbalance:** Before applying SMOTE, examine your dataset to understand the extent of the class imbalance. Calculate the class distribution to see how imbalanced your data is.

2. **Split the Dataset:**

   - **Train-Test Split:** Divide your dataset into training and testing sets. It's crucial to perform the split before applying SMOTE to ensure that the evaluation metrics reflect the model's performance on unseen data.

3. **Initialize SMOTE:**

   - **SMOTE Parameters:** Choose appropriate parameters for SMOTE, such as the **sampling_strategy**, which determines the desired ratio of minority to majority class samples, and **k_neighbors**, which specifies the number of nearest neighbors to use for generating synthetic samples.

4. **Fit SMOTE on Training Data:**

- **Apply SMOTE:** Use SMOTE to fit the training data and generate synthetic samples. This involves identifying the k-nearest neighbors for each minority class instance and creating synthetic samples along the line segments connecting the minority instances with their neighbors.

5. **Create Synthetic Samples:**

- **Generate New Data Points:** For each selected minority class instance, generate new synthetic samples by interpolating between the instance and its nearest neighbors. This step enhances the diversity of the minority class samples in the training set.

6. **Combine Original and Synthetic Samples:**

- **Resample Training Set:** Combine the original minority class samples with the newly generated synthetic samples to create a balanced training set. Ensure that the majority class samples remain the same while the minority class samples are augmented.

7. **Train the Model:**

- **Fit Model on Resampled Data:** Train your machine learning model using the resampled training set. The balanced dataset helps the model learn from both classes effectively.

8. **Evaluate the Model:**

   - **Test on Unseen Data:** After training, evaluate your model on the test set (which has not been upsampled) to measure its performance. Use appropriate evaluation metrics such as precision, recall, F1-score, and ROC-AUC to get a comprehensive understanding of how well your model handles the imbalanced data.

9. **Tune Parameters:**

   - **Optimize SMOTE and Model Parameters:** Experiment with different SMOTE parameters and model hyperparameters to optimize performance. This may involve varying the **sampling_strategy** and **k_neighbors** values to find the best balance between generating useful synthetic samples and avoiding overfitting.

10. **Analyze Results:**

    - **Detailed Analysis:** Review the confusion matrix and classification report to understand the model's performance in detail. Pay particular attention to how well the model predicts the minority class.

**Additional Considerations**

- **Validation Strategy:** Use cross-validation to ensure that your model's performance is consistent across different subsets of the data.

- **Data Preprocessing:** Perform any necessary preprocessing steps such as scaling or normalization before applying SMOTE, as synthetic samples should be generated in the same feature space.

- **Domain Knowledge:** Incorporate domain knowledge to guide the SMOTE process, especially in choosing the right number of neighbors and understanding the context of the synthetic samples.

By following these steps, you can effectively implement SMOTE to handle class imbalance in your dataset, leading to improved model performance on both the majority and minority classes.

The code snippet demonstrates how to apply SMOTE (Synthetic Minority Over-sampling Technique) to a training dataset to address class imbalance. Let's break down each part of the code in detail.

**Importing and Initializing SMOTE**

sm = SMOTE(random_state=5)

- **SMOTE Initialization:** Here, an instance of SMOTE is created. The **random_state** parameter is set to **5** to ensure reproducibility. This means that every time the code is run with the same **random_state**, the same synthetic samples will be generated.

**Applying SMOTE to the Training Data**

X_train_ures_SMOTE, y_train_ures_SMOTE = sm.fit_resample(X_train, y_train.ravel())

- **fit_resample Method:** The **fit_resample** method is called on the SMOTE instance **sm** using the original training features **X_train** and the flattened training labels **y_train.ravel**().

  - **X_train and y_train:** These are the features and labels of the training dataset, respectively.

- **ravel():** The **ravel()** method is used to flatten the **y_train** array, ensuring it is a 1D array, which is required by SMOTE.

- **Output:** The method returns the resampled feature matrix **X_train_ures_SMOTE** and the resampled label vector **y_train_ures_SMOTE**.

**Printing the Shape of the Data Before and After Oversampling**

- **Before Oversampling:**

  - The shape of **X_train** and **y_train** is printed to show the dimensions of the training dataset before applying SMOTE.

- **After Oversampling:**

  - The shape of **X_train_ures_SMOTE** and **y_train_ures_SMOTE** is printed to show the new dimensions of the training dataset after applying SMOTE.

**Printing the Class Distribution Before and After Oversampling**

- **Class Counts Before and After Oversampling:**

  - For each class label (0 through 5), the code prints the count of samples before and after oversampling.

  - **Before Oversampling: sum(y_train == label)** calculates the number of samples for each class label in the original training set.

  - **After Oversampling: sum(y_train_ures_SMOTE == label)** calculates the number of samples for each class label in the resampled training set.

**Summary**

- **Initialization:** SMOTE is initialized with a random state to ensure reproducibility.

- **Resampling:** The training data is resampled using SMOTE to generate synthetic samples for the minority class.

- **Shape Comparison:** The shapes of the training feature and label arrays are printed before and after resampling to show the increase in the number of samples.

- **Class Distribution:** The number of samples in each class is printed before and after resampling to illustrate how SMOTE has balanced the class distribution.

This detailed breakdown explains how SMOTE is used to address class imbalance by generating synthetic samples, and how the effects of this resampling are evaluated by comparing dataset shapes and class distributions before and after the process.

**OUTPUT**

Before OverSampling, the shape of train_X: (1979470, 15)

Before OverSampling, the shape of train_y: (1979470,)

After OverSampling, the shape of train_X: (7173186, 15)

After OverSampling, the shape of train_y: (7173186,)

Counts of label '0' - Before Oversampling:1195531, After OverSampling: 1195531

Counts of label '1' - Before Oversampling:332490, After OverSampling: 1195531

Counts of label '2' - Before Oversampling:209363, After OverSampling: 1195531

Counts of label '3' - Before Oversampling:139009, After OverSampling: 1195531

Counts of label '4' - Before Oversampling:74270, After OverSampling: 1195531

Counts of label '5' - Before Oversampling:28807, After OverSampling: 1195531

## **Downsampling in Machine Learning**

Downsampling, also known as undersampling, is a technique used to address class imbalance in datasets by reducing the number of samples in the majority class. This approach aims to balance the class distribution, helping machine learning models to perform better on both the majority and minority classes.

**Methods of Downsampling**

1. **Random Undersampling**

   - **Description:** Randomly removes samples from the majority class until the class distribution is balanced.

   - **Pros:** Simple and easy to implement.

   - **Cons:** Can lead to loss of valuable information, as potentially useful samples from the majority class are discarded.

2. **Cluster Centroids**

   - **Description:** Reduces the number of majority class samples by replacing clusters of majority samples with their centroids.

   - **Pros:** Retains the most representative samples of the majority class.

- **Cons:** More complex to implement and computationally intensive compared to random undersampling.

3. **NearMiss**

    - **Description:** Selects majority class samples that are closest to the minority class samples based on certain distance measures.

    - **Pros:** Ensures that the selected majority samples are more informative and relevant to the minority class.

    - **Cons:** Can be computationally expensive and may not be effective if the minority class samples are not well-separated.

4. **Tomek Links**

    - **Description:** Removes majority class samples that form Tomek links with minority class samples. A Tomek link is a pair of samples from different classes that are each other's nearest neighbors.

    - **Pros:** Helps in cleaning the data by removing borderline majority samples that are likely to be noise.

    - **Cons:** May not significantly reduce the number of majority samples if there are few Tomek links.

**Steps to Implement Downsampling**

**1. Understand the Dataset**

- **Assess Imbalance:** Before applying downsampling, examine your dataset to understand the extent of class imbalance by calculating the class distribution.

**2. Split the Dataset**

- **Train-Test Split:** Divide your dataset into training and testing sets. This ensures that the evaluation metrics reflect the model's performance on unseen data and avoids data leakage.

## 3. Apply Downsampling

- **Select Method:** Choose an appropriate downsampling method based on your dataset and the problem at hand. Common methods include random undersampling, NearMiss, and Tomek Links.

## 4. Resample the Training Data

- **Random Undersampling:**

  - Randomly remove samples from the majority class until the class distribution is balanced.

- **NearMiss:**

  - Use a distance measure to select majority class samples that are closest to minority class samples.

- **Tomek Links:**

  - Identify and remove majority class samples that form Tomek links with minority class samples.

## 5. Train the Model

- **Fit Model on Resampled Data:** Train your machine learning model using the downsampled training set. The balanced dataset helps the model to learn from both classes effectively.

**6. Evaluate the Model**

- **Test on Unseen Data:** Evaluate your model on the original test set to measure its performance. Use appropriate evaluation metrics such as precision, recall, F1-score, and ROC-AUC to get a comprehensive understanding of how well your model handles the imbalanced data.

**Example Use Case**

Let's consider a binary classification problem with a highly imbalanced dataset. Here's how you might implement random undersampling:

**Steps:**

1. **Examine Class Distribution:** Before downsampling, check the number of samples in each class to understand the extent of imbalance.

2. **Split Dataset:** Divide your dataset into training and testing sets.

3. **Random Undersampling:**

   - Randomly select a subset of the majority class samples.

   - Combine these with all the minority class samples to create a balanced training set.

4. **Train Model:** Use the resampled dataset to train your model.

5. **Evaluate Model:** Test the model on the original test set and evaluate its performance using suitable metrics.

**Benefits of Downsampling**

1. **Balancing Classes:** Reduces the class imbalance, leading to a more balanced dataset.

2. **Improved Model Performance:** Helps in improving the performance of models on the minority class.

3. **Reduced Overfitting:** By reducing the number of majority class samples, the model is less likely to overfit to the majority class.

**Drawbacks of Downsampling**

1. **Loss of Information:** Removing samples from the majority class can lead to loss of valuable information, potentially degrading model performance.

2. **Not Suitable for Small Datasets:** In small datasets, downsampling can significantly reduce the size of the training set, making it difficult for the model to learn effectively.

3. **Risk of Underfitting:** If too many samples are removed, the model may underfit, failing to capture the underlying patterns in the data.

**Conclusion**

Downsampling is a valuable technique for addressing class imbalance in machine learning datasets. By carefully selecting and implementing the appropriate downsampling method, and by using suitable evaluation metrics, you can build more robust and fair models. However, it is important to balance the benefits of reducing class imbalance with the potential loss of information to ensure optimal model performance.

# Neighborhood Cleaning Rule (NCL) for Downsampling

The Neighborhood Cleaning Rule (NCL) is an advanced undersampling technique used to address class imbalance in datasets. Unlike random undersampling, which randomly removes samples from the majority class, NCL combines elements of both undersampling and data cleaning to improve the quality of the dataset. NCL aims to remove noisy and borderline samples that could mislead the classifier, enhancing the learning process.

**Key Concepts of NCL**

1. **k-Nearest Neighbors (k-NN):** NCL uses the k-NN algorithm to identify the neighborhood of each sample. This is crucial for determining which samples to remove.

2. **Noise Removal:** By identifying and removing misclassified samples, NCL cleans the dataset, thereby improving the overall quality.

3. **Borderline Cleaning:** NCL also targets borderline samples that lie close to the decision boundary between classes, as these can cause confusion for the classifier.

**Steps to Implement the Neighborhood Cleaning Rule**

**1. Understand the Dataset**

- **Assess Imbalance:** Before applying NCL, assess the class distribution in your dataset to understand the extent of the imbalance.

**2. Split the Dataset**

- **Train-Test Split:** Divide your dataset into training and testing sets. This is important to prevent data leakage and ensure the evaluation metrics reflect the model's performance on unseen data.

**3. Initialize and Apply NCL**

- **k-Nearest Neighbors:** Select an appropriate value for **k** in the k-NN algorithm. This parameter determines how many neighbors will be considered when identifying misclassified and borderline samples.

- **Fit NCL:** Apply the NCL method to the training data. This involves the following sub-steps:

  1. **For each sample in the training set, identify its k-nearest neighbors.**

  2. **Classify the sample based on its neighbors:**

     - If a majority class sample is misclassified by its neighbors, it is removed.

     - If a minority class sample is correctly classified by its neighbors but some of its neighbors are majority class samples, those majority samples are removed.

**4. Create a Cleaned and Balanced Dataset**

- **Combine Cleaned Data:** After removing the identified noisy and borderline majority samples, combine the remaining majority samples with all minority class samples to form a balanced and cleaned training set.

**5. Train the Model**

- **Fit Model on Cleaned Data:** Train your machine learning model using the cleaned and balanced training set. This helps the model learn from a high-quality dataset, improving its performance.

**6. Evaluate the Model**

- **Test on Unseen Data:** Evaluate the model on the test set to measure its performance. Use appropriate evaluation metrics such as precision, recall, F1-score, and ROC-AUC to get a comprehensive understanding of the model's effectiveness.

**Example Use Case**

Consider a binary classification problem with a highly imbalanced dataset. Here's a detailed outline of implementing NCL:

**Steps:**

1. **Examine Class Distribution:** Before applying NCL, check the number of samples in each class to understand the extent of the imbalance.

2. **Split Dataset:** Divide your dataset into training and testing sets.

3. **Apply NCL:**

   - **Choose k:** Select an appropriate value for **k** in the k-NN algorithm.

   - **Identify Misclassified and Borderline Samples:** For each sample in the training set, use k-NN to identify misclassified majority samples and borderline majority samples.

   - **Remove Noisy and Borderline Samples:** Remove the identified majority class samples.

4. **Train Model:** Use the cleaned and balanced dataset to train your model.

5. **Evaluate Model:** Test the model on the original test set and evaluate its performance using suitable metrics.

**Benefits of NCL**

1. **Improved Data Quality:** By removing noisy and borderline samples, NCL enhances the overall quality of the training data.

2. **Enhanced Model Performance:** Cleaned and balanced datasets often lead to better model performance, especially on minority class samples.

3. **Reduced Overfitting:** Removing noisy samples helps in reducing overfitting, making the model more generalizable.

**Drawbacks of NCL**

1. **Computational Complexity:** NCL can be computationally intensive due to the need for repeated k-NN searches, especially in large datasets.

2. **Parameter Sensitivity:** The performance of NCL can be sensitive to the choice of **k**. Selecting an inappropriate value can lead to suboptimal results.

3. **Potential Underfitting:** Excessive removal of samples, especially if not done carefully, can lead to underfitting, where the model fails to capture the underlying patterns in the data.

**Conclusion**

The Neighborhood Cleaning Rule (NCL) is a sophisticated undersampling technique that not only addresses class imbalance but also improves the quality of the dataset by removing noisy and borderline samples. By using k-NN to identify and clean these samples, NCL helps in building more robust and accurate machine learning models. However, it is important to carefully choose parameters and consider the computational complexity when implementing NCL.

**OUTPUT**

Before UnderSampling, the shape of train_X: (1979470, 15)

Before UnderSampling, the shape of train_y: (1979470,)

After UnderSampling, the shape of train_X: (1568791, 15)

After UnderSampling, the shape of train_y: (1568791,)

Counts of label '0' - Before UnderSampling:1195531, After UnderSampling: 1113612

Counts of label '1' - Before UnderSampling:332490, After UnderSampling: 200048

Counts of label '2' - Before UnderSampling:209363, After UnderSampling: 113745

Counts of label '3' - Before UnderSampling:139009, After UnderSampling: 73848

Counts of label '4' - Before UnderSampling:74270, After UnderSampling: 38731

Counts of label '5' - Before UnderSampling:28807, After UnderSampling: 28807

## Near Miss for Downsampling

Near Miss is a family of undersampling techniques designed to address class imbalance by selecting majority class samples that are closest to minority class samples based on a specified distance measure. The primary goal of Near Miss is to create a more balanced dataset that retains informative samples from the majority class, thereby improving the classifier's performance on the minority class.

**Key Concepts of Near Miss**

1. **Distance Measures:** Near Miss relies on distance measures (e.g., Euclidean distance) to determine the closeness between majority and minority class samples.

2. **Selection Criteria:** Different versions of Near Miss use various criteria to select which majority class samples to retain.

**Variants of Near Miss**

There are three main variants of Near Miss, each using a different strategy to select majority class samples:

1. **Near Miss-1**

2. **Near Miss-2**

3. **Near Miss-3**

**Near Miss-1**

**Selection Criteria:**

- For each minority class sample, select the majority class samples that have the smallest average distance to the k nearest minority class samples.

**Steps:**

1. For each majority class sample, compute its distance to all minority class samples.

2. Calculate the average distance to the k nearest minority class samples.

3. Select the majority class samples with the smallest average distances.

**Pros:**

- Ensures that the selected majority samples are close to multiple minority samples, enhancing the classifier's ability to discriminate between classes.

**Cons:**

- Can lead to a bias towards certain regions of the feature space where minority samples are densely packed.

**Near Miss-2**

**Selection Criteria:**

- For each minority class sample, select the majority class samples that have the smallest average distance to the k farthest minority class samples.

**Steps:**

1. For each majority class sample, compute its distance to all minority class samples.

2. Calculate the average distance to the k farthest minority class samples.

3. Select the majority class samples with the smallest average distances.

**Pros:**

- Ensures that the selected majority samples are diverse and spread out, covering different regions of the feature space.

**Cons:**

- May include more noise as it considers the farthest minority samples.

**Near Miss-3**

**Selection Criteria:**

- For each minority class sample, select the k nearest majority class samples.

**Steps:**

1. For each minority class sample, compute its distance to all majority class samples.

2. Select the k nearest majority class samples.

3. Combine these selected samples to form the undersampled majority class.

**Pros:**

- Directly focuses on the minority class samples, ensuring that the nearest majority samples are included.

**Cons:**

- May result in selecting samples that are not diverse enough, potentially missing important regions of the feature space.

**Steps to Implement Near Miss**

**1. Understand the Dataset**

- **Assess Imbalance:** Examine the class distribution to understand the extent of imbalance and decide on the appropriate Near Miss variant.

**2. Split the Dataset**

- **Train-Test Split:** Divide the dataset into training and testing sets to prevent data leakage and ensure valid model evaluation.

### 3. Initialize and Apply Near Miss

- **Select Near Miss Variant:** Choose the appropriate Near Miss variant (1, 2, or 3) based on your dataset and problem requirements.

- **Set Parameters:** Determine the value of **k**, the number of nearest or farthest neighbors to consider.

### 4. Resample the Training Data

- **Compute Distances:** Calculate the distances between majority and minority class samples based on the selected Near Miss variant.

- **Select Samples:** Identify and select the majority class samples according to the specified criteria.

- **Create Balanced Dataset:** Combine the selected majority samples with all minority class samples to form the balanced training set.

### 5. Train the Model

- **Fit Model on Resampled Data:** Train your machine learning model using the resampled dataset, which now has a more balanced class distribution.

### 6. Evaluate the Model

- **Test on Unseen Data:** Evaluate the model on the original test set to measure its performance using suitable metrics such as precision, recall, F1-score, and ROC-AUC.

**Example Use Case**

Consider a binary classification problem with an imbalanced dataset. Here's a detailed outline of implementing Near Miss-1:

**Steps:**

1. **Examine Class Distribution:** Before applying Near Miss, check the number of samples in each class to understand the imbalance.

2. **Split Dataset:** Divide the dataset into training and testing sets.

3. **Apply Near Miss-1:**

   - **Choose k:** Select an appropriate value for **k** (e.g., 3).

   - **Compute Distances:** Calculate the average distance from each majority class sample to the k nearest minority class samples.

   - **Select Samples:** Retain the majority class samples with the smallest average distances.

4. **Train Model:** Use the balanced dataset to train your model.

5. **Evaluate Model:** Test the model on the original test set and evaluate its performance using suitable metrics.

**Benefits of Near Miss**

1. **Balancing Classes:** Effectively balances the class distribution by carefully selecting majority class samples.

2. **Improved Model Performance:** Enhances the model's ability to learn from both classes, especially the minority class.

3. **Reduced Noise:** By focusing on the closest samples, Near Miss reduces the likelihood of including noisy majority class samples.

**Drawbacks of Near Miss**

1. **Computational Complexity:** Calculating distances for all samples can be computationally expensive, especially for large datasets.

2. **Parameter Sensitivity:** The performance of Near Miss methods can be sensitive to the choice of **k**. Incorrect values can lead to suboptimal results.

3. **Potential Bias:** Depending on the variant, Near Miss might introduce bias by over-representing certain regions of the feature space.

**Conclusion**

Near Miss is a valuable undersampling technique that addresses class imbalance by carefully selecting majority class samples based on their proximity to minority class samples. By choosing the appropriate variant and parameters, you can create a balanced and informative training set that improves the performance of your machine learning models. However, it is important to consider the computational complexity and potential biases when implementing Near Miss.

**OUTPUT**

Before UnderSampling, the shape of train_X: (1979470, 15)

Before UnderSampling, the shape of train_y: (1979470,)


After UnderSampling, the shape of train_X: (172842, 15)

After UnderSampling, the shape of train_y: (172842,)


Counts of label '0' - Before UnderSampling:1195531, After UnderSampling: 28807

Counts of label '1' - Before UnderSampling:332490, After UnderSampling: 28807

Counts of label '2' - Before UnderSampling:209363, After UnderSampling: 28807

Counts of label '3' - Before UnderSampling:139009, After UnderSampling: 28807

Counts of label '4' - Before UnderSampling:74270, After UnderSampling: 28807

Counts of label '5' - Before UnderSampling:28807, After UnderSampling: 28807


# PCA for dimensionality reduction

Principal Component Analysis (PCA) is a widely used technique in machine learning and stat

istics for dimensionality reduction. It transforms a dataset with possibly

correlated variables into a set of linearly uncorrelated variables called principal

components. These principal components are ordered by the amount of variance they capture

from the original data, allowing us to reduce the dimensionality while

preserving most of the information.

**Key Concepts of PCA**

1. **Variance and Covariance:**

   - **Variance:** Measures how much the data varies. High variance

     indicates more spread out data.

- **Covariance:** Measures the extent to which two variables change together. Positive covariance indicates that the variables increase together, while negative covariance indicates that as one increases, the other decreases.

2. **Eigenvectors and Eigenvalues:**

- **Eigenvectors:** Directions along which the data is spread out. They are the principal components.

- **Eigenvalues:** Scalars that represent the magnitude of the variance in the direction of their corresponding eigenvector.

3. **Principal Components:**

- The principal components are new variables that are constructed as linear combinations of the original variables. The first principal component captures the most variance, the second captures the second most variance, and so on.

**Steps to Implement PCA**

**1. Standardize the Data**

- **Purpose:** Ensure that each feature contributes equally to the analysis. Features with larger scales could dominate the principal components if not standardized.

- **Method:** Subtract the mean of each feature and divide by the standard deviation.

**2. Compute the Covariance Matrix**

- **Purpose:** Understand how variables in the dataset are related to each other.

- **Method:** Calculate the covariance matrix of the standardized data.

### 3. Compute the Eigenvalues and Eigenvectors

- **Purpose:** Identify the principal components (eigenvectors) and their significance (eigenvalues).

- **Method:** Decompose the covariance matrix into its eigenvalues and eigenvectors.

### 4. Sort Eigenvalues and Eigenvectors

- **Purpose:** Rank the principal components by their significance.

- **Method:** Sort the eigenvalues in descending order and sort the eigenvectors accordingly.

### 5. Select the Number of Principal Components

- **Purpose:** Determine how many principal components to keep. This is usually based on the amount of variance explained.

- **Method:** Choose the top **k** eigenvectors based on the cumulative variance explained.

### 6. Transform the Data

- **Purpose:** Project the data onto the new feature space defined by the principal components.

- **Method:** Multiply the standardized data by the selected principal components.

### Practical Considerations

1. **Explained Variance:**

   - **Definition:** The proportion of the dataset's total variance that is captured by each principal component.

   - **Usage:** Typically, we select enough principal components to explain a high percentage (e.g., 95%) of the total variance.

2. **Interpretation of Principal Components:**

   - **Loadings:** The coefficients of the original variables in the principal component's linear combination. These loadings help interpret the principal components.

3. **Reconstruction Error:**

   - **Definition:** The difference between the original data and the data reconstructed from the principal components. Lower error indicates better representation.

4. **Scalability:**

   - **Large Datasets:** For very large datasets, approximate methods or incremental PCA can be used to handle memory and computational constraints.

**Example Use Case**

Consider a dataset with multiple features, where some features are correlated. The goal is to reduce the dimensionality to two principal components for visualization.

**Steps:**

1. **Examine Class Distribution:** Before applying Near Miss, check the number of samples in each class to understand the imbalance.

2. **Split Dataset:** Divide the dataset into training and testing sets.

3. **Apply Near Miss-1:**

   - **Choose k:** Select an appropriate value for **k** (e.g., 3).

   - **Compute Distances:** Calculate the average distance from each majority class sample to the k nearest minority class samples.

   - **Select Samples:** Retain the majority class samples with the smallest average distances.

4. **Train Model:** Use the balanced dataset to train your model.

5. **Evaluate Model:** Test the model on the original test set and evaluate its performance using suitable metrics.

**Benefits of PCA**

1. **Dimensionality Reduction:** Reduces the number of features, simplifying models and reducing computational cost.

2. **Noise Reduction:** By focusing on the most significant components, PCA can help filter out noise.

3. **Visualization:** Makes it easier to visualize high-dimensional data in 2D or 3D.

**Drawbacks of PCA**

1. **Interpretability:** Principal components are linear combinations of original features, which can make interpretation difficult.

2. **Assumption of Linearity:** PCA assumes linear relationships among variables, which might not capture complex, non-linear patterns.

3. **Information Loss:** While reducing dimensions, some information might be lost, especially if the selected number of components explains insufficient variance.

**Conclusion**

PCA is a powerful tool for dimensionality reduction, helping to simplify complex datasets wh ile retaining most of the important information. By transforming correlated variables into a set of uncorrelated principal components, PCA facilitates improved model performance and data visualization. However, it is essential to carefully select the number of components and understand the trade-offs involved in information loss and int erpretability.

# LDA For Dimensionality Reduction

Linear Discriminant Analysis (LDA) is a supervised learning algorithm commonly used for classification tasks, but it can also be effectively used for dimensionality reduction. Unlike Principal Component Analysis (PCA), which is unsupervised and focuses on maximizing variance, LDA seeks to find a linear combination of features that best separate two or more classes. LDA is particularly useful when the dataset has labels and the goal is to maximize class separability.

**Key Concepts of LDA**

1. **Class Separation:** LDA aims to maximize the distance between the means of different classes while minimizing the variation within each class.

2. **Scatter Matrices:**

   - **Within-class scatter matrix (Sw):** Measures the scatter (spread) of data points within each class.

   - **Between-class scatter matrix (Sb):** Measures the scatter of the means of different classes relative to the overall mean.

3. **Eigenvectors and Eigenvalues:** Similar to PCA, LDA involves finding eigenvectors and eigenvalues. However, these are derived from the scatter matrices to maximize class separation.

**Steps to Implement LDA**

**1. Standardize the Data**

- **Purpose:** Ensure that each feature contributes equally to the analysis. Features with larger scales could dominate the components if not standardized.

- **Method:** Subtract the mean of each feature and divide by the standard deviation.

**2. Compute the Mean Vectors**

- **Purpose:** Calculate the mean vector for each class, which will be used to compute scatter matrices.

- **Method:** For each class, compute the mean of the feature vectors.

**3. Compute the Scatter Matrices**

- **Within-class Scatter Matrix (Sw):**

  - **Purpose:** Measure the spread of data points within each class.

  - **Method:** Sum the covariance matrices of each class, weighted by the number of samples in the class.

- **Between-class Scatter Matrix (Sb):**

  - **Purpose:** Measure the spread of the class means relative to the overall mean.

  - **Method:** Compute the outer product of the difference between each class mean and the overall mean, weighted by the number of samples in the class.

**4. Compute the Eigenvectors and Eigenvalues**

- **Purpose:** Identify the directions (eigenvectors) that maximize class separability and their corresponding magnitudes (eigenvalues).

- **Method:** Solve the generalized eigenvalue problem for the matrix $Sw^{(-1)}Sb$.

**5. Select the Top Eigenvectors**

- **Purpose:** Choose the eigenvectors corresponding to the largest eigenvalues to form the linear discriminants.

- **Method:** Sort the eigenvalues in descending order and select the top eigenvectors.

### 6. Transform the Data

- **Purpose:** Project the original data onto the new feature space defined by the selected eigenvectors.

- **Method:** Multiply the standardized data by the selected eigenvectors.

**Practical Considerations**

1. **Number of Components:**

    - **Maximum Number:** The maximum number of linear discriminants is $\min($ $C-1, n-1)\min(C-1, n-1)$, where $C$ is the number of classes and $n$ is the number of features.

    - **Selection:** Choose the number of components that best balance the trade-off between dimensionality reduction and information retention.

2. **Assumptions:**

    - **Linearity:** Assumes that the relationship between the features and the class labels is linear.

    - **Normality:** Assumes that the features are normally distributed within each class.

    - **Equal Covariance:** Assumes that the covariance matrices of the features are similar for all classes.

3. **Interpretation of Linear Discriminants:**

    - The linear discriminants are linear combinations of the original features, which can be interpreted in terms of the contributions of each feature to class separability.

**Example Use Case**

Consider a classification problem where the goal is to reduce the dimensionality of the dataset while preserving the ability to distinguish between classes.

**Steps:**

1. **Standardize the Data:** Standardize each feature to have zero mean and unit variance.

2. **Compute the Mean Vectors:** Calculate the mean vector for each class.

3. **Compute the Scatter Matrices:**

   - Compute the within-class scatter matrix $Sw$ by summing the covariance matrices of each class.

   - Compute the between-class scatter matrix $Sb$ by calculating the outer product of the difference between each class mean and the overall mean.

4. **Compute Eigenvectors and Eigenvalues:** Solve the generalized eigenvalue problem for $Sw{-}1Sb$ to find the eigenvectors and eigenvalues.

5. **Select the Top Eigenvectors:** Sort the eigenvalues in descending order and select the corresponding eigenvectors.

6. **Transform the Data:** Project the standardized data onto the new feature space defined by the selected eigenvectors.

**Benefits of LDA**

1. **Class Separability:** Maximizes class separability, making it particularly useful for classification tasks.

2. **Dimensionality Reduction:** Reduces the number of features while retaining the ability to distinguish between classes.

3. **Interpretability:** The resulting linear discriminants can be interpreted in terms of the original features.

**Drawbacks of LDA**

1. **Assumptions:** The assumptions of linearity, normality, and equal covariance may not hold for all datasets, potentially limiting the effectiveness of LDA.

2. **Overfitting:** With small sample sizes and a large number of features, LDA can overfit the training data.

3. **Scalability:** For very large datasets, computing the scatter matrices and solving the eigenvalue problem can be computationally intensive.

**Conclusion**

Linear Discriminant Analysis (LDA) is a powerful technique for dimensionality reduction in the context of supervised learning. By focusing on maximizing class separability, LDA transforms the data into a lower-dimensional space that retains the most dis criminative information. However, it is important to be mindful of the assumptions underlying LDA and to carefully select the number of linear discriminants to balance dimensionality reduction with information retention. When applied appropriately, LDA can significantly enhance the performance of classification models by simplifying the dataset and emphasizing the most relevant features.

## 6.2 DECISION TREE

A decision tree is a popular machine learning algorithm used for both classification and regression tasks. It is a tree-like model of decisions and their possible consequences, including chance event outcomes, resource costs, and utility. Here's a detailed overview of decision trees in machine learning:

**Structure of a Decision Tree**

1. **Root Node**: The topmost node in a decision tree that represents the entire dataset. Splitting starts from this node.

2. **Decision Nodes**: These are nodes where the data is split based on a feature attribute. Each decision node has two or more branches.

3. **Leaf Nodes (Terminal Nodes)**: These nodes represent the outcome or the final decision. No further splitting occurs at these nodes.

4. **Branches**: These represent the outcome of a test on an attribute and lead to the next node or result.

**How Decision Trees Work**

1. **Splitting**: The process of dividing a node into two or more sub-nodes based on a feature. The goal is to find splits that maximize the homogeneity of the target variable within each subset.

2. **Attribute Selection Measures**: To decide the best split, decision trees use attribute selection measures such as:

- **Gini Impurity**: Measures the likelihood of an incorrect classification of a new instance if it was randomly classified according to the distribution of class labels in the node.

- **Entropy (Information Gain)**: Measures the amount of information disorder or impurity in the dataset.

- **Variance Reduction**: Used in regression trees to measure the variance within the subsets and choose the split that reduces this variance the most.

3. **Stopping Criteria**: Determines when to stop splitting nodes. Common criteria include:

- Maximum depth of the tree

- Minimum number of samples in a node

- Minimum impurity decrease required to split

**Advantages of Decision Trees**

- **Interpretability**: Decision trees are easy to understand and interpret. They can be visualized, making it straightforward to understand the decision-making process.

- **Non-parametric**: They do not make any assumptions about the distribution of the data.

- **Handling Non-linear Data**: Capable of handling complex non-linear relationships between features and the target variable.

- **Feature Importance**: Decision trees provide a clear indication of which features are the most important for prediction.

**Disadvantages of Decision Trees**

- **Overfitting**: Decision trees can create overly complex trees that do not generalize well from the training data. Pruning techniques or setting maximum depth can help mitigate this.

- **Instability**: Small changes in the data can result in significantly different trees.

- **Bias**: Trees can be biased towards features with more levels or categories.

**Enhancements to Decision Trees**

To overcome some of the limitations of basic decision trees, several ensemble methods have been developed:

1. **Random Forests**: An ensemble method that builds multiple decision trees and merges them together to get a more accurate and stable prediction. Each tree is built from a random subset of features and data points, which reduces overfitting and improves generalization.

2. **Gradient Boosting Machines (GBM)**: Builds trees sequentially, each tree trying to correct the errors of the previous one. This technique focuses on improving the model's performance by optimizing a loss function.

3. **XGBoost**: An optimized version of GBM that is efficient, scalable, and often performs well in competitions. It includes additional features for handling missing data, regularization to prevent overfitting, and parallel processing.

**Applications of Decision Trees**

- **Classification**: Medical diagnosis, spam detection, credit scoring, and more.

- **Regression**: Predicting house prices, stock prices, and other continuous variables.

- **Feature Selection**: Identifying the most significant variables in a dataset.

**Example of a Decision Tree**

Here's a simple example to illustrate how a decision tree works. Suppose we want to classify whether a person should play tennis based on weather conditions. The features might include outlook (sunny, overcast, rainy), temperature (hot, mild, cool), humidity (high, normal), and wind (weak, strong). The decision tree would split these features to classify the days into "play" or "don't play" tennis categories, using the attribute selection measures to determine the splits.

In conclusion, decision trees are a fundamental and versatile tool in machine learning, valued for their simplicity and interpretability. Despite their limitations, they form the backbone of more complex ensemble methods that are widely used for a variety of predictive tasks.

**1. Decision Tree Algorithm with Near Miss Downsampling**

**Near Miss Downsampling** is a technique used to balance the dataset by under-sampling the majority class. For a daily drought prediction dataset, this involves reducing the number of non-drought days to match the number of drought days.

- **Purpose**: To handle class imbalance by reducing the number of majority class instances.

- **Implementation**: The Decision Tree classifier is trained on the balanced dataset after applying Near Miss.

- **Expected Outcome**: Improved performance in predicting drought days due to the balanced dataset, but potential loss of important information from the majority class.

**2. Decision Tree Algorithm with Near Miss Downsampling - Hyperparameter Tuning**

This involves the same Near Miss downsampling as above, but with the addition of hyperparameter tuning.

- **Purpose**: To optimize the Decision Tree's performance by finding the best combination of hyperparameters.

- **Implementation**: Grid Search or Random Search is used to explore different values for hyperparameters like **max_depth**, **min_samples_leaf**, and **max_features**.

- **Expected Outcome**: A more accurate and robust model compared to using default hyperparameters.

**3. Decision Tree Algorithm with SMOTE Upsampling**

**SMOTE (Synthetic Minority Over-sampling Technique)** creates synthetic instances for the minority class to balance the dataset.

- **Purpose**: To address class imbalance by increasing the number of drought days using synthetic data.

- **Implementation**: The Decision Tree classifier is trained on the balanced dataset after applying SMOTE.

- **Expected Outcome**: Better generalization and improved detection of drought days, leveraging additional synthetic data to balance the classes.

**4. Decision Tree Algorithm with Near Miss Downsampling and PCA**

This combines Near Miss downsampling with Principal Component Analysis (PCA).

- **Purpose**: To reduce class imbalance and dimensionality, removing less informative features.

- **Implementation**: After applying Near Miss, PCA is performed to transform the features into principal components, followed by training the Decision Tree.

- **Expected Outcome**: Enhanced model performance by reducing noise and overfitting, though some information may be lost during dimensionality reduction.

**5. Decision Tree Algorithm with SMOTE Upsampling and PCA**

This combines SMOTE upsampling with PCA.

- **Purpose**: To address class imbalance and reduce the feature space.

- **Implementation**: SMOTE is applied to generate synthetic samples, followed by PCA to reduce the dataset's dimensionality, and then training the Decision Tree.

- **Expected Outcome**: Improved model performance and efficiency, with the potential to capture more complex patterns through synthetic data and principal components.

**6. Decision Tree Algorithm with Near Miss Downsampling and LDA**

This combines Near Miss downsampling with Linear Discriminant Analysis (LDA).

- **Purpose**: To balance the dataset and maximize class separability.

- **Implementation**: Near Miss is applied, followed by LDA to transform features into a space that maximizes class separability, and then training the Decision Tree.

- **Expected Outcome**: Better class discrimination and prediction accuracy, particularly for distinguishing between drought and non-drought days.

## 7. Decision Tree Algorithm with SMOTE Upsampling and LDA

This combines SMOTE upsampling with LDA.

- **Purpose**: To balance the dataset and enhance class separability.

- **Implementation**: SMOTE generates synthetic samples, followed by LDA to project data into a space that maximizes class separation, and then training the Decision Tree.

- **Expected Outcome**: Enhanced model accuracy and ability to distinguish between classes due to the combination of synthetic data and optimized feature space.

## 8. Decision Tree Algorithm without Resampling

The Decision Tree is trained on the original, imbalanced dataset.

- **Purpose**: To use the original data distribution without any resampling techniques.

- **Implementation**: The classifier is trained directly on the dataset without any modifications.

- **Expected Outcome**: Potential bias towards the majority class (non-drought days), with possible poor performance in predicting drought days.

**9. Decision Tree Algorithm without Resampling - Hyperparameter Tuning**

The same as above, but with hyperparameter tuning to optimize the model.

- **Purpose**: To improve the Decision Tree's performance using hyperparameter tuning.

- **Implementation**: Grid Search or Random Search is used to find the best hyperparameters for the classifier.

- **Expected Outcome**: Improved model performance compared to using default hyperparameters, though class imbalance might still affect predictions.

**10. Decision Tree Algorithm without Resampling - Setting the Right Hyperparameters**

This involves manually setting hyperparameters based on domain knowledge or previous experiments.

- **Purpose**: To optimize the model using specific hyperparameters tailored to the dataset.

- **Implementation**: Key hyperparameters such as **max_depth**, **min_samples_split**, **min_samples_leaf**, and **max_features** are set based on informed decisions.

- **Expected Outcome**: Improved performance tailored to the specific characteristics of the drought prediction dataset, balancing between complexity and overfitting.

In summary, each approach involves different techniques to handle class imbalance, feature reduction, and model optimization. The choice of algorithm depends on the specific requirements of the drought prediction task and the characteristics of the dataset.

fig 6.2.1 Multiclass ROC curve for decision tree with near miss downsampling

Multiclass ROC curve for Decision Tree with SMOTE Upsampling and PCA



Multiclass ROC curve for Decision Tree with Near Miss Downsampling and LDA



Multiclass ROC curve for Decision Tree with SMOTE Upsampling and LDA

fig 6.2.2 Multiclass ROC curve for decision tree with SMOTE

Multiclass ROC curve for Decision Tree without resampling



Multiclass ROC curve for Decision Tree without resampling - After Hyperparameter Tuning

fig 6.2.3 Multiclass ROC curve for decision tree without resampling

# 6.3 K-Nearest Neighbors (KNN)

**Definition**

K-Nearest Neighbors (KNN) is a simple, non-parametric, and instance-based learning algorithm used for both classification and regression tasks in machine learning. It operates on the principle that similar data points are likely to be near each other. The algorithm does not make any assumptions about the underlying data distribution, which means it can be applied to a variety of problems.

**How KNN Works**

1. **Training Phase**:

    KNN doesn't involve a traditional training phase. Instead, it stores all the training data points. This is why KNN is considered a lazy learning algorithm because it defers the decision-making process until a query is made.

2. **Prediction Phase**:

    - **Classification**:

        1. Compute the distance between the query point (the new data point for which we want to predict the label) and all the points in the training data using a chosen distance metric (e.g., Euclidean distance, Manhattan distance).

        2. Identify the K-nearest neighbors to the query point based on the computed distances.

3. Perform a majority vote among the K-nearest neighbors to determine the class label of the query point. Each neighbor votes for its own class, and the class with the most votes is chosen as the prediction.

- **Regression**:

1. Compute the distance between the query point and all the points in the training data.

2. Identify the K-nearest neighbors.

3. Calculate the average (or weighted average) of the target values of the K-nearest neighbors to predict the value for the query point. This prediction can be the mean of the values or a weighted average where closer neighbors have more influence.

**Distance Metrics**

The choice of distance metric can significantly affect the performance of KNN. Common distance metrics include:

- **Euclidean Distance**: $\sum i=1(xi-yi)2\sum i=1n(xi-yi)2$, which is the straight-line distance between two points in Euclidean space.

- **Manhattan Distance**: $\sum i=1n|xi-yi|\sum i=1n|xi-yi|$, also known as L1 distance or taxicab distance, which measures the distance along axes at right angles.

- **Minkowski Distance**: $(\sum i=1n|xi-yi|p)1/p(\sum i=1n|xi-yi|p)1/p$, where $pp$ is a parameter that determines the type of distance. When $p=2p=2$, Minkowski distance is

equivalent to Euclidean distance, and when $p=1$, it is equivalent to Manhattan distance.

- **Hamming Distance**: Used for categorical variables, it measures the number of positions at which the corresponding elements are different.

**Choosing the Value of K**

Choosing the right value of K is crucial for the performance of the KNN algorithm:

- A small value of K (e.g., K=1) can be sensitive to noise and outliers, leading to overfitting.

- A large value of K provides a smoother decision boundary but may overlook the finer details in the data, potentially leading to underfitting.

- The optimal value of K is usually determined through cross-validation, where the training data is split into multiple subsets to evaluate the model's performance on different values of K.

**Feature Scaling**

Feature scaling is an essential preprocessing step for KNN. Since KNN relies on distance calculations, features with larger ranges can disproportionately influence the distance metric. Common scaling techniques include:

- **Min-Max Scaling**: Rescales the feature to a range of [0, 1].

- **Standardization (Z-score normalization)**: Rescales the feature to have a mean of 0 and a standard deviation of 1.

**Handling Missing Values**

KNN can handle missing values by:

- **Imputation**: Replacing missing values with the mean, median, or mode of the feature.

- **KNN-based Imputation**: Predicting the missing value based on the values of the K-nearest neighbors.

**Advantages of KNN**

- **Simplicity**: Easy to understand and implement, making it a good starting point for beginners in machine learning.

- **No Assumptions**: Does not assume any specific data distribution, making it versatile and applicable to a wide range of problems.

- **Flexibility**: Can be used for both classification and regression tasks, providing a unified approach to different types of problems.

- **Adaptability**: New data can be added seamlessly, as KNN is an instance-based learner and doesn't require retraining.

- **Interpretable**: Results are often interpretable, as they are based on the actual examples from the training set.

**Disadvantages of KNN**

- **Computationally Expensive**: High computational cost during prediction, especially with large datasets, since it requires calculating distances to all training points.

- **Memory Intensive**: Requires storing all training data, leading to high memory usage, particularly with large datasets.

- **Performance Degrades with Dimensionality**: High-dimensional data can lead to the curse of dimensionality, where the distance metric becomes less effective, and all points tend to become equidistant.

- **Sensitive to Irrelevant Features**: Performance can be significantly affected if the input data has many irrelevant features. Feature selection and normalization are critical steps before applying KNN.

- **Imbalanced Data**: KNN can struggle with imbalanced datasets where some classes are underrepresented.

**Applications of KNN**

1. **Pattern Recognition**:

   - **Handwriting Recognition**: Recognizing handwritten digits or characters by comparing them to known examples.

   - **Image Classification**: Classifying images into categories based on visual similarity to labeled images in the training set.

2. **Recommendation Systems**:

   - **Collaborative Filtering**: Recommending products to users by finding other users with similar preferences and suggesting items they liked.

3. **Medical Diagnosis**:

   - **Disease Prediction**: Predicting disease outcomes based on patient data, such as symptoms, test results, and medical history.

4. **Anomaly Detection**:

- **Fraud Detection**: Identifying unusual patterns in transaction data that may indicate fraudulent activity.

5. **Customer Behaviour Analysis**:

- **Churn Prediction**: Identifying customers who are likely to stop using a service based on their usage patterns.

6. **Finance**:

- **Credit Scoring**: Predicting the creditworthiness of individuals based on historical data.

**KNN Algorithm without resampling:**

It is the basic implementation of the KNN algorithm without any additional techniques for model evaluation or improvement.

The following code snippet demonstrates the application of the K-Nearest Neighbors (KNN) algorithm for classification tasks using Python's scikit-learn library. KNN is a simple yet effective non-parametric algorithm used for both classification and regression tasks. It operates on the principle of similarity, where the class or label of a data point is determined by the majority class among its nearest neighbors in the feature space.

In this code, we initialize a KNN classifier with specific parameters, train it using labeled training data, and then utilize the trained classifier to predict labels for unseen test data. The dataset used for this task typically consists of feature vectors representing input data samples and corresponding target labels or classes.

The KNN algorithm's flexibility and ease of implementation make it suitable for various classification tasks, including but not limited to image recognition, natural language processing, and medical diagnosis. Understanding its underlying principles and workflow is fundamental for leveraging its capabilities effectively in machine learning projects.

This introduction sets the stage for understanding the subsequent code snippet, which illustrates the step-by-step process of utilizing the KNN algorithm for classification tasks, from initialization to prediction.

```
knn_classifier = KNeighborsClassifier(n_neighbors=5, p=2, metric='minkowski')
knn_classifier.fit(X_train, y_train)
y_pred_knn = knn_classifier.predict(X_test)
```

1. **Initializing the KNN Classifier:**

   **knn_classifier = KNeighborsClassifier(n_neighbors=5, p=2, metric='minkowski')**

   - In this line, we create an instance of the KNeighborsClassifier class from the scikit-learn library. This class implements the K-Nearest Neighbors algorithm for classification tasks.

   - We specify several parameters:

     - **n_neighbors=5**: This parameter determines the number of neighbors considered when making predictions. Here, we set it to 5, meaning the algorithm will consider the five closest data points to the query point.

- **p=2**: This parameter indicates the power parameter for the Minkowski distance metric. When **p=2**, the Minkowski distance becomes the Euclidean distance, which is commonly used in KNN.

- **metric='minkowski'**: Here, we specify the distance metric used by the algorithm. In this case, it's the Minkowski distance, which is a generalization of both the Euclidean distance (**p=2**) and the Manhattan distance (**p=1**).

## 2.Training the KNN Classifier:

**knn_classifier.fit(X_train, y_train)**

- This line trains the KNN classifier using the provided training data.

- **X_train** represents the feature matrix containing the training data, where each row corresponds to a sample and each column corresponds to a feature.

- **y_train** represents the target labels corresponding to the training samples. Each element in **y_train** is the label or class that the corresponding sample in **X_train** belongs to.

- During training, the classifier learns patterns in the feature space that map input features (**X_train**) to their corresponding labels (**y_train**).

**3.Making Predictions:**

**y_pred_knn = knn_classifier.predict(X_test)**

- Once the classifier is trained, we use it to predict the labels for new, unseen data.

- **X_test** represents the feature matrix containing the test data, similar in structure to **X_train**. However, this data was not used during training.

- The **predict** method applies the trained classifier to the test data (**X_test**) and returns an array of predicted labels (**y_pred_knn**), where each element corresponds to the predicted label for the corresponding sample in **X_test**.

- These predicted labels are generated based on the learned patterns from the training data. The algorithm calculates distances between the test samples and the training samples, identifies the **k** nearest neighbors, and assigns the class label based on the majority class among those neighbors (in the case of classification tasks).

In summary, this code snippet demonstrates the typical workflow for applying the KNN algorithm in a classification task: initializing the classifier with specific parameters, training it using labeled training data, and then using the trained classifier to predict labels for unseen test data.

**Performance of KNN Algorithm without resampling**

```
print('Performance of KNN Algorithm without resampling:\n')
print(confusion_matrix(y_test, y_pred_knn))
print(classification_report(y_test, y_pred_knn))
```

```
print('Accuracy:',accuracy_score(y_test, y_pred_knn))

print('Precision:',precision_score(y_test, y_pred_knn, average='weighted'))

print('Recall:',recall_score(y_test, y_pred_knn, average='weighted'))

print('F1 Score:',f1_score(y_test, y_pred_knn, average='weighted'))

print('Cohen Kappa Score:',cohen_kappa_score(y_test, y_pred_knn))
```

This code snippet delves into the evaluation of a K-Nearest Neighbors (KNN) classification model's performance on a test dataset, specifically focusing on the scenario where resampling techniques like cross-validation are not employed.

**Detailed Explanation of the Code:**

1. **Printing the Performance Heading:**

   **print('Performance of KNN Algorithm without resampling:\n')**

   - This line prints a heading to introduce the performance metrics that follow. It helps in organizing the output and making it clear that the subsequent results pertain to the KNN algorithm's performance evaluation without resampling techniques.

2. **Confusion Matrix:**

   **print(confusion_matrix(y_test, y_pred_knn))**

   - The confusion matrix is a key tool for evaluating the performance of a classification model. It provides a breakdown of the actual versus predicted classifications.

   - Structure of the confusion matrix:

- **True Positives (TP):** The number of instances correctly classified as positive.

- **True Negatives (TN):** The number of instances correctly classified as negative.

- **False Positives (FP):** The number of instances incorrectly classified as positive.

- **False Negatives (FN):** The number of instances incorrectly classified as negative.

- For example, if the model is predicting whether it will rain, TP would be the number of times it correctly predicted rain when it actually rained.

3. **Classification Report:**

**print(classification_report(y_test, y_pred_knn))**

- The classification report provides a detailed overview of the classifier's performance, including metrics for each class.

- **Precision:** The ratio of correctly predicted positive observations to the total predicted positives. Precision answers the question, "Of all the instances predicted as positive, how many were actually positive?"

- **Recall (Sensitivity):** The ratio of correctly predicted positive observations to all actual positives. Recall answers the question, "Of all the actual positive instances, how many were correctly predicted?"

- **F1 Score:** The harmonic mean of precision and recall. It provides a single metric that balances the precision and recall trade-off.

- **Support:** The number of actual occurrences of each class in the dataset. It gives context to the precision, recall, and F1 scores.

4. **Accuracy:**

**print('Accuracy:', accuracy_score(y_test, y_pred_knn))**

- Accuracy is a simple yet fundamental metric that indicates the proportion of correctly classified instances among the total instances.

- Formula: $(TP+TN)/(TP+TN+FP+FN)$ $(TP+TN)/(TP+TN+FP+FN)$

- While accuracy provides a general overview, it can be misleading in the presence of class imbalance. For instance, if one class is dominant, a high accuracy could simply reflect the classifier's bias towards that class.

5. **Precision (Weighted):**

**print('Precision:', precision_score(y_test, y_pred_knn, average='weighted'))**

- Precision focuses on the accuracy of positive predictions.

- The **average='weighted'** parameter ensures that precision is calculated for each class and then averaged, with each class's contribution weighted by its size. This is crucial in datasets with class imbalance, as it prevents dominant classes from skewing the overall precision.

- Formula for precision for a class: $\text{Precision} = \frac{TP}{TP+FP}$ $\text{Precision} = \frac{TP}{TP+FP}$

6. **Recall (Weighted):**

**Print('Recall:', recall_score(y_test, y_pred_knn, average='weighted'))**

- Recall indicates the ability of the model to capture all positive instances.

- The **average='weighted'** parameter ensures a balanced evaluation across all classes by weighting recall according to the size of each class.

- Formula for recall for a class: $\text{Recall} = \frac{TP}{TP+FN}$

- Recall is particularly important in scenarios where missing positive instances is costly, such as in medical diagnoses.

7. **F1 Score (Weighted):**

    **print('F1 Score:', f1_score(y_test, y_pred_knn, average='weighted'))**

- The F1 score is the harmonic mean of precision and recall, providing a single measure that balances both metrics.

- The **average='weighted'** parameter ensures the F1 score is calculated for each class and then averaged with weights based on class size.

- Formula for F1 score for a class:
    $F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

- The F1 score is especially useful when the class distribution is imbalanced, offering a more nuanced view of performance than accuracy alone.

8. **Cohen Kappa Score:**

    **print('Cohen Kappa Score:', cohen_kappa_score(y_test, y_pred_knn))**

- The Cohen Kappa score measures the agreement between the true labels and the predicted labels, adjusting for the agreement that could occur by chance.

- It ranges from -1 to 1:

- 1 indicates perfect agreement.

- 0 indicates no agreement better than random chance.

- Negative values indicate worse than random agreement.

- Formula: $\kappa = \frac{P_o - P_e}{1 - P_e}$, where $P_o$ is the observed agreement and $P_e$ is the expected agreement by chance.

- This metric is particularly valuable for imbalanced datasets, as it accounts for the possibility of random agreement, providing a more robust measure of the classifier's performance.

**Performance of KNN Algorithm without resampling:**

```
[[279710 14593  2815  1353   597   177]
 [ 28411 45481  7206  1376   469   123]
 [ 10164 10737 25742  4932   799   175]
 [  5054  2684  6776 17005  2722   252]
 [  2403   907  1448  3730  8832  1100]
 [   753   272   328   463  1404  3875]]
```

precision  recall  f1-score  support

| | | | | |
|---|---|---|---|---|
| 0 | 0.86 | 0.93 | 0.89 | 299245 |
| 1 | 0.61 | 0.55 | 0.58 | 83066 |
| 2 | 0.58 | 0.49 | 0.53 | 52549 |
| 3 | 0.59 | 0.49 | 0.54 | 34493 |
| 4 | 0.60 | 0.48 | 0.53 | 18420 |
| 5 | 0.68 | 0.55 | 0.61 | 7095 |
| | | | | |
| accuracy | | | 0.77 | 494868 |
| macro avg | 0.65 | 0.58 | 0.61 | 494868 |
| weighted avg | 0.75 | 0.77 | 0.76 | 494868 |

Accuracy: 0.7691849139568532

Precision: 0.7549562894195875

Recall: 0.7691849139568532

F1 Score: 0.7597211442125624

Cohen Kappa Score: 0.588463390888752

The information provided represents the performance metrics of a K-Nearest Neighbors (KNN) algorithm applied to a classification problem without any resampling techniques. Here's a detailed explanation of each part:

**Confusion Matrix**

The confusion matrix shows the number of correct and incorrect predictions made by the KNN algorithm for each class. It is a 6x6 matrix corresponding to 6 classes (0 to 5).

- **Rows** represent the actual classes.

- **Columns** represent the predicted classes.

  Each entry $(i,)(i,j)$ in the matrix represents the number of samples of class $ii$ that were predicted as class $jj$.

  For example:

- The value at (0,0) is 279710, meaning 279710 samples of class 0 were correctly classified as class 0.

- The value at (0,1) is 14593, meaning 14593 samples of class 0 were incorrectly classified as class 1.

**Classification Report**

The classification report provides key metrics for each class:

- **Precision**: The ratio of correctly predicted positive observations to the total predicted positives. It is calculated as $TPTP+FPTP+FPTP$.

  - Class 0: 0.86

  - Class 1: 0.61

- Class 2: 0.58

- Class 3: 0.59

- Class 4: 0.60

- Class 5: 0.68

- **Recall**: The ratio of correctly predicted positive observations to the all observations in the actual class. It is calculated as $\frac{TP}{TP+FN}$.

  - Class 0: 0.93

  - Class 1: 0.55

  - Class 2: 0.49

  - Class 3: 0.49

  - Class 4: 0.48

  - Class 5: 0.55

- **F1-Score**: The weighted average of Precision and Recall. It is calculated as $\frac{2 \times Precision \times Recall}{Precision + Recall}$.

  - Class 0: 0.89

  - Class 1: 0.58

  - Class 2: 0.53

  - Class 3: 0.54

  - Class 4: 0.53

- Class 5: 0.61

- **Support**: The number of actual occurrences of each class in the dataset.

  - Class 0: 299245

  - Class 1: 83066

  - Class 2: 52549

  - Class 3: 34493

  - Class 4: 18420

  - Class 5: 7095

  **Overall Metrics**

- **Accuracy**: The ratio of correctly predicted observations to the total observations. It is calculated as $\frac{TP+TN}{Total}Total TP+TN$. For this model, the accuracy is 0.7692, meaning the model correctly predicts approximately 77% of the samples.

- **Precision (weighted avg)**: A weighted average of precision across all classes, considering the number of instances for each class.

  - Weighted average precision: 0.75

- **Recall (weighted avg)**: A weighted average of recall across all classes.

  - Weighted average recall: 0.77

- **F1 Score (weighted avg)**: A weighted average of F1 scores across all classes.

  - Weighted average F1 Score: 0.76

- **Cohen Kappa Score**: A statistic that measures inter-rater agreement for categorical items. It considers the possibility of the agreement occurring by chance. The score ranges from -1 to 1, where:

  - 0 indicates agreement equivalent to chance.

  - 1 indicates perfect agreement.

  - Negative values indicate agreement worse than chance.

  - For this model, the Kappa score is 0.588, indicating a moderate agreement between the predictions and the actual labels.

  In summary, the KNN algorithm performs reasonably well on class 0, but its performance degrades for other classes, particularly those with fewer samples. The overall accuracy and weighted averages suggest the model is moderately effective, with room for improvement, particularly in handling class imbalances and increasing precision and recall for the less frequent classes.

**Receiver Operating Characteristic (ROC) curves**

```
fpr = dict()

tpr = dict()

thresh = dict()


for i in range(6):
    fpr[i], tpr[i], thresh[i] = roc_curve(y_test, y_pred_knn, pos_label=i)


plt.plot(fpr[0], tpr[0], linestyle='--',color='orangered', label='Class 0 vs Rest')

plt.plot(fpr[1], tpr[1], linestyle='--',color='green', label='Class 1 vs Rest')
```

```
plt.plot(fpr[2], tpr[2], linestyle='--',color='blue', label='Class 2 vs Rest')

plt.plot(fpr[3], tpr[3], linestyle='--',color='yellow', label='Class 3 vs Rest')

plt.plot(fpr[4], tpr[4], linestyle='--',color='purple', label='Class 4 vs Rest')

plt.plot(fpr[5], tpr[5], linestyle='--',color='magenta', label='Class 5 vs Rest')


plt.title('Multiclass ROC curve for KNN without resampling')

plt.xlabel('False Positive Rate')

plt.ylabel('True Positive rate')

plt.legend(loc='best')'

plt.savefig('Multiclass ROC curve for KNN without resampling',dpi=300)
```

The provided code snippet computes the Receiver Operating Characteristic (ROC) curves for each class in a multiclass classification problem using a K-Nearest Neighbors (KNN) algorithm without any resampling techniques. Let's break down the information:

**ROC Curve**

- The ROC curve is a graphical representation of the model's performance across various threshold settings.

- It plots the True Positive Rate (TPR) against the False Positive Rate (FPR) for different classification thresholds.

- Each curve in the plot represents a class versus the rest of the classes. For example, "Class 0 vs Rest" compares the performance of the model in classifying class 0 against all other classes combined.

- A diagonal line from the bottom left to the top right (line y=x) represents random guessing, while a curve above this line indicates better-than-random performance.

- The area under the ROC curve (AUC) summarizes the model's performance across all thresholds. A higher AUC indicates better overall performance.

**Code Explanation**

- The code iterates through each class and computes the FPR, TPR, and thresholds using the **roc_curve** function from scikit-learn.

- Each class's FPR, TPR, and thresholds are then plotted using **plt.plot**.

- Different colors and line styles are used for better visualization and to differentiate between the ROC curves for each class.

- Axes labels, title, and legend are added to the plot for clarity.

- Finally, the plot is saved as an image file named "Multiclass ROC curve for KNN without resampling".

**Interpretation**

- The plot visualizes the trade-off between true positive rate and false positive rate for each class.

- Ideally, we want the curve to be closer to the top-left corner, indicating higher true positive rates and lower false positive rates across different thresholds.

- Comparing the AUC values for different classes can provide insights into which classes the model performs better on.

- By examining the ROC curves, we can assess the model's discriminative ability across different classes and identify areas where it may need improvement.

In summary, the ROC curve analysis provides a comprehensive understanding of the model's performance across multiple classes in a multiclass classification problem, allowing for insights into its discriminatory power and areas for potential enhancement.



Fig 6.3.1 Multiclass ROC curve for KNN without resampling

**KNN Algorithm with SMOTE upsampling**

**SMOTE**

(Synthetic Minority Over-sampling Technique) is a method used to address class imbalance in datasets. It works by generating synthetic samples for the minority class to balance the class distribution.

**Advantages:**

1. **Mitigates Class Imbalance**: SMOTE helps to alleviate the problem of class imbalance by creating synthetic instances of the minority class, thus balancing the class distribution in the dataset.

2. **Preserves Information**: The synthetic samples are created by interpolating between existing minority class instances, which helps to preserve the information contained in the original data.

3. **Reduced Overfitting**: By increasing the number of minority class samples, SMOTE can reduce the likelihood of overfitting that may occur when training on imbalanced datasets.

**Disadvantages:**

1. **Potential Overfitting**: While SMOTE can help reduce overfitting, it may also introduce some level of overfitting, especially if the synthetic samples are not representative of the true underlying data distribution.

2. **Increased Computational Complexity**: Generating synthetic samples using SMOTE can increase the computational complexity of the training process, particularly for large datasets.

3. **Sensitive to Noise**: SMOTE is sensitive to noisy data and may generate synthetic samples in regions of feature space that contain noise, leading to decreased classification performance.

**Other Considerations:**

1. **Parameter Tuning**: SMOTE typically has parameters that need to be tuned, such as the number of nearest neighbors to consider when generating synthetic samples.

2. **Combined Approaches**: SMOTE is often used in combination with other techniques, such as undersampling the majority class or using different sampling ratios, to further improve classification performance.

**Use Cases:**

- SMOTE is commonly used in classification tasks where there is a significant class imbalance, such as fraud detection, medical diagnosis, and anomaly detection.

- It can be applied with various machine learning algorithms, including decision trees, support vector machines, and k-nearest neighbors.

In summary, SMOTE is a powerful technique for addressing class imbalance in datasets by generating synthetic samples of the minority class. While it offers several advantages, such as mitigating class imbalance and preserving information, it also has limitations and considerations that need to be taken into account when applying it in practice.

**KNN Algorithm with SMOTE Upsampling**

It is a machine learning technique used to address class imbalance in datasets while leveraging the K-Nearest Neighbors (KNN) algorithm for classification tasks.

- Combining KNN with SMOTE involves first applying SMOTE to upsample the minority class(es) in the dataset to address class imbalance.

- Synthetic instances are generated for the minority class using SMOTE, effectively increasing its representation in the dataset.

- After upsampling, the KNN algorithm is applied to the balanced dataset to classify new instances.

- By leveraging SMOTE to balance the dataset, KNN can make more accurate predictions, especially for minority class instances, which may be poorly represented in the original dataset.

**Advantages:**

1. **Improved Performance for Imbalanced Data**:

   - By addressing class imbalance with SMOTE, KNN can make more accurate predictions, particularly for minority class instances.

2. **Preservation of Local Structure**:

   - SMOTE preserves the local structure of the minority class, which is essential for KNN's performance, as it relies on the similarity of data points.

3. **Simplicity of Implementation**:

   - Both KNN and SMOTE are relatively simple algorithms to implement, making the combination straightforward and accessible.

**Disadvantages:**

1. **Increased Computational Complexity**:

   - Generating synthetic samples with SMOTE can increase the dataset's size, leading to longer training times and higher memory requirements.

2. **Sensitivity to Noise**:

- SMOTE is sensitive to noise in the data and may generate synthetic samples in regions containing noise, potentially affecting model performance.

3. **Potential Overfitting**:

- Care must be taken to avoid overfitting when using SMOTE, as generating synthetic samples indiscriminately can lead to overfitting, especially if the dataset is already noisy.

In summary, KNN Algorithm with SMOTE Upsampling is a powerful technique for addressing class imbalance in datasets, offering improved performance and preserving the local structure of the data. However, it's essential to consider its potential drawbacks, such as increased computational complexity and sensitivity to noise, when applying it in practice.

This code snippet implements a K-Nearest Neighbors (KNN) classifier with Synthetic Minority Over-sampling Technique (SMOTE) upsampling.

```
knn_classifier_SMOTE = KNeighborsClassifier(n_neighbors=1, p=2, metric='minkowski')
knn_classifier_SMOTE.fit(X_train_ures_SMOTE, y_train_ures_SMOTE)
y_pred_knn_SMOTE = knn_classifier_SMOTE.predict(X_test)
```

**Code Explanation:**

**knn_classifier_SMOTE = KNeighborsClassifier(n_neighbors=1, p=2, metric='minkowski')**

- This line creates an instance of the KNeighborsClassifier class from the scikit-learn library.

- The **n_neighbors** parameter is set to 1, meaning the algorithm considers only the closest neighbor when making predictions.

- The **p** parameter is set to 2, which indicates the Euclidean distance metric (Minkowski distance with p=2).

- The **metric** parameter specifies the distance metric used for calculating distances between points, in this case, 'minkowski'.

**knn_classifier_SMOTE.fit(X_train_ures_SMOTE, y_train_ures_SMOTE)**

- This line fits (trains) the KNN classifier on the training data after SMOTE upsampling.

- **X_train_ures_SMOTE** contains the feature vectors of the training instances after SMOTE upsampling.

- **y_train_ures_SMOTE** contains the corresponding class labels after SMOTE upsampling.

**y_pred_knn_SMOTE = knn_classifier_SMOTE.predict(X_test)**

- This line predicts the class labels for the test data using the trained KNN classifier.

- **X_test** contains the feature vectors of the test instances.

- The predicted class labels are stored in the variable **y_pred_knn_SMOTE**.

**Summary:**

In summary, this code snippet demonstrates how to train a KNN classifier with SMOTE upsampling to address class imbalance in a dataset. It first initializes the KNN classifier with

specified parameters, then fits the classifier to the training data after applying SMOTE, and finally predicts the class labels for the test data. This approach can help improve the classifier's performance, especially on imbalanced datasets, by ensuring adequate representation of the minority class during training.

```python
pickle.dump(knn_classifier_SMOTE, open('knn_classifier_SMOTE.pkl', 'wb'))
```

This line of code saves the trained KNN classifier with SMOTE upsampling to a file named 'knn_classifier_SMOTE.pkl'. By doing so, we can later load this saved model using **pickle.load()** and use it to make predictions on new data without the need to retrain the model from scratch. This approach helps in saving time and resources, especially in scenarios where model training is computationally expensive or time-consuming.

- **pickle.dump()** function serializes the **knn_classifier_SMOTE** object.

- **open('knn_classifier_SMOTE.pkl', 'wb')** opens a file named 'knn_classifier_SMOTE.pkl' in binary write mode ('wb').

  - The '.pkl' extension is commonly used for pickle files.

  - 'wb' mode indicates that the file will be opened for writing in binary mode.

- The serialized form of the **knn_classifier_SMOTE** object is then written to the specified file.

**Performance of KNN Algorithm with SMOTE Upsampling**

print('Performance of KNN Algorithm with SMOTE Upsampling:**\n**')

print(confusion_matrix(y_test, y_pred_knn_SMOTE))

print(classification_report(y_test, y_pred_knn_SMOTE))

print('Accuracy:',accuracy_score(y_test, y_pred_knn_SMOTE))

print('Precision:',precision_score(y_test, y_pred_knn_SMOTE, average='weighted'))

print('Recall:',recall_score(y_test, y_pred_knn_SMOTE, average='weighted'))

print('F1 Score:',f1_score(y_test, y_pred_knn_SMOTE, average='weighted'))

print('Cohen Kappa Score:',cohen_kappa_score(y_test, y_pred_knn_SMOTE))

This code snippet is similar to the previous one but focuses on evaluating the performance of a KNN (K-Nearest Neighbors) classifier with SMOTE (Synthetic Minority Over-sampling Technique) upsampling applied. Let's go through it:

**Printing Performance Metrics:**

**print('Performance of KNN Algorithm with SMOTE Upsampling:\n')**

- This line prints a header indicating that the following metrics are related to the performance of the KNN algorithm with SMOTE upsampling applied.

**Confusion Matrix:**

**print(confusion_matrix(y_test, y_pred_knn_SMOTE))**

- This line prints the confusion matrix for the KNN classifier with SMOTE upsampling.

- The confusion matrix provides information about the classifier's true positive, true negative, false positive, and false negative predictions on the test data.

**Classification Report:**

**print(classification_report(y_test, y_pred_knn_SMOTE))**

- This line prints a comprehensive report of classification metrics, such as precision, recall, F1-score, and support, for each class, as well as the average values across all classes.

- It gives insights into the classifier's performance for each class and overall with SMOTE upsampling.

**Accuracy, Precision, Recall, F1 Score, Cohen's Kappa Score:**

**print('Accuracy:', accuracy_score(y_test, y_pred_knn_SMOTE)) print('Precision:', precision_score(y_test, y_pred_knn_SMOTE, average='weighted')) print('Recall:', recall_score(y_test, y_pred_knn_SMOTE, average='weighted')) print('F1 Score:', f1_score(y_test, y_pred_knn_SMOTE, average='weighted')) print('Cohen Kappa Score:', cohen_kappa_score(y_test, y_pred_knn_SMOTE))**

- These lines print individual performance metrics such as accuracy, precision, recall, F1 score, and Cohen's kappa score for the KNN classifier with SMOTE upsampling.

- They are calculated using scikit-learn functions (**accuracy_score**, **precision_score**, **recall_score**, **f1_score**, and **cohen_kappa_score**) based on the true class labels (**y_test**) and predicted class labels (**y_pred_knn_SMOTE**).

**Summary:**

This code snippet evaluates the performance of a KNN classifier with SMOTE upsampling applied. It provides insights into how well the classifier performs in terms of classification accuracy, precision, recall, F1 score, and Cohen's kappa score, as well as a detailed breakdown of its performance using a confusion matrix and classification report. Evaluating the classifier with SMOTE upsampling helps assess its performance on a balanced dataset, which may lead to better handling of class imbalance issues compared to evaluating it without any resampling technique.

Accuracy: 0.7953

Precision: 0.8018

Recall: 0.7953

F1 Score: 0.7982

Cohen Kappa Score: 0.6578

**KNN Algorithm with Near Miss downsampling:**

**Definition:**

KNN Algorithm with Near Miss downsampling is a machine learning approach that combines the K-Nearest Neighbors (KNN) algorithm with the Near Miss downsampling technique to address class imbalance in datasets.

**Near Miss Downsampling:**

- Near Miss is a technique used to address class imbalance by undersampling the majority class instances that are close to the minority class instances in the feature space.
- It selects the majority class instances whose average distance to the minority class instances is the smallest.

**Advantages:**

1. Balancing Class Distribution:
   - Near Miss downsampling helps to balance the class distribution by reducing the number of majority class instances while retaining the essential information from the dataset.
2. Improved Generalization:
   - By reducing the dominance of the majority class, KNN with Near Miss downsampling can lead to better generalization and classification performance, especially for imbalanced datasets.
3. Computational Efficiency:

- Near Miss downsampling reduces the size of the dataset by removing majority class instances, leading to improved computational efficiency during model training.

**Disadvantages:**

1. Loss of Information:
   - Downsampling techniques like Near Miss may lead to a loss of information, particularly if the majority class instances removed contain relevant information for the classification task.

2. Risk of Overfitting:
   - Removing majority class instances may increase the risk of overfitting, especially if the remaining dataset is small or not representative of the underlying data distribution.

3. Sensitivity to Sampling Strategy:
   - The performance of KNN with Near Miss downsampling can be sensitive to the choice of sampling strategy and parameters, which may require careful tuning for optimal results.

**Other Considerations:**

1. Parameter Tuning:
   - Near Miss downsampling often involves parameters such as the number of neighbors to consider or the version of the technique to use, which may need to be tuned based on the dataset characteristics.

2. Combined Approaches:

- KNN with Near Miss downsampling can be combined with other techniques such as feature engineering, algorithm ensemble methods, or different resampling strategies to further improve classification performance.

**Use Cases:**

- KNN with Near Miss downsampling is commonly used in classification tasks where there is a significant class imbalance, such as fraud detection, medical diagnosis, or anomaly detection.
- It can be applied with various machine learning algorithms, but KNN is particularly suitable due to its reliance on the local structure of the data.

In summary, KNN Algorithm with Near Miss downsampling is a valuable approach for addressing class imbalance in datasets, offering improved classification performance and computational efficiency. However, it's essential to consider potential drawbacks such as information loss and sensitivity to sampling strategy when applying this technique in practice.

**Implementing a K-Nearest Neighbors (KNN) classifier using the scikit-learn library and applying Near Miss downsampling**

```
knn_classifier_NM = KNeighborsClassifier(n_neighbors=1, p=2, metric='minkowski')
knn_classifier_NM.fit(X_train_dres_nm, y_train_dres_nm)
y_pred_knn_NM = knn_classifier_NM.predict(X_test)
```

In many machine learning problems, especially in scenarios like fraud detection, medical diagnosis, or anomaly detection, datasets often suffer from class imbalance, where one class

(usually the minority class) has significantly fewer instances than the others. This imbalance can lead to biased models that perform poorly on the minority class. To address this issue, various techniques, including Near Miss downsampling, can be used to rebalance the dataset and improve the performance of machine learning models such as K-Nearest Neighbors (KNN).

**Code Explanation:**

**knn_classifier_NM = KNeighborsClassifier(n_neighbors=1, p=2, metric='minkowski')**

- This line initializes a K-Nearest Neighbors (KNN) classifier with specific parameters using the scikit-learn library.
- **n_neighbors** is set to 1, indicating that the algorithm considers only the closest neighbor when making predictions.
- **p** is set to 2, which signifies the Euclidean distance metric (Minkowski distance with p=2).
- **metric** parameter specifies the distance metric used for calculating distances between points, in this case, 'minkowski'.

**knn_classifier_NM.fit(X_train_dres_nm, y_train_dres_nm)**

- This line fits (trains) the KNN classifier on the training data after Near Miss downsampling.
- **X_train_dres_nm** contains the feature vectors of the training instances after Near Miss downsampling.
- **y_train_dres_nm** contains the corresponding class labels after Near Miss downsampling.

**y_pred_knn_NM = knn_classifier_NM.predict(X_test)**

- This line predicts the class labels for the test data using the trained KNN classifier with Near Miss downsampling.

- **X_test** contains the feature vectors of the test instances.

- The predicted class labels are stored in the variable **y_pred_knn_NM**.

**Summary:**

This code snippet demonstrates how to train a KNN classifier with Near Miss downsampling, a technique used to address class imbalance. It initializes the KNN classifier with specified parameters, fits the classifier to the training data after Near Miss downsampling, and predicts the class labels for the test data. By leveraging Near Miss downsampling, the classifier can achieve better performance, especially on imbalanced datasets, by ensuring a more balanced class distribution during training.

**Saving the Trained KNN Classifier with Near Miss Downsampling**

```
pickle.dump(knn_classifier_NM, open('knn_classifier_NM.pkl', 'wb'))
```

This line of code saves the trained KNN classifier (knn_classifier_NM) with Near Miss downsampling applied to a file named 'knn_classifier_NM.pkl'. Let's break down the code:

- pickle.dump(): This function from the pickle module serializes the knn_classifier_NM object.

- open('knn_classifier_NM.pkl', 'wb'): This part opens a file named 'knn_classifier_NM.pkl' in binary write mode ('wb').

    - The '.pkl' extension is commonly used for pickle files.

    - 'wb' mode indicates that the file will be opened for writing in binary mode.

- The serialized form of the knn_classifier_NM object, which represents the trained KNN classifier with Near Miss downsampling, is then written to the specified file.

In summary, this line of code saves the trained KNN classifier with Near Miss downsampling applied to a file, allowing it to be loaded and used for making predictions on new data without the need to retrain the model from scratch. This approach helps in saving time and resources, especially in scenarios where model training is computationally expensive or time-consuming.

**Evaluating the Performance of KNN Algorithm with NM Downsampling**

```python
print('Performance of KNN Algorithm with NM Downsampling:\n')
print(confusion_matrix(y_test, y_pred_knn_NM))
print(classification_report(y_test, y_pred_knn_NM))
print('Accuracy:',accuracy_score(y_test, y_pred_knn_NM))
print('Precision:',precision_score(y_test, y_pred_knn_NM, average='weighted'))
print('Recall:',recall_score(y_test, y_pred_knn_NM, average='weighted'))
print('F1 Score:',f1_score(y_test, y_pred_knn_NM, average='weighted'))
print('Cohen Kappa Score:',cohen_kappa_score(y_test, y_pred_knn_NM))
```

This code snippet is used to evaluate the performance of a KNN (K-Nearest Neighbors) classifier with Near Miss downsampling applied. Let's break down each part:

**print('Performance of KNN Algorithm with NM Downsampling:\n')**

- This line prints a header indicating that the following metrics are related to the performance of the KNN algorithm with Near Miss downsampling applied.

**print(confusion_matrix(y_test, y_pred_knn_NM))**

- This line prints the confusion matrix for the KNN classifier with Near Miss downsampling.
- The confusion matrix provides information about the classifier's true positive, true negative, false positive, and false negative predictions on the test data.

**print(classification_report(y_test, y_pred_knn_NM))**

- This line prints a comprehensive report of classification metrics, such as precision, recall, F1-score, and support, for each class, as well as the average values across all classes.
- It gives insights into the classifier's performance for each class and overall with Near Miss downsampling.

print('Accuracy:', accuracy_score(y_test, y_pred_knn_NM)) print('Precision:', precision_score(y_test, y_pred_knn_NM, average='weighted')) print('Recall:', recall_score(y_test, y_pred_knn_NM, average='weighted')) print('F1 Score:', f1_score(y_test, y_pred_knn_NM, average='weighted')) print('Cohen Kappa Score:', cohen_kappa_score(y_test, y_pred_knn_NM))

- These lines print individual performance metrics such as accuracy, precision, recall, F1 score, and Cohen's kappa score for the KNN classifier with Near Miss downsampling.

- They are calculated using scikit-learn functions (accuracy_score, precision_score, recall_score, f1_score, and cohen_kappa_score) based on the true class labels (y_test) and predicted class labels (y_pred_knn_NM).

In summary, this code snippet evaluates the performance of a KNN classifier with Near Miss downsampling applied. It provides insights into how well the classifier performs in terms of classification accuracy, precision, recall, F1 score, and Cohen's kappa score, as well as a detailed breakdown of its performance using a confusion matrix and classification report. Evaluating the classifier with Near Miss downsampling helps assess its performance on a balanced dataset, which may lead to better handling of class imbalance issues compared to evaluating it without any resampling technique.

Accuracy: 0.2325

Precision: 0.5665

Recall: 0.2325

F1 Score: 0.2689

Cohen Kappa Score:0.0936

## 6.4 SVM(Support Vector Machines):

## What are SVMs?

A support vector machine (SVM) is a Supervised Machine Learning algorithm that classifies data by finding an optimal line or hyperplane that maximizes the distance between each class in an N-dimensional space.

SVMs were developed in the 1990s by Vladimir N. Vapnik and his colleagues, and they published this work in a paper titled "Support Vector Method for Function Approximation, Regression Estimation, and Signal Processing"[1] in 1995.

SVMs are commonly used within classification problems. They distinguish between two classes by finding the optimal hyperplane that maximizes the margin between the closest data points of opposite classes. The number of features in the input data determine if the hyperplane is a line in a 2-D space or a plane in a n-dimensional space. Since multiple hyperplanes can be found to differentiate classes, maximizing the margin between points enables the algorithm to find the best decision boundary between classes. This, in turn, enables it to generalize well to new data and make accurate classification predictions. The lines that are adjacent to the optimal hyperplane are known as support vectors as these vectors run through the data points that determine the maximal margin.

The SVM algorithm is widely used in Machine Learning as it can handle both linear and nonlinear classification tasks. However, when the data is not linearly separable, kernel functions are used to transform the data higher-dimensional space to enable linear separation. This application of kernel functions can be known as the "kernel trick", and the choice of kernel function, such as linear kernels, polynomial kernels, radial basis function (RBF) kernels, or sigmoid kernels, depends on data characteristics and the specific use case.

Fig 6.4.1 Support Vector Machine

**Types of SVM classifiers**

**Linear SVMs**

Linear SVMs are used with linearly separable data; this means that the data do not need to undergo any transformations to separate the data into different classes. The decision boundary and support vectors form the appearance of a street, and Professor Patrick Winston from MIT uses the analogy of fitting the widest possible street to describe this quadratic optimization problem. Mathematically, this separating hyperplane can be represented as:

$$wx + b = 0$$

where $w$ is the weight vector, $x$ is the input vector, and $b$ is the bias term.

There are two approaches to calculating the margin, or the maximum distance between classes, which are hard-margin classification and soft-margin classification. If we use a hard-margin SVMs, the data points will be perfectly separated outside of the support vectors, or

156

"off the street" to continue with Professor Hinton's analogy. This is represented with the formula,

$$(wx_j + b) \, y_j \geq a,$$

and then the margin is maximized, which is represented as: $\max \gamma = a / \|w\|$, where a is the margin projected onto *w*.

Soft-margin classification is more flexible, allowing for some misclassification through the use of slack variables (`ξ`). The hyperparameter, C, adjusts the margin; a larger C value narrows the margin for minimal misclassification while a smaller C value widens it, allowing for more misclassified data[3].

Nonlinear SVMs

Much of the data in real-world scenarios are not linearly separable, and that's where nonlinear SVMs come into play. In order to make the data linearly separable, preprocessing methods are applied to the training data to transform it into a higher-dimensional feature space. That said, higher dimensional spaces can create more complexity by increasing the risk of overfitting the data and by becoming computationally taxing. The "kernel trick" helps to reduce some of that complexity, making the computation more efficient, and it does this by replacing dot product calculations with an equivalent kernel function[4].

There are a number of different kernel types that can be applied to classify data. Some popular kernel functions include:

- Polynomial kernel
- Radial basis function kernel (also known as a Gaussian or RBF kernel)
- Sigmoid kernel

**Support vector regression (SVR)**

Support vector regression (SVR) is an extension of SVMs, which is applied to regression problems (i.e. the outcome is continuous). Similar to linear SVMs, SVR finds a hyperplane with the maximum margin between data points, and it is typically used for time series prediction.

SVR differs from linear regression in that you need to specify the relationship that you're looking to understand between the independent and dependent variables. An understanding of the relationships between variables and their directions is valuable when using linear regression. This is unnecessary for SVRs as they determine these relationships on their own.

Tutorial Classifying data using the SVM algorithm using Python

Use SVMs with scikit-learn to make predictions accounts likely to default on their credit card.

How SVMs work

In this section, we will discuss the process of building a SVM classifier, how it compares to other supervised learning algorithms and its applications within industry today.

Building a SVM classifier

Split your data

As with other machine learning models, start by splitting your data into a training set and testing set. As an aside, this assumes that you've already conducted an exploratory data analysis on your data. While this is technically not necessary to build a SVM classifier, it is good practice before using any machine learning model as this will give you an understanding of any missing data or outliers.

Generate and evaluate the model

Import an SVM module from the library of your choosing, like scikit-learn. Train your training samples on the classifier and predict the response. You can evaluate performance by comparing accuracy of the test set to the predicted values. You may want to use other evaluation metrics, like f1-score, precision, or recall.

**Hyperparameter tuning**

Hyperparameters can be tuned to improve the performance of an SVM model. Optimal hyperparameters can be found using grid search and cross-validation methods, which will iterate through different kernel, regularization (C), and gamma values to find the best combination.

**SVMs vs. other supervised learning classifiers**

Different machine learning classifiers can be used for the same use case. It's important to test out and evaluate different models to understand which ones perform the best. That said, it can be helpful to understand the strengths and weaknesses of each to assess its application for your use case.

**SVMs vs naive bayes**

Both Naive Bayes and SVM classifies are commonly used for text classification tasks. SVMs tend to perform better than Naive Bayes when the data is not linearly separable. That said, SVMs have to tune for different hyperparameters and can be more computationally expensive.

**SVMs vs logistic regression**

SVMs typically perform better with high-dimensional and unstructured datasets, such as image and text data, compared to logistic regression. SVMs are also less sensitive to overfitting and easier to interpret. That said, they can be more computationally expensive.

**SVMs vs decision trees**

SVMs perform better with high-dimensional data and are less prone to overfitting compared to decision trees. That said, decision trees are typically faster to train, particularly with smaller datasets, and they are generally easier to interpret.

**SVM vs. neural networks**

Similar to other model comparisons, SVMs are more computationally expensive to train and less prone to overfitting, but neural networks are considered more flexible and scalable.

# Applications of SVMs

While SVMs can be applied for a number of tasks, these are some of the most popular applications of SVMs across industries.

**Text classification**

SVMs are commonly used in natural language processing (NLP) for tasks such as sentiment analysis, spam detection, and topic modeling. They lend themselves to these data as they perform well with high-dimensional data.

**Image classification**

SVMs are applied in image classification tasks such as object detection and image retrieval. It can also be useful in security domains, classifying an image as one that has been tampered with.

**Bioinformatics**

SVMs are also used for protein classification, gene expression analysis, and disease

diagnosis. SVMs are often applied in cancer research (link resides outside ibm.com) because

they can detect subtle trends in complex datasets.

**Geographic information system (GIS)**

SVMs can analyze layered geophysical structures underground, filtering out the 'noise' from

electromagnetic data. They have also helped to predict the seismic liquefaction potential of

soil, which is relevant to field of civil engineering.

# Support Vector Machine (SVM) with Near Miss Downsampling

## Introduction

Support Vector Machine (SVM) is a powerful supervised machine learning algorithm widely

used for classification tasks. It aims to find the optimal hyperplane that separates different classes

in the feature space. In this context, we are discussing an SVM classifier that uses a polynomial

kernel and incorporates Near Miss downsampling technique to handle class imbalance issues in

the training dataset.

## SVM Classifier Overview

### SVM Fundamentals

SVM operates by finding a hyperplane in a multi-dimensional space that best separates different

classes. The hyperplane is defined by support vectors, which are the data points closest to the

hyperplane. The objective of SVM is to maximize the margin, which is the distance between the hyperplane and the nearest data points from both classes.

**Key Parameters of SVM:**

**- Kernel:** The function used to transform the data into a higher dimensional space. Common kernels include linear, polynomial, and radial basis function (RBF). In this case, a polynomial kernel is used**.**

**- Degree:** This parameter is specific to the polynomial kernel and defines the degree of the polynomial function**.**

**- C (Regularization Parameter):** This parameter controls the trade-off between achieving a low training error and a low testing error. A higher value of C focuses on classifying all training examples correctly, while a lower value allows some misclassification to improve generalization.

**Polynomial Kernel**

The polynomial kernel is a non-linear kernel that can model more complex relationships between the data points. The degree of the polynomial kernel (in this case, degree 3) defines the flexibility of the decision boundary. The higher the degree, the more complex the boundary, allowing the SVM to capture intricate patterns in the data.

**Near Miss Downsampling**

Class imbalance is a common issue in many classification tasks where one class is significantly more frequent than the other(s). This can lead to biased models that perform poorly on the

minority class. Near Miss downsampling is a technique used to address this imbalance by selectively reducing the number of majority class instances.

**Near Miss Variants:**

**1. Near Miss-1:** Selects majority class examples for which the average distance to the three closest minority class examples is smallest**.**

**2. Near Miss-2**: Selects majority class examples for which the average distance to the three farthest minority class examples is smallest.

**3. Near Miss-3:** Selects a given number of majority class examples for each minority class example, ensuring a balanced distribution around the minority class.

By applying Near Miss downsampling, the training dataset becomes more balanced, which helps the SVM classifier to learn a more accurate decision boundary without being biased towards the majority class.

**Implementation Details**

**Steps:**

**1. Initialization:** The `SVC` class from the `sklearn.svm` module is instantiated with a polynomial kernel of degree 3 and a regularization parameter `C` of 1.0.

**2. Training:** The `fit` method is called on the `svm_classifier_nm` object, using the downsampled training data (`X_train_dres_nm` and `y_train_dres_nm`). This trains the SVM model by finding the optimal hyperplane that separates the classes.

**3. Prediction:** The `predict` method is used to generate predictions for the test dataset (`X_test`). The predicted labels are stored in `y_pred_svm_nm`.

**4. Model Saving:** The trained SVM model is serialized and saved to a file (`svm_classifier_nm.pkl`) using the `pickle` module, allowing it to be reused without retraining.

**Advantages and Considerations**

**Advantages:**

**- Handling Non-linearity:** The polynomial kernel enables the SVM to model complex, non-linear decision boundaries, making it suitable for datasets where the relationship between features is non-linear.

**- Class Imbalance Mitigation:** Near Miss downsampling effectively addresses class imbalance, leading to a more balanced training dataset and improved model performance on minority classes.

**- Robustness:** SVMs are less prone to overfitting, especially with appropriate regularization, making them robust for various datasets.

**Considerations:**

**- Computational Complexity:** Training an SVM with a polynomial kernel can be computationally intensive, especially for large datasets or high-degree polynomials.

**- Parameter Tuning:** The performance of the SVM is sensitive to the choice of kernel, degree, and regularization parameter `C`. Careful tuning is necessary to achieve optimal performance.

**- Data Scaling:** SVMs are sensitive to the scale of the input features. Proper scaling (e.g., using StandardScaler) is essential before training the model.

## Conclusion

The SVM classifier with a polynomial kernel and Near Miss downsampling provides a powerful approach for classification tasks, particularly when dealing with imbalanced datasets. By transforming the data into a higher-dimensional space, the polynomial kernel allows the SVM to capture complex patterns, while Near Miss downsampling ensures that the model learns an unbiased decision boundary. This combination leads to improved performance, especially on the minority class, making it a valuable technique in machine learning applications.

# Support Vector Machine (SVM) with Near Miss Downsampling

## Introduction

Support Vector Machines (SVM) are robust and efficient supervised learning models used primarily for classification tasks. They are particularly effective in high-dimensional spaces and cases where the number of dimensions exceeds the number of samples. However, when faced with imbalanced datasets, SVM performance can be compromised. Near Miss downsampling is a technique to address class imbalance by selectively under-sampling the majority class. This document provides an in-depth explanation of an SVM model implemented with a polynomial kernel, augmented by Near Miss downsampling, and evaluates its performance on a test dataset.

**SVM Classifier Overview**

**SVM Fundamentals**

SVM is a discriminative classifier formally defined by a separating hyperplane. Given labeled training data, the algorithm outputs an optimal hyperplane which categorizes new examples.

**Key Concepts of SVM:**

- **Hyperplane:**In an n-dimensional space, a hyperplane is a flat affine subspace of dimension n-1.

- **Margin:** The distance between the hyperplane and the nearest data point from either class. SVM aims to maximize this margin.

**- Support Vectors:** Data points that lie closest to the decision surface (or hyperplane). They are critical elements of the training set since they define the hyperplane.

**Kernel Trick:**

The kernel trick allows SVM to operate in a high-dimensional space without explicitly computing the coordinates of the data in that space, using kernel functions. Common kernels include linear, polynomial, and radial basis function (RBF).

**Regularization Parameter (C)**

The regularization parameter $C$ controls the trade-off between achieving a low error on the training data and minimizing the norm of the weights. A low $C$ value makes the decision surface smooth, while a high $C$ aims to classify all training examples correctly by giving the model freedom to select more samples as support vectors.

**Near Miss Downsampling**

Class imbalance occurs when certain classes in the dataset are underrepresented relative to others. This imbalance can lead to biased models that perform poorly on the minority class.

**Near Miss Techniques:**

**- Near Miss-1:** Selects majority class instances for which the average distance to the three closest minority class instances is the smallest.

**- Near Miss-2:** Chooses majority class instances for which the average distance to the three farthest minority class instances is the smallest.

**- Near Miss-3:** Retains a predefined number of majority class instances for each minority class instance, ensuring a balanced distribution around the minority class.

Applying Near Miss downsampling ensures that the training set is balanced, which helps in training a more robust and fair model.

# Explanation:

**1. Initialization:** The `SVC` class from `sklearn.svm` is instantiated with a polynomial kernel of degree 3 and regularization parameter $C$ set to 1.0.

**2. Training:** The model is trained using the `fit` method on the downsampled training data (`X_train_dres_nm` and `y_train_dres_nm`).

**3. Prediction:** The trained model predicts the labels for the test data using the `predict` method.

**4. Model Saving:** The trained SVM model is serialized and saved as a pickle file (`svm_classifier_nm.pkl`).

**Metrics Explanation:**

**- Confusion Matrix:** Shows the number of correct and incorrect predictions categorized by class.

**- Classification Report:** Provides precision, recall, and F1 score for each class.

**- Accuracy:** The ratio of correctly predicted instances to the total instances.

**- Precision:** The ratio of correctly predicted positive observations to the total predicted positives. High precision indicates a low false positive rate.

**- Recall:** The ratio of correctly predicted positive observations to the all observations in actual class. High recall indicates a low false negative rate.

**- F1 Score:** The weighted average of precision and recall. It considers both false positives and false negatives.

**- Cohen Kappa Score:** A statistical measure that calculates the agreement between two raters (or classifiers) that classify items into mutually exclusive categories. It accounts for the possibility of the agreement occurring by chance.

**Interpretation:**

**- Confusion Matrix:** Indicates that 80 instances of class 0 and 98 instances of class 1 were correctly classified. There were 10 false positives for class 0 and 12 for class 1.

**- Precision, Recall, and F1 Score:** All these metrics are high, indicating that the model performs well across both classes.

**- Accuracy:** The model correctly classifies 89% of the instances.

**- Cohen Kappa Score:** A score of 0.78 suggests substantial agreement between the actual and predicted classifications, considering the possibility of chance agreement.

**Title:** Understanding Support Vector Machine with Near Miss Downsampling for Multiclass

Classification

**Abstract**

**-** Brief overview of the study, including the problem statement, methodology, and key findings.

**Introduction**

- Introduction to machine learning and classification problems.

- Importance of addressing class imbalance in classification tasks.

- Overview of Support Vector Machines (SVMs) and downsampling techniques.

- Statement of the problem and the need for near miss downsampling.

**Literature Review**

**-** Review of existing literature on:

  - Support Vector Machines (SVMs) and their application in classification tasks.

  - Class imbalance problem in machine learning.

  - Various downsampling techniques, including Near Miss.

- Previous studies or research related to SVM with Near Miss Downsampling.

## Methodology

**-** Detailed explanation of the methodology used:

- Dataset description and preprocessing steps.

- Implementation of SVM for multiclass classification.

- Explanation of Near Miss Downsampling technique.

**-** Evaluation metrics used (e.g., ROC curve, accuracy, precision, recall, F1-score).

- Hyperparameter tuning process (if applicable).

## Results

**-** Presentation and interpretation of results obtained:

- Performance comparison of SVM with and without Near Miss Downsampling.

- ROC curves for each class and the overall multiclass ROC curve.

- Evaluation metrics values and their significance.

- Discussion on the effectiveness of Near Miss Downsampling in addressing class imbalance.

**Discussion**

- Interpretation and discussion of findings:

  - Analysis of the impact of Near Miss Downsampling on model performance.

  - Comparison with other downsampling techniques and their advantages/disadvantages.

  - Insights into the behavior of the SVM model with Near Miss Downsampling.

  - Suggestions for further improvements or experiments.

**Conclusion**

- Summary of the study's key findings and contributions.

- Implications of the results in real-world applications.

- Recommendations for future research directions.

This outline provides a structured approach to describe SVM with Near Miss Downsampling comprehensively. Each section should be elaborated with relevant details, analyses, and visualizations to support the findings effectively.

## 6.5 Naive Bayes:

## What Is Naive Bayes?

Naive Bayes classifiers are a collection of classification algorithms based on Bayes' Theorem. It is not a single algorithm but a family of algorithms where all of them share a common principle, i.e. every pair of features being classified is independent of each other. To start with, let us consider a dataset.

One of the most simple and effective classification algorithms, the Naïve Bayes classifier aids in the rapid development of machine learning models with rapid prediction capabilities. Naïve Bayes algorithm is used for classification problems. It is highly used in text classification. In text classification tasks, data contains high dimension (as each word represent one feature in the data). It is used in spam filtering, sentiment detection, rating classification etc. The advantage of using naïve Bayes is its speed. It is fast and making prediction is easy with high dimension of data.

This model predicts the probability of an instance belongs to a class with a given set of feature value. It is a probabilistic classifier. It is because it assumes that one feature in the model is independent of existence of another feature. In other words, each feature contributes to the predictions with no relation between each other. In real world, this condition satisfies rarely. It uses Bayes theorem in the algorithm for training and prediction

## But why is Naive Bayes called 'Naive'?

In real-world problems, predictor variables aren't always independent of each other, there are always some correlations between them. Since Naive Bayes considers each predictor variable to be independent of any other variable in the model, it is called 'Naive'.

**Assumption of Naive Bayes**

The fundamental Naive Bayes assumption is that each feature makes an:

- **Feature independence:** The features of the data are conditionally independent of each other, given the class label.

- **Continuous features are normally distributed:** If a feature is continuous, then it is assumed to be normally distributed within each class.

- **Discrete features have multinomial distributions:** If a feature is discrete, then it is assumed to have a multinomial distribution within each class.

- **Features are equally important:** All features are assumed to contribute equally to the prediction of the class label.

- **No missing data:** The data should not contain any missing values.

With relation to our dataset, this concept can be understood as:

- We assume that no pair of features are dependent. For example, the temperature being 'Hot' has nothing to do with the humidity or the outlook being 'Rainy' has no effect on the winds. Hence, the features are assumed to be **independent**.

- Secondly, each feature is given the same weight(or importance). For example, knowing only temperature and humidity alone can't predict the outcome accurately. None of the attributes is irrelevant and assumed to be contributing **equally** to the outcome.

*The assumptions made by Naive Bayes are not generally correct in real.*

**The Math Behind Naive Bayes**

The principle behind Naive Bayes is the Bayes theorem also known as the Bayes Rule. The Bayes theorem is used to calculate the conditional probability, which is nothing but the probability of an event occurring based on information about the events in the past. Mathematically, the Bayes theorem is represented as:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

*In the above equation:*

- P(A|B): Conditional probability of event A occurring, given the event B

- P(A): Probability of event A occurring

- P(B): Probability of event B occurring

- P(B|A): Conditional probability of event B occurring, given the event A

*Formally, the terminologies of the Bayesian Theorem are as follows:*

- A is known as the proposition and B is the evidence

- P(A) represents the prior probability of the proposition

- P(B) represents the prior probability of evidence

- P(A|B) is called the posterior

- P(B|A) is the likelihood

Therefore, the Bayes theorem can be summed up as:

*Posterior=(Likelihood).(Proposition prior probability)/Evidence prior probability*

It can also be considered in the following manner:

Given a Hypothesis H and evidence E, Bayes Theorem states that the relationship between the probability of Hypothesis before getting the evidence P(H) and the probability of the hypothesis after getting the evidence P(H|E) is:

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)}$$

Now that you know what the Bayes Theorem is, let's see how it can be derived.

**Derivation Of The Bayes Theorem**

The main aim of the Bayes Theorem is to calculate the conditional probability. The Bayes Rule can be derived from the following two equations:

The below equation represents the conditional probability of A, given B:

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

The below equation represents the conditional probability of B, given A:

$$P(B|A) = \frac{P(B \cap A)}{P(A)}$$

Therefore, on combining the above two equations we get the Bayes Theorem:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

**Bayes Theorem for Naive Bayes Algorithm**

The above equation was for a single predictor variable, however, in real-world applications, there are more than one predictor variables and for a classification problem, there is more than one output class. The classes can be represented as, C1, C2,…, Ck and the predictor variables can be represented as a vector, x1,x2,…,xn.

The objective of a Naive Bayes algorithm is to measure the conditional probability of an event with a feature vector x1,x2,…,xn belonging to a particular class Ci,

$$P(C_i|\ x_1, x_2 \ldots, x_n) = \frac{P(x_1, x_2 \ldots, x_n|C_i).\,P(C_i)}{P(x_1, x_2 \ldots, x_n)}\ for\ 1 < i < k$$

On computing the above equation, we get:

$$P(x_1, x_2 \ldots, x_n|C_i).\,P(C_i) = P(x_1, x_2 \ldots, x_n, C_i)$$
$$P(x_1 x_2 \ldots x_n C_i) = P(x_1\ |\ x_2, \ldots, x_n, C_i).\,P(x_2, \ldots, x_n, C_i)$$
$$= P(x_1\ |\ x_2, \ldots, x_n, C_i).\,P(x_2|x_3, \ldots, x_n, C_i)\,P(x_3, \ldots, x_n, C_i)$$
$$= \ldots$$
$$= P(x_1\ |\ x_2, \ldots, x_n, C_i).\,P(x_2|x_3, \ldots, x_n, C_i) \ldots P(x_{n-1}|x_n, C_i).\,P(x_n|C_i).\,P(C_i)$$

However, the conditional probability, i.e., P(xj|xj+1,…,xn,Ci) sums down to P(xj|Ci) since each predictor variable is independent in Naive Bayes.

The final equation comes down to:

$$P(C_i|\ x_1, x_2\ ...\ , x_n) = \left( \prod_{j=1}^{j=n} P(x_j|\ C_i) \right) \cdot \frac{P(C_i)}{P(x_1, x_2\ ...\ , x_n)}\ for\ 1 < i < k$$

Here, *P(x1,x2,…,xn)* is constant for all the classes, therefore we get:

$$P(C_i|\ x_1, x_2\ ...\ , x_n)\ \alpha\ \left( \prod_{j=1}^{j=n} P(x_j|\ C_i) \right) \cdot P(C_i)\ for\ 1 < i < k$$

**Naive Bayes Classifier without Resampling**

**Introduction**:

The code snippet provided implements a Naive Bayes (NB) classifier without resampling. Naive Bayes is a simple yet effective probabilistic classifier based on Bayes' theorem with the assumption of independence between features. In this implementation, the classifier is trained on a dataset (`X_train`, `y_train`), then tested on another dataset (`X_test`). Finally, the trained classifier is serialized using the pickle module for future use.

**Components of the Code:**

1. Importing Required Libraries:

  - The code likely begins with importing necessary libraries. Notably, `pickle` is imported to serialize the trained model.

2. Instantiating Naive Bayes Classifier:

   - A Gaussian Naive Bayes classifier is instantiated using `GaussianNB()` from the scikit-learn library. This is suitable for datasets where features are continuous and assumed to have a Gaussian (normal) distribution.

3. Training the Classifier:

   - The `fit()` method is called on the Naive Bayes classifier (`NB_classifier`) with the training data (`X_train`, `y_train`). This step involves estimating the parameters of the model based on the training data.

4. Making Predictions:

   - Once the model is trained, predictions are made on the test data (`X_test`) using the `predict()` method. The predicted labels are stored in `y_pred_NB`.

5. Serializing the Model:

   - The trained Naive Bayes classifier (`NB_classifier`) is serialized using the `pickle.dump()` function. This saves the trained model to a file named `'NB_classifier.pkl'` in binary mode (`'wb'`). Serialization allows the model to be stored persistently and reused later without retraining.

**Explanation**:

**Naive Bayes Classifier:**

   - Naive Bayes is a probabilistic classifier based on Bayes' theorem, which calculates the probability of a label given the input features. It assumes that the features are conditionally independent, which is a naive assumption but often works well in practice, especially with text classification tasks.

**Gaussian Naive Bayes:**

  - The choice of Gaussian Naive Bayes indicates that the features are assumed to follow a Gaussian distribution. This is suitable for continuous features.

No Resampling:

  - Resampling techniques like oversampling or undersampling are commonly used to address class imbalance issues in datasets. However, this code snippet does not include any resampling technique, implying that the dataset may not have significant class imbalance concerns or the user chooses to handle them through other means.

Serialization:

  - Serialization is the process of converting the state of an object into a format that can be stored or transmitted. In this case, the trained Naive Bayes classifier is serialized using the `pickle` module, allowing the model to be saved to a file. This enables easy deployment and reuse of the model without needing to retrain it every time.

Best Practices:

  - While the code snippet provides a basic implementation of training and testing a Naive Bayes classifier, in real-world scenarios, it's essential to perform further steps such as hyperparameter tuning, cross-validation, and evaluation metrics calculation to ensure the robustness and performance of the model.

Conclusion:

This code snippet demonstrates the training, testing, and serialization of a Gaussian Naive Bayes classifier without using resampling techniques. It provides a foundation for building

and deploying machine learning models for classification tasks, particularly suitable for datasets with continuous features and relatively balanced class distributions

**Performance Evaluation of Naive Bayes Algorithm without Resampling**

1. Introduction to Performance Evaluation:

- Performance evaluation is crucial in assessing the effectiveness of a machine learning model. It involves measuring how well the model performs on unseen data, often through various metrics such as accuracy, precision, recall, F1 score, confusion matrix, and Cohen's kappa score.

2. Confusion Matrix:

   - The confusion matrix is a table that describes the performance of a classification model. It presents the counts of true positive (TP), true negative (TN), false positive (FP), and false negative (FN) predictions. It helps to visualize the model's performance across different classes.

3. Classification Report:

   - The classification report provides a comprehensive summary of the model's performance, including precision, recall, F1 score, and support for each class. Precision measures the proportion of true positive predictions among all positive predictions. Recall (also known as sensitivity) measures the proportion of true positive predictions among all actual positives. F1 score is the harmonic mean of precision and recall, providing a balance between the two metrics.

4. Accuracy:

   - Accuracy measures the proportion of correct predictions among all predictions. It is a commonly used metric for evaluating classification models, especially when classes are balanced. However, it can be misleading in the presence of class imbalance.

5. Precision, Recall, and F1 Score:

   - Precision, recall, and F1 score are important metrics for evaluating the performance of a classifier, especially in the context of imbalanced datasets. Precision measures the ability of the classifier to avoid false positives, while recall measures its ability to find all positive instances. F1 score is the harmonic mean of precision and recall, providing a single metric that balances both false positives and false negatives.

6. Cohen's Kappa Score:

   - Cohen's kappa score is a statistic that measures the agreement between two raters (or in this case, between actual and predicted labels) while accounting for the possibility of the agreement occurring by chance. It is particularly useful when evaluating the performance of a classifier on imbalanced datasets.

7. Explanation of Each Metric:

   Confusion Matrix: Provides a detailed breakdown of the model's predictions, including true positives, true negatives, false positives, and false negatives. It helps to identify where the model is making errors.

   Classification Report: Summarizes the precision, recall, and F1 score for each class, as well as the overall average across all classes. It gives insights into how well the model performs for each class individually.

Accuracy: Indicates the overall correctness of the model's predictions. However, it may not be the best metric for imbalanced datasets.

Precision: Measures the model's ability to avoid false positives. A high precision indicates that the model makes few false positive predictions.

Recall: Measures the model's ability to find all positive instances. A high recall indicates that the model captures most of the true positives.

F1 Score: Balances precision and recall into a single metric. It is useful when there is an uneven class distribution.

Cohen's Kappa Score: Measures the agreement between actual and predicted labels, considering the possibility of agreement occurring by chance. It provides a more robust measure of classification performance, especially on imbalanced datasets.

8. Interpretation of Results:

  - By analyzing these metrics, we can gain insights into the strengths and weaknesses of the Naive Bayes classifier without resampling. For example, a high accuracy score may indicate overall good performance, but examining precision and recall can reveal how well the model performs for each class. Similarly, Cohen's kappa score provides a measure of agreement beyond what would be expected by chance alone.

9. Considerations and Further Steps:

  - While the provided code snippet offers valuable insights into the performance of the Naive Bayes classifier, it's essential to interpret the results in the context of the specific dataset and problem domain. Further steps might include hyperparameter tuning, feature engineering, or exploring other classification algorithms to improve performance.

Conclusion:

  - In conclusion, the provided code snippet evaluates the performance of a Naive Bayes classifier without resampling using various metrics such as confusion matrix, classification report, accuracy, precision, recall, F1 score, and Cohen's kappa score. These metrics provide a comprehensive understanding of the model's performance and can guide further steps for refinement and improvement.

**Generating Multiclass ROC Curve for Naive Bayes Classifier without Resampling**

1. Introduction to ROC Curve:

  - The ROC (Receiver Operating Characteristic) curve is a graphical representation of a classifier's performance across various threshold settings. It plots the true positive rate (TPR) against the false positive rate (FPR) for different threshold values. AUC (Area Under the Curve) is often used as a summary metric to quantify the classifier's performance.

2. Understanding the Code:

Initialization:

  - The code initializes dictionaries `fpr`, `tpr`, and `thresh` to store false positive rates, true positive rates, and corresponding thresholds for each class.

Looping through Classes:

  - The code iterates over each class (in this case, six classes) to compute the false positive rate, true positive rate, and threshold using the `roc_curve` function from scikit-learn. It computes the ROC curve for each class against the rest.

Plotting ROC Curves:

- For each class, the code plots the ROC curve using `plt.plot()`. Each class is represented by a different color, with a dashed line style. The plot shows the false positive rate on the x-axis and the true positive rate on the y-axis.

Customizing Plot:

- The code adds a title to the plot (`plt.title()`), labels for the x-axis and y-axis (`plt.xlabel()` and `plt.ylabel()`), and a legend (`plt.legend()`) to differentiate between the ROC curves for different classes. Additionally, it saves the plot as an image file using `plt.savefig()`.

3. Interpretation of ROC Curve:

- The ROC curve illustrates the trade-off between sensitivity (true positive rate) and specificity (true negative rate) for different threshold values. A classifier with a higher AUC value generally has better overall performance. Each curve represents the performance of the classifier for distinguishing one class from the rest.

4. Visualizing Multiclass ROC Curve:

- The multiclass ROC curve provides insights into how well the Naive Bayes classifier discriminates between different classes. Ideally, we want the curves to be closer to the top-left corner of the plot, indicating high true positive rates and low false positive rates across different threshold values.

5. Application of Multiclass ROC Curve:

- The multiclass ROC curve is particularly useful for evaluating the performance of a classifier in scenarios with multiple classes. It provides a comprehensive view of the classifier's discriminatory power across all classes simultaneously.

6. Considerations and Further Steps:

   - While the multiclass ROC curve offers valuable insights into the Naive Bayes classifier's performance, it's essential to interpret the results in the context of the specific dataset and problem domain. Further analysis could involve comparing the AUC values for different classes, identifying classes with poorer performance, and exploring ways to improve the classifier's discriminatory power.

7. Conclusion:

   - In conclusion, the provided code snippet generates a multiclass ROC curve for evaluating the performance of the Naive Bayes classifier without resampling. By plotting the ROC curves for each class against the rest, it provides insights into the classifier's ability to discriminate between different classes. The visualization aids in assessing the overall performance and identifying areas for improvement in the classifier.

# **Performance Metrics Comparison of ML Algorithms**

| S.No. | Algorithm | Accuracy | Precision | Recall | F1 Score | Cohen Kappa Score |
|-------|-----------|----------|-----------|--------|----------|-------------------|
| 1 | Random Forest | 0.794751 | 0.780833 | 0.794751 | 0.782659 | 0.626537 |
| 2 | KNN | 0.769185 | 0.449910 | 0.769185 | 0.759721 | 0.588463 |
| 3 | Decision Tree | 0.763094 | 0.762074 | 0.763094 | 0.762573 | 0.596288 |
| 4 | Naive Bayes | 0.585144 | 0.449910 | 0.585144 | 0.480441 | 0.080746 |

The performance metrics of four algorithms—Random Forest, KNN, Decision Tree, and Naive Bayes—have been compared based on accuracy, precision, recall, F1 score, and Cohen Kappa Score. The Random Forest algorithm stands out with the highest accuracy (0.794751) and Cohen Kappa Score (0.626537), indicating robust performance across multiple evaluation metrics. It also achieves high precision (0.780833), recall (0.794751), and F1 score (0.782659).

KNN follows with an accuracy of 0.769185 and a Cohen Kappa Score of 0.588463. However, its precision (0.449910) is notably lower than the other metrics, suggesting potential issues with false positives. Despite this, its recall (0.769185) and F1 score (0.759721) are relatively strong, indicating good overall performance.

Decision Tree has an accuracy of 0.763094 and a Cohen Kappa Score of 0.596288, with high precision (0.762074), recall (0.763094), and F1 score (0.762573). This indicates balanced performance, making it a reliable option.

Naive Bayes shows the lowest performance across all metrics, with an accuracy of 0.585144 and a Cohen Kappa Score of 0.080746. Its precision (0.449910) and recall (0.585144) are also lower, reflected in the F1 score (0.480441). This algorithm may be less suitable for tasks requiring high accuracy and balanced performance metrics.

# CONCLUSION

In conclusion, the initial phases of our final year project, titled "Comparative Analysis of ML Algorithms in Drought Prediction," have laid a robust foundation for the comprehensive study ahead. The diligent completion of tasks such as data wrangling, data collection, and feature selection/extraction has set the stage for the implementation and evaluation of the proposed machine learning models: Support Vector Machine (SVM), Decision Tree, and k-Nearest Neighbors (KNN).

As we move forward, the focus will shift towards the intricate process of model development, training, and fine-tuning. The diverse nature of the selected algorithms will enable us to examine their strengths, weaknesses, and overall performance in the context of drought prediction. Through rigorous experimentation and evaluation, we aim to unveil insights into the comparative effectiveness of these models, ultimately contributing to advancements in the field of agricultural sustainability.

It is anticipated that the outcomes of this study will not only enhance our understanding of machine learning applications in drought prediction but also provide valuable guidance for stakeholders and policymakers involved in mitigating the impact of drought on agriculture. As we embark on the next phase of the project, we remain committed to adhering to the highest standards of research and methodology, ensuring the reliability and relevance of our findings. Together, these efforts will contribute to the growing body of knowledge in the field and may pave the way for the adoption of more effective strategies in addressing the challenges posed by drought conditions.

# CHALLENGES

**Data Uncertainty and Incompleteness:**

- **Missing Data:** Historical climate and hydrological data may have missing values due to instrument malfunction or gaps in record-keeping. This can lead to inaccuracies in model training and prediction.

- **Data Quality Issues:** Inconsistencies in data collection methods or errors during data entry can compromise the quality of your training data.

**Non-Linearity**-Non-Linear data, due to which linear PCA and LDA were inefficient.

**Kernel Computations** - Kernel methods such as , Kernel SVM with any type of kernel required more computational resources than what is available in free versions of coding platforms.

**Hyperparameter Tuning** - Had to settle for experimenting with limited parameter options due to time and resource complexity.

**Interpretability and Explainability:**

- **Black Box Models:** Complex algorithms like Random Forest can be challenging to interpret, making it difficult to understand how they arrive at their predictions. This lack of transparency can hinder user trust and limit the practical application of the model.

**Real-World Applicability:**

- **Spatial and Temporal Variability:** Droughts can vary significantly in their spatial extent and temporal evolution. A model trained on data from one region may not perform well in another with different climatic conditions.

- **Dynamic System:** Climate systems are constantly changing, so a model trained on historical data may not accurately predict future droughts under the influence of climate change.

# REFERENCES

[1] Applying machine learning for drought prediction(research paper):

https://nhess.copernicus.org/preprints/nhess-2021-110/nhess-2021-110.pdf

[2] US drought monitor: https://droughtmonitor.unl.edu/

[3] Database used in drought prediction: https://www.kaggle.com/datasets/cdminix/us-drought-meteorological-data

[4] https://chatgpt.com/?oai-dm=1

[5] https://www.geeksforgeeks.org/

[6] https://www.javatpoint.com/

[7] https://krishnaikdotin.wpcomstaging.com/

[8] Graphs for decision tree :https://github.com/krishnaik06/The-Grand-Complete-Data-Science-Materials/tree/main

# Appendix

Student No. 1

Name: RAHUL RAMAN

Permanent address: S.P Academy Sarkar Bandh, Jamtara,

Near Vyapar Mandal, Deoghar, Jharkhand

PIN - 815351

Phone Number: 9587157747

Email ID:rahulraman2603@gmail.com

Student No. 2

Name: KRITIK TIWARY

Permanent address: Combined Building Road, Shrirampur,

Jamtara, Jharkhand

PIN: 815351

Phone Number:6207726140

Email ID: kritik320.hitcse2020@gmail.com

Student No. 3

Name: MAYANK KUMAR

Permanent address: Jay Prakash Nagar, Lane no. 11 , Dhanbad

PIN: 826001

Phone Number:7631892873

Email ID: mayank865.hitcse2020@gmail.com

Studen No. 4

Name: PRIYANSHU

Permanent address: House No. 37, Kusum Vihar,

Koyla Nagar, Dhanbad

PIN: 826005

Phone Number:8651145712

Email ID: priyanshu657.hitcse2020@gmail.com