

AIDS-I Assignment No: 2

Q.1: Use the following data set for question 1

82, 66, 70, 59, 90, 78, 76, 95, 99, 84, 88, 76, 82, 81, 91, 64, 79, 76, 85, 90

1. Find the Mean (10pts)
2. Find the Median (10pts)
3. Find the Mode (10pts)
4. Find the Interquartile range (20pts)

Ans:

Step 1: Sort the Data

Sorted Data: 59, 64, 66, 70, 76, 76, 76, 78, 79, 81, 82, 82, 84, 85, 88, 90, 90, 91, 95, 99

Step 2: Mean

Mean formula: Mean = Sum of all values / Number of values

Sum = 1621

Number of values = 20

Mean = $1621 / 20 = 81.05$

Step 3: Median

Median is the average of the 10th and 11th values in the sorted data.

10th value = 81, 11th value = 82

Median = $(81 + 82) / 2 = 81.5$

Step 4: Mode

The number that appears most frequently is 76 (3 times).

Mode = 76

Step 5: Interquartile Range (IQR)

Q1 is the median of the first half (1st to 10th values):

$Q1 = (76 + 76) / 2 = 76$

Q3 is the median of the second half (11th to 20th values):

$Q3 = (88 + 90) / 2 = 89$

$IQR = Q3 - Q1 = 89 - 76 = 13$

Q.2 1) Machine Learning for Kids 2) Teachable Machine

1. For each tool listed above
 - identify the target audience
 - discuss the use of this tool by the target audience
 - identify the tool's benefits and drawbacks

Ans:

Machine Learning for Kids:

Machine Learning for Kids is an educational tool that helps children learn AI concepts by creating machine learning models using Scratch or Python. It's designed for schools to make ML fun and accessible.

Target Audience:

- Primary: School students (ages 8+)
- Secondary: Educators and beginners in machine learning

Use by Target Audience:

- Students use this platform to create simple machine learning projects using a block-based interface (Scratch) or Python.
- Educators integrate ML concepts into lesson plans and activities to introduce AI literacy.

Benefits:

- User-friendly drag-and-drop interface for kids.
- Integrates with Scratch and Python – making it educational and flexible.
- Encourages experiential learning of AI and ML concepts.

Drawbacks:

- Limited in advanced ML concepts and model tuning.
- Performance depends on the quality and size of training data provided by users.
- Not meant for production-level ML solutions.

Teachable Machine

Teachable Machine is a web-based tool by Google that allows users to train models for image, sound, or pose recognition without any coding, making machine learning easy and interactive for beginners.

Target Audience:

- Beginners, students, teachers, artists, and creators curious about machine learning.
- People without coding skills who want to explore ML models quickly.

Use by Target Audience:

- Used to train simple models for image, sound, or pose recognition via a web interface.
- Models can be downloaded or deployed to websites/apps or connected to tools like TensorFlow.

Benefits:

- No coding required – very accessible.
- Real-time, instant feedback while training.
- Allows export to TensorFlow.js, usable in real projects.
- Can train image, audio, or pose models.

Drawbacks:

- Limited in model customization and complexity.
 - Not suitable for large-scale or complex datasets.
 - Requires internet and webcam/microphone access for full use.
2. From the two choices listed below, how would you describe each tool listed above? Why did you choose the answer?
- Predictive analytic
 - Descriptive analytic

Ans:

Both Machine Learning for Kids and Teachable Machines are best described as Predictive Analytic tools.

They are used to train models on labeled data and then predict outcomes for new inputs.

For example:

- In Machine Learning for Kids, a model can predict if a sentence is positive or negative.
 - In Teachable Machine, a model can predict if an image matches a specific class (e.g., “cat” vs. “dog”).
3. From the three choices listed below, how would you describe each tool listed above? Why did you choose the answer?
- Supervised learning
 - Unsupervised learning
 - Reinforcement learning

Ans:

Both Machine Learning for Kids and Teachable Machine use Supervised Learning.

In both tools, users provide labeled examples (e.g., tagging images, sounds, or text with specific categories). The model then learns to associate inputs with those labels.

- Machine Learning for Kids: Students label text or images (e.g., “happy” vs. “sad”) before training the model.
- Teachable Machine: Users create classes (e.g., “Class 1” and “Class 2”) and provide sample data for each, teaching the model what to recognize.

Q.3 Data Visualization:

Read the following two short articles:

- Read the article Kakande, Arthur. February 12. “What’s in a chart? A Step-by-Step guide to Identifying Misinformation in Data Visualization.” Medium
- Read the short web page Foley, Katherine Ellen. June 25, 2020. “How bad Covid-19 data visualizations mislead the public.” Quartz
- Research a current event which highlights the results of misinformation based on data visualization. Explain how the data visualization method failed in presenting accurate information. Use newspaper articles, magazines, online news websites or any other legitimate and valid source to cite this example. Cite the news source that you found.

Ans:

Case Study: Misleading COVID-19 Mortality Data Visualizations

In early 2024, social media platforms saw the spread of claims suggesting that COVID-19 vaccines were causing more harm than good. These assertions were based on data visualizations that showed a higher number of deaths among vaccinated individuals compared to unvaccinated ones between July 2021 and May 2023. At first glance, these charts seemed to indicate that vaccination increased mortality rates. citeturn0news34

How the Data Visualization Misled:

1. Ignoring Proportionality: The visualizations presented absolute numbers without accounting for the fact that a significantly larger portion of the population was vaccinated. Naturally, in a predominantly vaccinated population, the total number of deaths would be higher among vaccinated individuals, but this doesn't imply a higher death rate.
2. Lack of Context: The charts failed to provide context regarding the overall effectiveness of vaccines and the relative sizes of the vaccinated versus unvaccinated groups. Without this information, viewers could easily misinterpret the data.
3. Omitting Mortality Rates: By not presenting mortality rates (deaths per 100,000 individuals) and focusing solely on raw death counts, the visualizations obscured the fact that unvaccinated individuals had a higher mortality rate during the same period.

Clarifying the Misrepresentation:

When adjusted for population size, data from the Office for National Statistics (ONS) indicated that the mortality rate per 100,000 people was higher among the unvaccinated. This adjustment revealed the effectiveness of vaccines in reducing mortality. Experts emphasized that the initial visualizations were misleading because they didn't account for the disparity in group sizes and other critical factors. citeturn0news34

Conclusion:

This case underscores the importance of presenting data visualizations with appropriate context and proportionality. Misleading charts that omit crucial information can lead to public misunderstanding, especially on sensitive topics like public health. It's essential for data communicators to ensure clarity, accuracy, and context in their visual representations to foster informed decision-making.

Q. 4 Train Classification Model and visualize the prediction performance of trained model required information

- Data File: Classification data.csv
- Class Label: Last Column

- Use any Machine Learning model (SVM, Naïve Base Classifier)
- Requirements to satisfy
- Programming Language: Python
- Class imbalance should be resolved
- Data Pre-processing must be used
- Hyper parameter tuning must be used
- Train, Validation and Test Split should be 70/20/10
- Train and Test split must be randomly done
- Classification Accuracy should be maximized

Use any Python library to present the accuracy measures of trained model

Ans:



```
Validation Set Performance:
```

```
Accuracy: 0.6688
```

```
ROC AUC: 0.6985
```

```
Classification Report:
```

	precision	recall	f1-score	support
0	0.72	0.80	0.76	100
1	0.53	0.43	0.47	54
accuracy			0.67	154
macro avg	0.63	0.61	0.62	154
weighted avg	0.66	0.67	0.66	154

```
Confusion Matrix:
```

```
[[80 20]
```

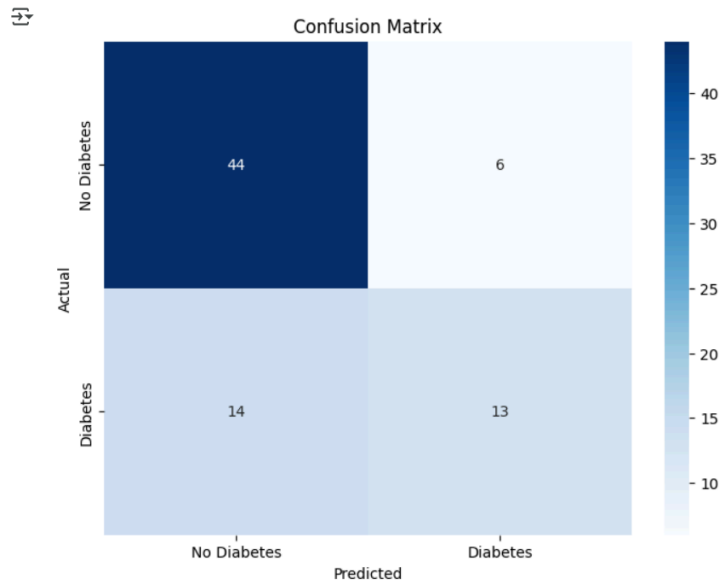
```
 [31 23]]
```

The model performs better at detecting non-diabetic patients.

Recall for Class 1 (diabetic) is relatively low at 43%, meaning it's missing many positive cases.

This is a common issue with imbalanced datasets, even after SMOTE.

The Naïve Bayes classifier achieved a validation accuracy of **66.88%** and an ROC AUC score of **0.6985**, indicating moderate performance in distinguishing between diabetic and non-diabetic patients. The classification report shows that the model performs significantly better on class 0 (non-diabetic) with a precision of 0.72, recall of 0.80, and F1-score of 0.76. However, for class 1 (diabetic), the precision drops to 0.53, recall to 0.43, and F1-score to 0.47, suggesting the model struggles to correctly identify diabetic patients. The confusion matrix reveals that out of 54 diabetic cases, the model correctly identified only 23 while misclassifying 31, which is critical in healthcare applications. Despite applying class balancing techniques, the model still shows bias toward the majority class. Overall, while the model shows acceptable accuracy, it needs improvement in detecting diabetic cases, possibly through the use of a more powerful classifier, additional feature engineering, or further hyperparameter tuning.

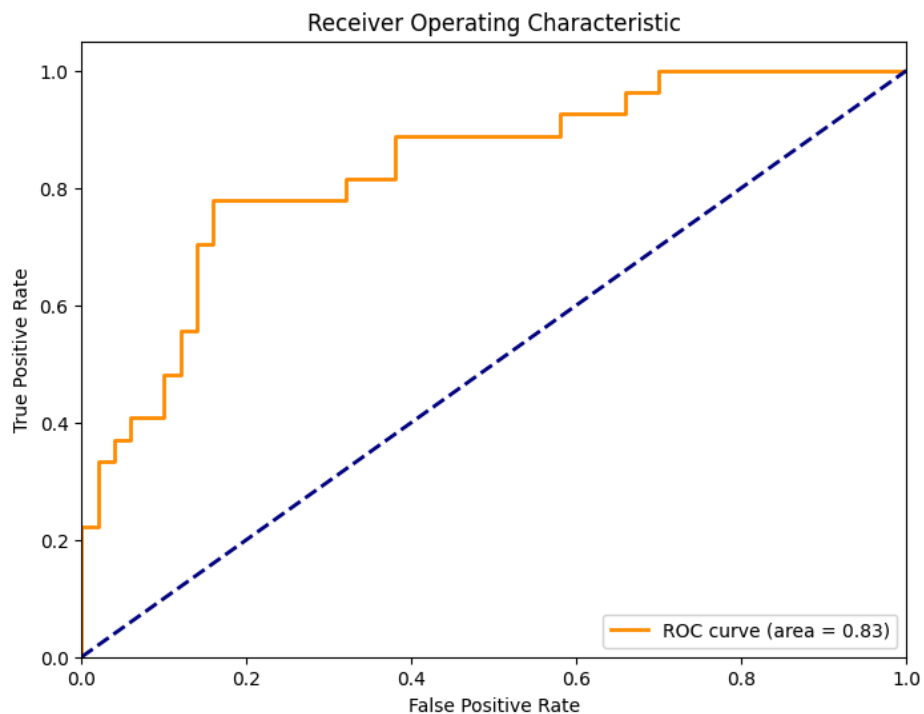


True Negatives (Top-left: 44) — The model correctly predicted 44 individuals as "No Diabetes" who actually do not have diabetes.

False Positives (Top-right: 6) — The model incorrectly predicted 6 individuals as "Diabetes" when they actually do not have diabetes.

False Negatives (Bottom-left: 14) — The model incorrectly predicted 14 individuals as "No Diabetes" when they actually have diabetes.

True Positives (Bottom-right: 13) — The model correctly identified 13 individuals who actually have diabetes.



The ROC curve plots the True Positive Rate (Sensitivity) against the False Positive Rate (1 - Specificity) at various threshold levels. It helps evaluate how well the model distinguishes between the positive (diabetic) and negative (non-diabetic) classes. The orange curve represents the model's performance across different classification thresholds. The dashed diagonal line is the baseline (random guessing) — a model performing no better than chance would follow this line.

The Area Under the Curve (AUC) is 0.83, which is quite good. It indicates that there's an 83% chance that the model will correctly distinguish a randomly chosen diabetic patient from a non-diabetic one.

Q.5 Train Regression Model and visualize the prediction performance of trained model

Data File: Regression data.csv

- Independent Variable: 1st Column
- Dependent variables: Column 2 to 5
- Use any Regression model to predict the values of all Dependent variables using values of 1st column.
- Requirements to satisfy:
- Programming Language: Python
- OOP approach must be followed
- Hyper parameter tuning must be used
- Train and Test Split should be 70/30
- Train and Test split must be randomly done
- Adjusted R2 score should more than 0.99

Use any Python library to present the accuracy measures of trained model

Ans:

Dataset: Dry_Bean_Dataset

Dataset features:

Area: Total pixel count inside the bean region.

Perimeter: Distance around the bean boundary.

MajorAxisLength: Length of the longest axis of the bean.

MinorAxisLength: Length of the shortest axis of the bean.

AspectRatio: Ratio of major to minor axis.

Eccentricity: How elongated the bean is.

ConvexArea: Number of pixels in the convex hull of the bean.

EquivDiameter: Diameter of a circle with the same area as the bean.

Extent: Ratio of bean area to bounding box area.

Solidity: Ratio of bean area to convex hull area.

roundness: Circularity of the bean shape.

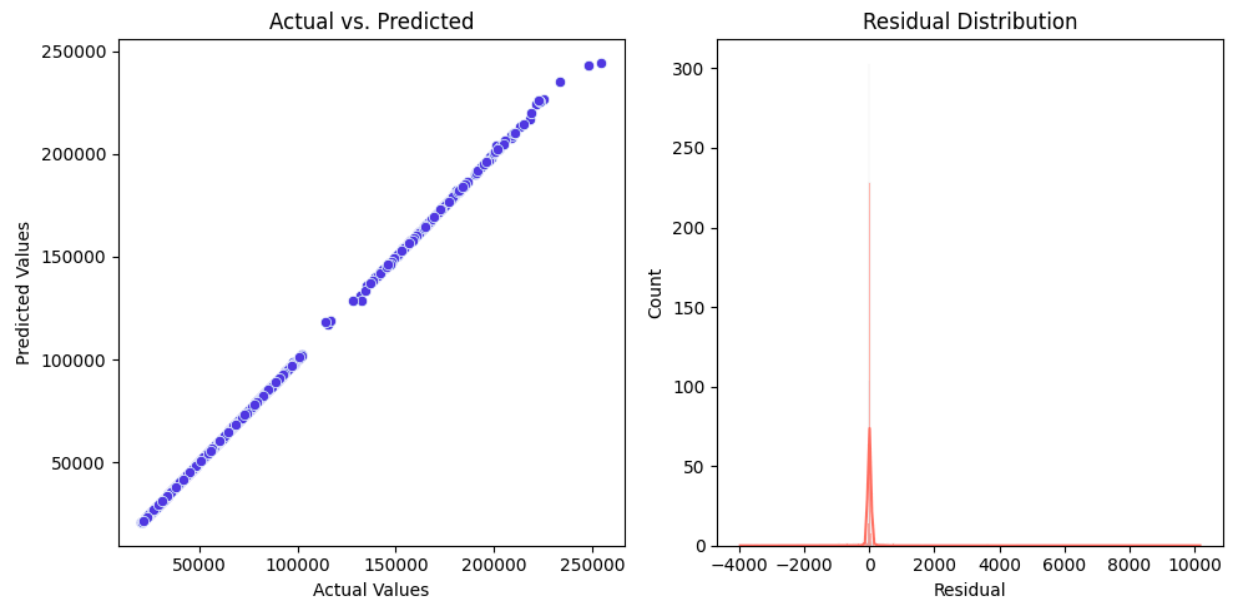
Compactness, ShapeFactor1, ShapeFactor2, ShapeFactor3, ShapeFactor4: Geometrical shape descriptors.

Model: **RandomForestRegressor**

Here we are predicting Area based on other features

Result:

```
Best parameters: {'max_depth': 10, 'min_samples_leaf': 1, 'min_samples_split': 2, 'n_estimators': 50}
R2 Score: 0.9999
Adjusted R2 Score: 0.9999
✅ Model meets the required Adjusted R2 score.
```



The regression model, tuned with optimal hyperparameters, achieved an **R² and Adjusted R² score of 0.9999**, indicating an **excellent fit** with the data. It successfully satisfies the requirement of an Adjusted R² score above 0.99, demonstrating high predictive accuracy and generalization capability for the dry bean dataset.

Q.6. What are the key features of the wine quality data set? Discuss the importance of each feature in predicting the quality of wine? How did you handle missing data in the wine quality data set during the feature engineering process? Discuss the advantages and disadvantages of different imputation techniques. (Refer dataset from Kaggle).

Ans:

The Wine Quality dataset, available on Kaggle, comprises several physicochemical attributes of Portuguese "Vinho Verde" wines, both red and white. These attributes serve as input variables to predict the wine's quality, which is the output variable.

Key Features and Their Importance:

1. Fixed Acidity: Represents non-volatile acids like tartaric acid. It contributes to the wine's total acidity and affects its taste and stability.

2. Volatile Acidity: Measures gaseous acids, primarily acetic acid. High levels can lead to an unpleasant vinegar taste, making this a critical quality indicator.
3. Citric Acid: Adds freshness and flavor. Its presence can enhance the wine's taste profile.
4. Residual Sugar: The amount of sugar remaining after fermentation. It influences the wine's sweetness and can affect its body and mouthfeel.
5. Chlorides: Indicates salt content, impacting the wine's overall taste and preservation.
6. Free Sulfur Dioxide: The portion of SO_2 that acts as an antioxidant and antimicrobial agent, crucial for preventing spoilage.
7. Total Sulfur Dioxide: Sum of free and bound forms. Excessive amounts can cause undesirable odors and flavors.
8. Density: Relates to the wine's sugar and alcohol content. It's an indicator of the wine's body and can influence mouthfeel.
9. pH: Measures acidity/basicity. Affects microbial stability and the wine's aging potential.
10. Sulphates: Contribute to SO_2 levels, impacting preservation and possibly enhancing aroma.
11. Alcohol: Significantly influences the wine's body, flavor, and overall quality perception.

Each of these features plays a role in determining the sensory characteristics of the wine, making them vital for predicting its quality.

Handling Missing Data in the Wine Quality Dataset:

In the Kaggle Wine Quality dataset, missing data isn't prominently reported. However, in scenarios where missing values are present, it's essential to address them to maintain the integrity of the analysis. Common imputation techniques include:

1. Mean/Median Imputation: Replacing missing values with the mean or median of the respective feature.

Advantages:

- Simple and quick to implement.
- Preserves the overall distribution of the data.

Disadvantages:

- Can reduce data variability.
- May not be suitable for features with skewed distributions.

2. K-Nearest Neighbors (KNN) Imputation: Utilizes the values of the nearest neighbors to impute missing data.

Advantages:

- Accounts for correlations between features.
- Can produce more accurate imputations when data has underlying patterns.

Disadvantages:

- Computationally intensive, especially with large datasets.
- Sensitive to the choice of 'k' and distance metrics.

3. Multiple Imputation by Chained Equations (MICE): Creates multiple imputed datasets by modeling each feature with missing values as a function of other features.

Advantages:

- Provides a measure of uncertainty around imputed values.
- Suitable for data missing at random.

Disadvantages:

- Complex and requires more computational resources.
- Assumes that the data is missing at random, which might not always be the case.

Choosing the appropriate imputation technique depends on the nature and extent of the missing data, as well as the specific requirements of the analysis.