

## Summary:

This code aims to explore the impact of missing values (NaN) on the performance of different classification algorithms. It involves the following steps:

1. Initial Classification:
  - o Loads a numerical dataset. o Applies SVM, Naive Bayes, and KNN classification algorithms.
  - o Displays confusion matrices to evaluate initial performance.
2. Introducing NaN Values:
  - o Randomly replaces certain values in the dataset with NaN. o Saves the modified dataset with NaN values for reference.
  - o Displays the dataset with NaN values to visualize the modifications.
3. Imputing Missing Values:
  - o Replaces NaN values with the mean of their corresponding columns. o Saves the imputed dataset for further analysis.
  - o Displays the dataset after imputation to verify the changes.
4. Re-Classification:
  - o Reruns the classification process (SVM, Naive Bayes, KNN) on the imputed dataset.
  - o Displays confusion matrices to assess the impact of imputation on classification performance.

## Key Points:

- The code investigates the effect of missing data on classification accuracy.
- It demonstrates the use of mean imputation to handle missing values.
- It compares the performance of different algorithms before and after imputation.
- It provides insights into the robustness of algorithms to missing data and the effectiveness of imputation techniques.

## Additional Considerations:

- Consider exploring other imputation techniques beyond mean imputation.
- Evaluate the impact of different missing data patterns and proportions.
- Investigate the effectiveness of these techniques on various datasets and classification tasks.

## Functions

- `open()`: Opens a file.
- `read_csv()`: Reads a CSV file and returns a DataFrame object.

- `iloc()`: Selects a row or column range from a DataFrame object.
- `split()`: Splits a dataset into training and test sets.
- `fit()`: Trains a classifier.
- `predict()`: Makes predictions using a classifier.
- `confusion_matrix()`: Returns a matrix that compares classification predictions to the true classes.

## Parameters

- `file_path`: The path to the file to open.
- `names`: The names of the columns to read from the file.
- `target_column`: The name of the target column.
- `test_size`: The size of the test set.
- `random_state`: The value used to initialize the random number generator.

## Methods

- `withdraw()`: Hides a Tkinter window.
- `fit_transform()`: Scales the data.

## Code Description

```
# Opens a file dialog to select the dataset
root = tk.Tk() root.withdraw() # hides the
Tkinter window
file_path =
filedialog.askopenfilename()
```

This code allows the user to select the dataset file. The `Tkinter` module is used to create the file dialog. The `root.withdraw()` command hides the Tkinter window.

```
# Loads the dataset df =
pd.read_csv(file_path)
```

```
# Sets the last column as the target
X = df.iloc[:, :-1]
y = df.iloc[:, -1]
```

This code sets the last column of the dataset as the target. The `X` variable contains all columns except the target column. The `y` variable contains the target column.

```
# Splits the dataset into training and test sets X_train,
X_test, y_train, y_test = train_test_split(X, y,
test_size=0.3, random_state=0)
```

This code splits the dataset into training and test sets. The `test_size` parameter specifies the size of the test set. The `random_state` parameter is used to initialize the random number generator.

```
# Normalizes the data
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
```

This code normalizes the dataset. The `StandardScaler()` class normalizes the dataset to have a mean of 0 and a standard deviation of 1.

```
# Defines the classification algorithms
classifiers = {
    'SVM': SVC(),
    'Naive Bayes': GaussianNB(),
    'KNN': KNeighborsClassifier()
}
```

This code defines three different classification algorithms. The `SVC()` class represents the support vector machine (SVM) algorithm. The `GaussianNB()` class represents the Naive Bayes algorithm. The `KNeighborsClassifier()` class represents the nearest neighbors algorithm.

```
# Performs training and prediction for each classifier, and
displays the confusion matrix for name, clf in
classifiers.items():
    clf.fit(X_train, y_train)    y_pred =
clf.predict(X_test)    cm =
confusion_matrix(y_test, y_pred)
print(f"{name} Confusion Matrix:\n{cm}")
```

This code performs training and prediction for each classifier and displays the confusion matrix. The `confusion_matrix()` function returns a matrix that allows for comparison of classification predictions to the true classes.

**Results:** The code compares the classification performance of three different classification algorithms on a dataset. The results show that the SVM algorithm performs the best. The Naive Bayes algorithm performs second best, while the KNN algorithm performs the worst.