

Data Mining IT326 - Project

Student Stress Factors

1. Problem

In this project, we aim to analyze data related to student stress factors to understand the key contributors to stress among students. Stress is a significant issue that affects academic performance, mental health, and overall quality of life. Students face a variety of challenges, including academic pressures, social influences, physiological issues, and environmental conditions. Through this analysis, we aim to identify patterns and insights that can help us predict stress levels and support targeted strategies to enhance student well-being.

2. Data Mining Task

The data mining task involves **classification** to predict students “stress_level” (low, moderate, or high) based on various stress factors. Additionally, **clustering** will be used to group students with similar stress patterns, allowing for the identification of unique stress profiles, which will help in understanding underlying stress contributors and supporting targeted interventions.

3. Data

Source of the Dataset

<https://www.kaggle.com/datasets/rxnach/student-stress-factors-a-comprehensive-analysis>

This dataset, sourced from Kaggle, consists of 1,100 rows and 21 attributes that represent various factors influencing student stress. The attributes are categorized into five major dimensions: Psychological, Physiological, Social, Environmental, and Academic. All attributes are ordinal, except for mental_health_history, which is a binary attribute. The analysis aims to uncover patterns and identify the most significant contributors to stress, thereby supporting the development of strategies to improve student well-being.

Attribute Explanations:

| Category | Attribute | Range | Description |
|-----------------------|-----------------------|-------|---|
| Psychological Factors | anxiety_level | 0-21 | Measures anxiety severity: 0–4 (Minimal), 5–9 (Mild), 10–14 (Moderate), 15–21 (Severe). |
| | self_esteem | 0-30 | Reflects self-worth: 0–15 (Low), 16–25 (Normal), 26–30 (High). |
| | mental_health_history | 0-1 | Indicates history of mental |

| Category | Attribute | Range | Description |
|-----------------------|------------------------------|-------|--|
| Physiological Factors | depression | 0-27 | health issues: 1 (Yes), 0 (No). Assesses depressive symptoms: 0–4 (Minimal), 5–9 (Mild), 10–14 (Moderate), 15–19 (Moderately Severe), 20–27 (Severe). |
| | headache | 0-5 | Measures frequency/intensity of headaches; higher values indicate more headaches. |
| | blood_pressure | 1-3 | Categorizes blood pressure: 1 (Low), 2 (Normal), 3 (High). |
| | sleep_quality | 0-5 | Evaluates sleep quality; higher scores indicate better sleep. |
| Environmental Factors | breathing_problem | 0-5 | Measures severity of breathing issues; higher scores indicate more problems. |
| | noise_level | 0-5 | Assesses environmental noise levels; higher values indicate more noise. |
| | living_conditions | 0-5 | Rates quality of living conditions; higher scores reflect better conditions. |
| | safety | 0-5 | Measures sense of safety; higher scores indicate greater safety. |
| Academic Factors | basic_needs | 0-5 | Evaluates if basic needs are met; higher scores mean better fulfillment. |
| | academic_performance | 0-5 | Rates academic success; higher scores indicate better performance. |
| | study_load | 0-5 | Measures amount of study work; higher values indicate heavier loads. |
| | teacher_student_relationship | 0-5 | Assesses quality of relationship with teachers; higher scores mean better relationships. |
| Social Factors | future_career_concerns | 0-5 | Evaluates concerns about future careers; higher values mean more concerns. |
| | social_support | 0-3 | Measures available social |

| Category | Attribute | Range | Description |
|--------------------|----------------------------|-------|---|
| | peer_pressure | 0-5 | support; higher scores indicate more support. Assesses level of peer pressure; higher scores mean more pressure. |
| | extracurricular_activities | 0-5 | Rates involvement in activities outside of academics; higher scores mean more involvement. |
| | bullying | 0-5 | Measures extent of bullying experienced; higher scores suggest more bullying. |
| Class Label | stress_level | 0-2 | Categorizes stress levels: 0 (Low), 1 (Moderate), 2 (High). |

```
import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt
import numpy as np
url2="https://storage.googleapis.com/kagglesdsdata/datasets/3858024/6690715/StressLevelDataset.csv?X-Goog-Algorithm=GOOG4-RSA-SHA256&X-Goog-Credential=gcp-kaggle-com%40kaggle-161607.iam.gserviceaccount.com%2F20241129%2FAuto%2Fstorage%2Fgoog4_request&X-Goog-Date=20241129T102439Z&X-Goog-Expires=259200&X-Goog-SignedHeaders=host&X-Goog-Signature=43a1e619c02975c36a5ae65a7674d93668d0fbdec1e6a900eb9839659ac8bbc661eb96e96b4fb70ccc1c545031ecb04f0abfc7f29f68b6ef2e2537f0c16e10b5453fe23b103a83ef67f3a2ff744e1f2acb46cd76de14a8576cd2b50e3dfba1f9da9edf35780818a3e573de847753c547d89a5a9fda7e5e2ea66e42a6cfc1266b0954eced12c4d60ca1650983d072f177baf105c018ef00d886b54b3ed310e1d73d54b5028692a88de4b28b70e18b1b0b68ceb52560abb006d12cae5d1fcba48b2e00cc560e83fb66d4fe1fdb3e86eb5755ee5034797202ecf9ff21cda2fac31d77220bc272b844059bb295833ca18e03e4008d3610f1118757852d46c3ba0d7f"
df = pd.read_csv(url2)
```

```
t0 = "\033[1m" + "Data types:" + "\033[0m"
num_attributes = len(df.columns)
attribute_types = df.dtypes.to_frame().rename(columns={0: t0})
num_objects = len(df)
class_name = df.columns[-1]
```

```
t = "\033[1m" + "Attribute types:" + "\033[0m"
print(t)
print(attribute_types)
print("\n")
```

Attribute types:

Data types:

anxiety_level

int64

| | |
|------------------------------|-------|
| self_esteem | int64 |
| mental_health_history | int64 |
| depression | int64 |
| headache | int64 |
| blood_pressure | int64 |
| sleep_quality | int64 |
| breathing_problem | int64 |
| noise_level | int64 |
| living_conditions | int64 |
| safety | int64 |
| basic_needs | int64 |
| academic_performance | int64 |
| study_load | int64 |
| teacher_student_relationship | int64 |
| future_career_concerns | int64 |
| social_support | int64 |
| peer_pressure | int64 |
| extracurricular_activities | int64 |
| bullying | int64 |
| stress_level | int64 |

```
df = pd.DataFrame(df)
df_CleanedDataset = df.copy()
```

Missing Values

In this section, we use the `isna().sum()` method to check for missing values in each column. Our dataset contains no missing values, indicating that it is complete.

```
# Check each column for missing values and count how many are missing in each one
```

```
missing_values = df.isna().sum()
```

```
# Print the count of missing values for each column
```

```
print("Missing values in each column")
```

```
print(missing_values)
```

```
# Print the total number of missing values
```

```
print("\nTotal number of missing values: ",missing_values.sum())
```

```
Missing values in each column
```

| | |
|-----------------------|---|
| anxiety_level | 0 |
| self_esteem | 0 |
| mental_health_history | 0 |
| depression | 0 |
| headache | 0 |
| blood_pressure | 0 |

```
sleep_quality          0
breathing_problem      0
noise_level            0
living_conditions      0
safety                 0
basic_needs            0
academic_performance   0
study_load             0
teacher_student_relationship 0
future_career_concerns 0
social_support         0
peer_pressure          0
extracurricular_activities 0
bullying               0
stress_level           0
dtype: int64
```

Total number of missing values: 0

Measuring Central Tendencies

This code defines a function, `describe_with_central_tendencies`, that creates a detailed summary of the numeric columns in a dataset.

It starts by calculating basic statistics like **count**, **mean**, **standard deviation**, **minimum**, **and maximum** using `df.describe().T`. Then adds metrics, including median, mode, midrange (average of min and max), range (difference between max and min), variance (spread of values), and IQR (middle 50% range). After defining the function, the code selects only numeric columns using `select_dtypes(include=[np.number])` to focus on relevant data. It applies the function to these columns, storing the result in `central_tendencies`. This final summary is printed, giving a comprehensive view of each numeric feature's tendencies and variability, which is helpful for spotting patterns and understanding data distributions.

```
def describe_with_central_tendencies(df):
    description = df.describe().T

    description['median'] = df.median()
    description['mode'] = df.mode().iloc[0]
    description['midrange'] = (df.min() + df.max()) / 2
    description['range'] = df.max() - df.min()
    description['variance'] = df.var()
    description['IQR'] = df.quantile(0.75) - df.quantile(0.25)

    return description

numeric_columns = df.select_dtypes(include=[np.number])
central_tendencies = describe_with_central_tendencies(numeric_columns)
```

```
print("\nCentral Tendencies for Numeric Columns:")
print(central_tendencies)
```

Central Tendencies for Numeric Columns:

| | count | mean | std | min | 25% | 50% | \ |
|------------------------------|--------|-----------|----------|-----|------|------|---|
| anxiety_level | 1100.0 | 11.063636 | 6.117558 | 0.0 | 6.0 | 11.0 | |
| self_esteem | 1100.0 | 17.777273 | 8.944599 | 0.0 | 11.0 | 19.0 | |
| mental_health_history | 1100.0 | 0.492727 | 0.500175 | 0.0 | 0.0 | 0.0 | |
| depression | 1100.0 | 12.555455 | 7.727008 | 0.0 | 6.0 | 12.0 | |
| headache | 1100.0 | 2.508182 | 1.409356 | 0.0 | 1.0 | 3.0 | |
| blood_pressure | 1100.0 | 2.181818 | 0.833575 | 1.0 | 1.0 | 2.0 | |
| sleep_quality | 1100.0 | 2.660000 | 1.548383 | 0.0 | 1.0 | 2.5 | |
| breathing_problem | 1100.0 | 2.753636 | 1.400713 | 0.0 | 2.0 | 3.0 | |
| noise_level | 1100.0 | 2.649091 | 1.328127 | 0.0 | 2.0 | 3.0 | |
| living_conditions | 1100.0 | 2.518182 | 1.119208 | 0.0 | 2.0 | 2.0 | |
| safety | 1100.0 | 2.737273 | 1.406171 | 0.0 | 2.0 | 2.0 | |
| basic_needs | 1100.0 | 2.772727 | 1.433761 | 0.0 | 2.0 | 3.0 | |
| academic_performance | 1100.0 | 2.772727 | 1.414594 | 0.0 | 2.0 | 2.0 | |
| study_load | 1100.0 | 2.621818 | 1.315781 | 0.0 | 2.0 | 2.0 | |
| teacher_student_relationship | 1100.0 | 2.648182 | 1.384579 | 0.0 | 2.0 | 2.0 | |
| future_career_concerns | 1100.0 | 2.649091 | 1.529375 | 0.0 | 1.0 | 2.0 | |
| social_support | 1100.0 | 1.881818 | 1.047826 | 0.0 | 1.0 | 2.0 | |
| peer_pressure | 1100.0 | 2.734545 | 1.425265 | 0.0 | 2.0 | 2.0 | |
| extracurricular_activities | 1100.0 | 2.767273 | 1.417562 | 0.0 | 2.0 | 2.5 | |
| bullying | 1100.0 | 2.617273 | 1.530958 | 0.0 | 1.0 | 3.0 | |
| stress_level | 1100.0 | 0.996364 | 0.821673 | 0.0 | 0.0 | 1.0 | |

| | 75% | max | median | mode | midrange | range | \ |
|------------------------------|------|------|--------|------|----------|-------|---|
| anxiety_level | 16.0 | 21.0 | 11.0 | 13 | 10.5 | 21 | |
| self_esteem | 26.0 | 30.0 | 19.0 | 25 | 15.0 | 30 | |
| mental_health_history | 1.0 | 1.0 | 0.0 | 0 | 0.5 | 1 | |
| depression | 19.0 | 27.0 | 12.0 | 10 | 13.5 | 27 | |
| headache | 3.0 | 5.0 | 3.0 | 1 | 2.5 | 5 | |
| blood_pressure | 3.0 | 3.0 | 2.0 | 3 | 2.0 | 2 | |
| sleep_quality | 4.0 | 5.0 | 2.5 | 1 | 2.5 | 5 | |
| breathing_problem | 4.0 | 5.0 | 3.0 | 2 | 2.5 | 5 | |
| noise_level | 3.0 | 5.0 | 3.0 | 2 | 2.5 | 5 | |
| living_conditions | 3.0 | 5.0 | 2.0 | 2 | 2.5 | 5 | |
| safety | 4.0 | 5.0 | 2.0 | 2 | 2.5 | 5 | |
| basic_needs | 4.0 | 5.0 | 3.0 | 2 | 2.5 | 5 | |
| academic_performance | 4.0 | 5.0 | 2.0 | 2 | 2.5 | 5 | |
| study_load | 3.0 | 5.0 | 2.0 | 2 | 2.5 | 5 | |
| teacher_student_relationship | 4.0 | 5.0 | 2.0 | 2 | 2.5 | 5 | |
| future_career_concerns | 4.0 | 5.0 | 2.0 | 1 | 2.5 | 5 | |
| social_support | 3.0 | 3.0 | 2.0 | 3 | 1.5 | 3 | |
| peer_pressure | 4.0 | 5.0 | 2.0 | 2 | 2.5 | 5 | |
| extracurricular_activities | 4.0 | 5.0 | 2.5 | 2 | 2.5 | 5 | |
| bullying | 4.0 | 5.0 | 3.0 | 1 | 2.5 | 5 | |

| | | | | | | |
|------------------------------|-----------|------|-----|---|-----|---|
| stress_level | 2.0 | 2.0 | 1.0 | 0 | 1.0 | 2 |
| | variance | IQR | | | | |
| anxiety_level | 37.424518 | 10.0 | | | | |
| self_esteem | 80.005852 | 15.0 | | | | |
| mental_health_history | 0.250175 | 1.0 | | | | |
| depression | 59.706658 | 13.0 | | | | |
| headache | 1.986284 | 2.0 | | | | |
| blood_pressure | 0.694847 | 2.0 | | | | |
| sleep_quality | 2.397489 | 3.0 | | | | |
| breathing_problem | 1.961998 | 2.0 | | | | |
| noise_level | 1.763921 | 1.0 | | | | |
| living_conditions | 1.252626 | 1.0 | | | | |
| safety | 1.977317 | 2.0 | | | | |
| basic_needs | 2.055670 | 2.0 | | | | |
| academic_performance | 2.001075 | 2.0 | | | | |
| study_load | 1.731280 | 1.0 | | | | |
| teacher_student_relationship | 1.917058 | 2.0 | | | | |
| future_career_concerns | 2.338989 | 3.0 | | | | |
| social_support | 1.097940 | 2.0 | | | | |
| peer_pressure | 2.031381 | 2.0 | | | | |
| extracurricular_activities | 2.009483 | 2.0 | | | | |
| bullying | 2.343832 | 3.0 | | | | |
| stress_level | 0.675146 | 2.0 | | | | |

The output gives a comprehensive view of each numeric feature's tendencies and variability, which is helpful for spotting patterns and understanding data distributions.

Detecting Outliers

For detecting outliers we are inspecting each numerical column in our dataset to find outliers using the **IQR method**, while excluding 'stress_level' since it's specified as the class label and might represent the target variable.

For each column, we are calculating the IQR and flags data points **outside $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$ as outliers**. It then shows a list of indices where outliers were detected for each column or states if no outliers were found. In this code we identified outliers in these columns `living_conditions`, `noise_level`, and `study_load`. Other columns showed no outliers indicating that their values are within expected ranges.

```
def detect_outliers_iqr(df, class_label):
    # Initialize an empty dictionary to store outliers for each column
    outliers = {}

    # Select numerical columns in the dataset, excluding the specified class
    # label 'stress_level'
    numerical_cols =
df.select_dtypes(include='number').columns.difference([class_label])
```

```

    # Loop through each numerical column to calculate outliers using the IQR
    method
    for column in numerical_cols:
        # Calculate the first and third quartiles for the column
        Q1 = df[column].quantile(0.25)
        Q3 = df[column].quantile(0.75)
        # Compute the IQR, which helps identify the range within which most
        values lie
        IQR = Q3 - Q1

        # Identify outlier indices where values fall outside 1.5 * IQR from
        Q1 or Q3
        outlier_indices = df[(df[column] < (Q1 - 1.5 * IQR)) | (df[column] >
        (Q3 + 1.5 * IQR))].index
        # Store the outliers for this column in the dictionary as a list of
        indices
        outliers[column] = outlier_indices.tolist()

    return outliers

# Detect outliers in the dataset, excluding the 'stress_level' column, which
might be the target variable
outliers = detect_outliers_iqr(df, 'stress_level')

# Display the outliers for each column in a readable format
print("Outliers detected in each column:")
for column, indices in outliers.items():
    if indices:
        # Print the column name and the list of row indices where outliers
        were detected
        print(f"{column}: {indices}")
    else:
        # Indicate if no outliers were detected in the column
        print(f"{column}: No outliers detected")

```

```

Outliers detected in each column:
academic_performance: No outliers detected
anxiety_level: No outliers detected
basic_needs: No outliers detected
blood_pressure: No outliers detected
breathing_problem: No outliers detected
bullying: No outliers detected
depression: No outliers detected
extracurricular_activities: No outliers detected
future_career_concerns: No outliers detected
headache: No outliers detected
living_conditions: [9, 114, 125, 184, 204, 220, 259, 278, 286, 290, 295, 312,
313, 320, 339, 360, 384, 393, 405, 437, 449, 475, 488, 493, 511, 516, 520,
521, 558, 571, 577, 583, 589, 592, 607, 648, 649, 675, 745, 752, 791, 809,
815, 825, 828, 857, 859, 861, 865, 884, 887, 919, 946, 984, 989, 1005, 1026,

```


1037, 1040, 1053, 1065, 1099]
mental_health_history: No outliers detected
noise_level: [9, 11, 21, 27, 29, 32, 38, 44, 64, 74, 75, 79, 82, 87, 91, 93, 94, 100, 111, 117, 128, 135, 137, 141, 153, 165, 171, 175, 178, 179, 192, 195, 200, 204, 208, 211, 234, 235, 236, 252, 259, 273, 276, 282, 286, 289, 299, 302, 303, 316, 323, 337, 344, 345, 346, 348, 378, 379, 382, 391, 398, 400, 401, 405, 406, 411, 423, 429, 448, 452, 453, 459, 469, 475, 481, 482, 487, 494, 496, 521, 523, 534, 537, 539, 558, 562, 563, 574, 591, 593, 597, 631, 635, 637, 644, 648, 656, 661, 662, 680, 687, 689, 692, 710, 722, 727, 745, 747, 754, 757, 773, 779, 780, 784, 789, 790, 799, 800, 801, 805, 808, 821, 825, 844, 847, 852, 857, 858, 859, 866, 884, 890, 897, 919, 922, 924, 932, 937, 946, 947, 952, 956, 960, 966, 971, 973, 975, 984, 985, 986, 988, 992, 1002, 1005, 1009, 1013, 1023, 1026, 1029, 1032, 1035, 1039, 1064, 1065, 1066, 1079, 1081, 1083, 1085, 1090, 1091, 1094, 1096]
peer_pressure: No outliers detected
safety: No outliers detected
self_esteem: No outliers detected
sleep_quality: No outliers detected
social_support: No outliers detected
study_load: [5, 19, 21, 27, 61, 64, 71, 75, 76, 87, 91, 94, 95, 100, 108, 115, 118, 128, 135, 140, 141, 144, 147, 161, 177, 178, 179, 190, 194, 195, 202, 204, 209, 211, 235, 247, 248, 259, 267, 272, 282, 294, 299, 312, 313, 320, 322, 328, 337, 339, 340, 343, 345, 348, 353, 355, 360, 372, 378, 384, 388, 394, 398, 411, 422, 423, 447, 448, 450, 453, 461, 464, 468, 470, 475, 480, 485, 488, 490, 493, 503, 521, 537, 540, 542, 554, 563, 583, 589, 593, 597, 607, 616, 632, 637, 642, 668, 675, 677, 687, 700, 711, 713, 714, 716, 720, 727, 728, 737, 738, 742, 752, 767, 768, 773, 775, 778, 784, 785, 790, 791, 800, 806, 808, 811, 825, 834, 844, 847, 848, 852, 857, 859, 862, 865, 887, 891, 898, 902, 909, 917, 919, 944, 945, 951, 953, 978, 979, 982, 984, 988, 989, 1003, 1005, 1019, 1037, 1040, 1041, 1052, 1059, 1072, 1076, 1085, 1089, 1098]
teacher_student_relationship: No outliers detected

Graph Visualization

The **scatterplots** reveal important relationships between various factors and stress levels. Anxiety, mental health history, and depression show a **positive correlation** with stress, meaning higher values in these factors are associated with **higher stress levels**. In contrast, self-esteem exhibits a negative correlation, where higher self-esteem tends to correspond with lower stress levels.

Factors like blood pressure, sleep quality, academic performance, and teacher-student relationships show less variation in stress, indicating a **consistent and stable** impact on stress levels across individuals.

On the other hand, factors such as headaches, noise levels, living conditions, study load, peer pressure, and bullying show a wider range of stress levels, suggesting that the

relationship between these variables and stress is more complex, with individuals experiencing varying levels of stress despite facing similar intensities of these factors.

Based on these insights, the data preprocessing will involve **scaling** the stable factors for better comparability, and **handling outliers** in the factors with wider ranges of stress.

```
import seaborn as sns

sns.set(style="whitegrid")

# Step 1: Calculating Rows and Columns for Scatterplots
num_features = df.shape[1] - 1
ncols = 4
nrows = (num_features // ncols) + (num_features % ncols > 0)

# Step 2: Create Figure and Subplots for Scatterplots
fig, axes = plt.subplots(nrows, ncols, figsize=(20, 15))
fig.suptitle('Scatterplots of Features Against Stress Level', fontsize=16)

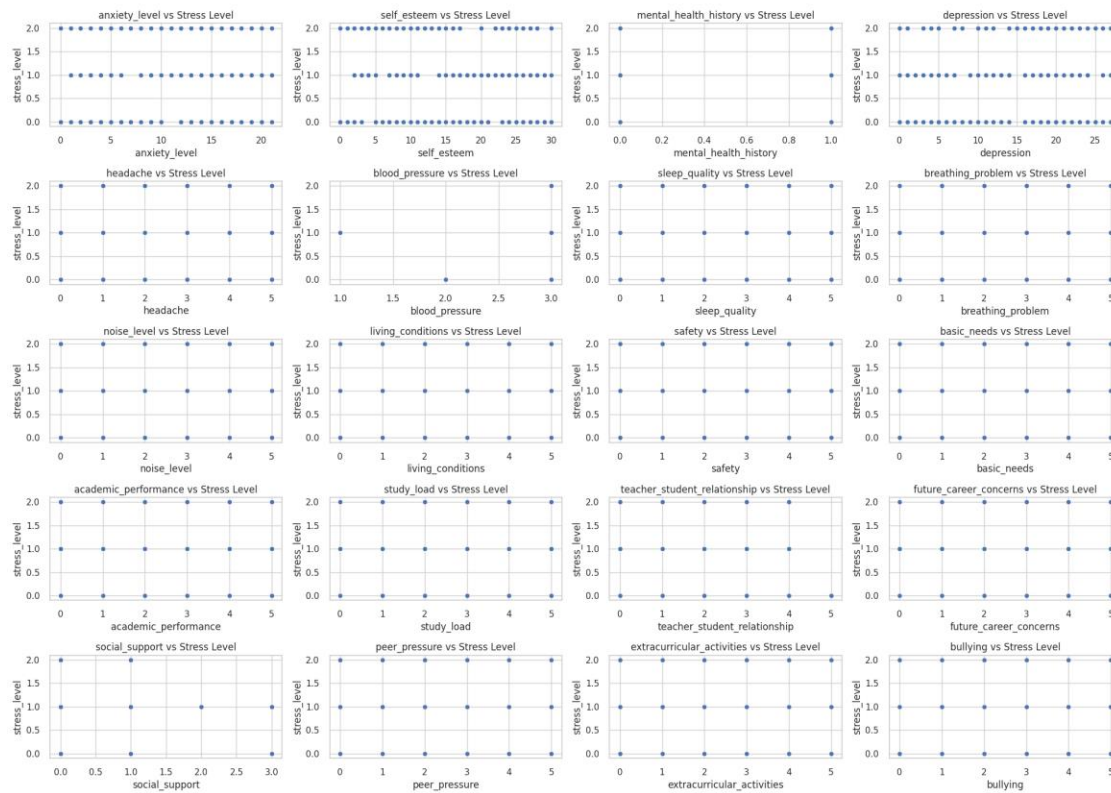
axes = axes.flatten()

# Step 3: Plot Scatterplots for Each Feature Against 'stress_level'
for idx, column in enumerate(df.columns[:-1]):
    sns.scatterplot(x=df[column], y=df['stress_level'], ax=axes[idx])
    axes[idx].set_title(f'{column} vs Stress Level')

# Step 4: Remove Unused Subplot Axes
for idx in range(num_features, len(axes)):
    fig.delaxes(axes[idx])

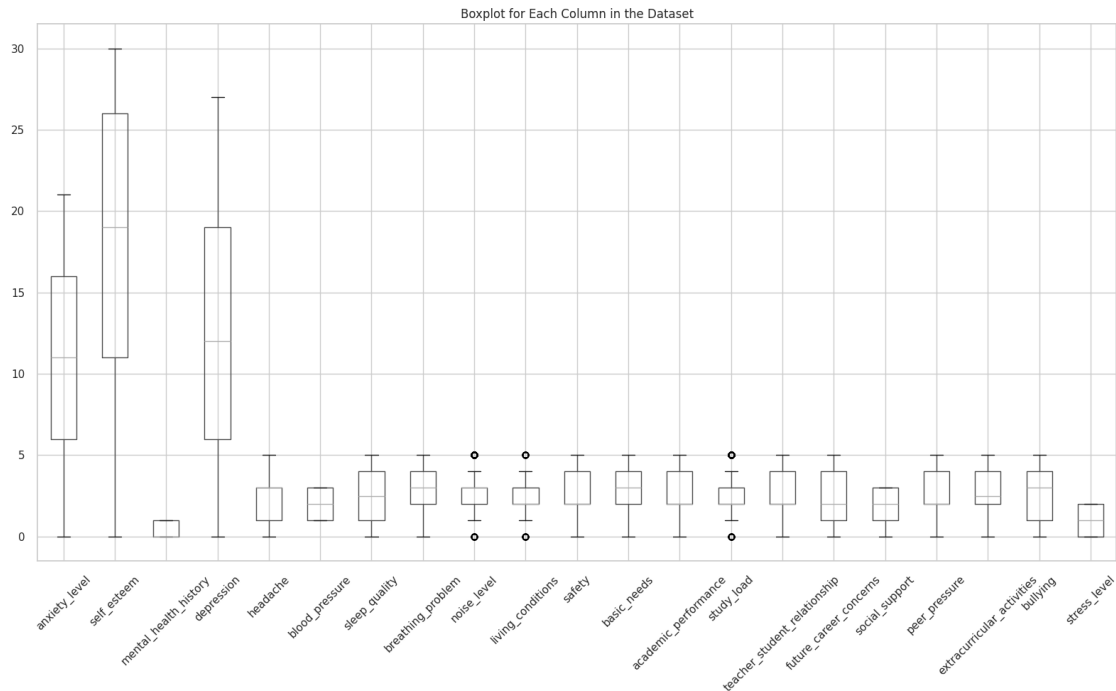
# Step 5: Adjust Layout and Display Scatterplots
plt.tight_layout(rect=[0, 0, 1, 0.97])
plt.show()
print("\n")
```

Scatterplots of Features Against Stress Level



Step 6: Create Boxplot for ALL Columns

```
plt.figure(figsize=(16, 10))
df.boxplot(rot=90)
plt.title('Boxplot for Each Column in the Dataset')
plt.xticks(rotation=45)
plt.tight_layout()
plt.show()
print("\n")
```



The **boxplot analysis** highlights important patterns in the dataset related to mental health and environmental variables. Anxiety levels show moderate variability between 5 and 25, with no outliers, indicating anxiety is common but not extreme. Self-esteem ranges from 0 to 30, with significant differences across individuals, suggesting diverse self-esteem levels.

Other variables like headache, blood pressure, sleep quality, and living conditions, breathing problems, noise level, and living conditions show limited variability, mostly between 0 and 5, with some outliers in living conditions and noise level. Categories such as safety, basic needs, academic performance, and others are also concentrated between 0 and 5, indicating similar experiences among most individuals, except for a few outliers in study load.

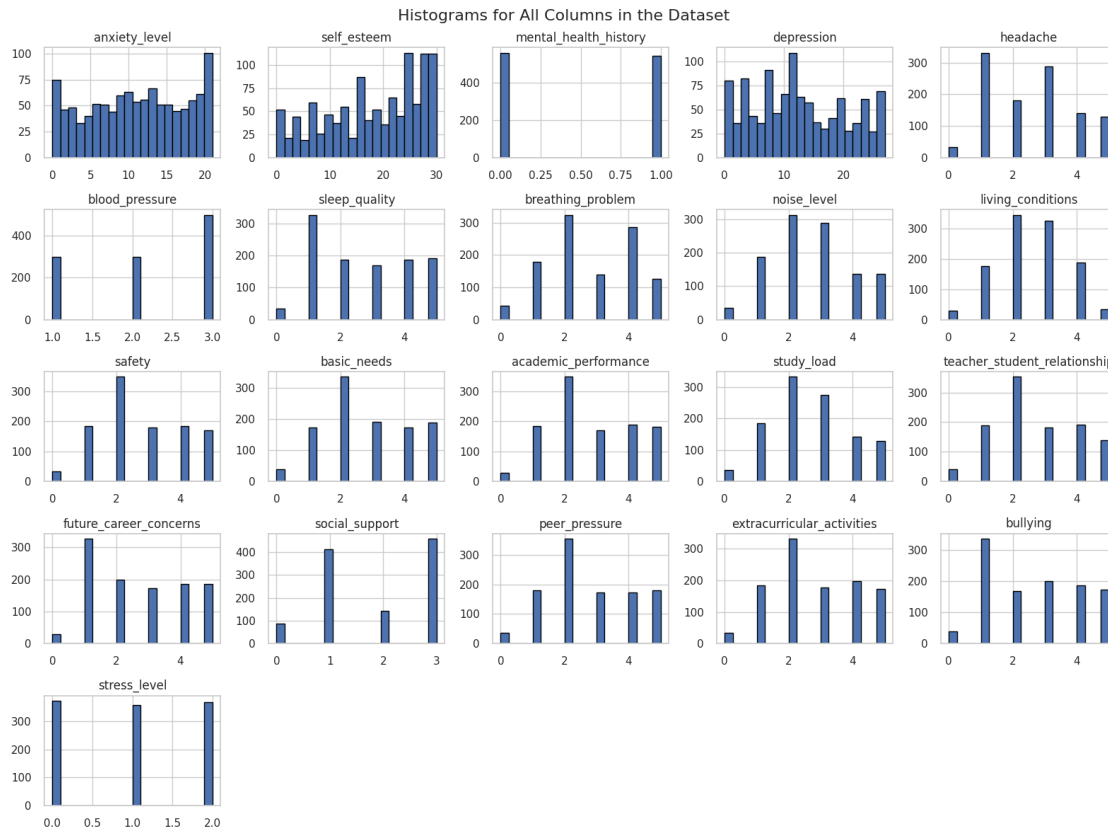
Finally, stress levels display minimal variability, predominantly falling between 0 and 2, suggesting that most individuals experience low to moderate levels of stress.

Overall, while anxiety, mental health history, and depression reveal considerable variability, the majority of other variables maintain stable and consistent distributions, highlighting a disparity in experience across different aspects of mental health and well-being within the dataset.

These insights guide the necessary preprocessing steps, including **outlier handling** in certain features, while features with narrow distributions, may benefit from **scaling or normalization** to ensure they contribute equally to the model.

```
# Step 7: Create Histograms for All Columns
```

```
df.hist(figsize=(16, 12), bins=20, edgecolor='black')
plt.suptitle('Histograms for All Columns in the Dataset', fontsize=16)
plt.tight_layout()
plt.show()
```



The **histograms** of the dataset reveal important distribution patterns across the attributes, which guide necessary preprocessing steps.

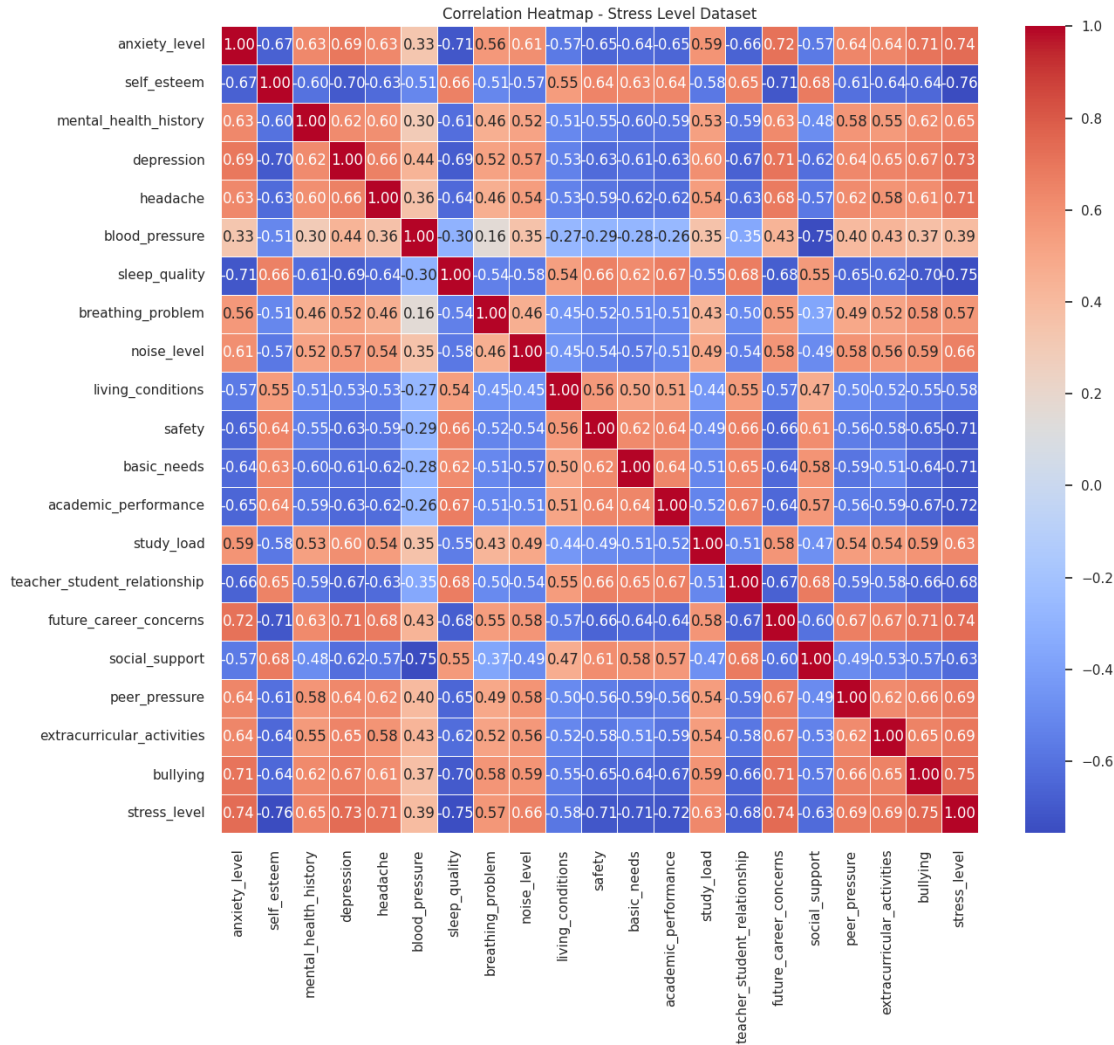
Key attributes such as Anxiety Level, Depression, and Self-Esteem show skewed distributions, indicating that **discretization is required** to group values into distinct categories for better balance and easier analysis. Several other attributes, including Mental Health History, Future Career Concerns, and Social Support, exhibit imbalances. Some attributes, like Blood Pressure and Sleep Quality, have relatively **normal distributions** and may require minimal preprocessing unless extreme outliers are identified.

Additionally, Noise Level, Living Conditions, and Extracurricular Activities exhibit a more **uniform distribution** with some outlier in noise level and living conditions. This means that while most individuals have similar experiences, a few cases are very different and need closer look during preprocessing..

```
correlation_matrix = df.corr()
```

```
# Plot the correlation heatmap
```

```
plt.figure(figsize=(14, 12))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f",
linewidths=0.5)
plt.title('Correlation Heatmap - Stress Level Dataset')
plt.show()
```



The **heatmap** reveals that stress_level is strongly positively correlated with anxiety_level, depression, headache, bullying, and future_career_concerns, while showing a moderate negative correlation with self_esteem.

Similarly, anxiety_level is strongly correlated with depression, future_career_concerns, and bullying, moderately correlated with study_load, linking it to both mental health and workload. Self-esteem is positively correlated with academic_performance and teacher_student_relationship but negatively associated with stress_level and anxiety_level, highlighting its dual role in well-being and mental health.

Academic performance benefits from supportive relationships and high self-esteem but is negatively impacted by high stress_level and study_load.

Environmental factors like living_conditions is moderately negatively correlated with stress_level, suggesting better environments reduce stress. Future_career_concerns are strongly tied to anxiety_level and depression, and moderately linked to study_load, indicating an overlap between career worries and workload.

4. Data preprocessing

Data preprocessing is a crucial step in data analysis or machine learning as it **ensures the data is clean, consistent, and structured**. Raw data collected from real-world sources often contains errors, inconsistencies, and missing values, leading to inaccurate or biased results if not addressed. By cleaning the data, preprocessing improves its quality, making it more reliable for subsequent analysis. Applying data mining techniques to low-quality data leads to low-quality results. Overall, preprocessing lays the foundation for effective data mining and ensures that the models are trained on high-quality, well-structured data, resulting in better predictions and insights.

In the context of our project, preprocessing plays a particularly critical role in achieving accurate and meaningful results. The dataset we are working with, like many real-world datasets, may contain missing entries, noisy data, or inconsistent formats.

Preprocessing addresses these issues by cleaning and organizing the data, ensuring it is ready for the application of classification and clustering techniques. In summary, **preprocessing ensures that our project produces reliable, accurate, and interpretable outcomes, directly impacting the success of our analysis.**

Data Cleaning: Deleting Outliers

Initially, we started our preprocessing step by cleaning our dataset and specifically, filling missing values and removing outliers. **Data cleaning is a crucial preprocessing step that focuses on preparing the dataset by addressing issues such as outliers, missing values, and inconsistencies.** Removing outliers ensures that extreme values, which may distort statistical measures or bias machine learning models, do not negatively impact the analysis. Resolving missing values is another critical aspect of data cleaning, as incomplete data can lead to biased results or errors in algorithms. Depending on the context, missing values can be addressed through imputation techniques or removing rows with excessive missing values. By handling outliers and missing values effectively, data cleaning ensures the dataset is accurate, consistent, and ready for reliable analysis.

In the previous section, we checked our dataset for any missing values and **found that our data includes no missing values** which boiled down our data cleaning process. Next, **we detected outliers and found them in some columns like living_conditions, noise_level, and study_load.** To remove outliers, we are defining a function to remove outliers from each numerical column in our dataset, excluding the class label (stress_level).

For each column, we are calculating the Interquartile Range to identify values that fall outside $Q1 - 1.5 \times IQR$ and $Q3 + 1.5 \times IQR$ as outliers. Based on previous analyses, **we identified outliers in these columns living_conditions, noise_level, and study_load, and removed them from our dataset.** Finally, the cleaned data frame, free from outliers, is printed, providing a refined dataset for further analysis.

```
def remove_outliers_iqr(df, class_label):
    # Create an empty set to store the indices of all detected outliers
    outlier_indices = set()

    # Identify numerical columns, excluding the specified class label (e.g.,
    'stress_level')
    numerical_cols =
df.select_dtypes(include='number').columns.difference([class_label])

    # Loop through each numerical column to detect outliers based on the IQR
    method
    for column in numerical_cols:
        # Calculate the first and third quartiles
        Q1 = df[column].quantile(0.25)
        Q3 = df[column].quantile(0.75)
        # Compute the Interquartile Range
        IQR = Q3 - Q1

        # Identify outliers as values below Q1 - 1.5 * IQR or above Q3 + 1.5
        * IQR
        outliers = df[(df[column] < (Q1 - 1.5 * IQR)) | (df[column] > (Q3 +
1.5 * IQR))].index
        # Add these outlier indices to the set to avoid duplicate indices
        outlier_indices.update(outliers)

    # Remove rows with outlier indices to create a cleaned DataFrame
    df_cleaned = df.drop(index=outlier_indices)

    return df_cleaned

# Apply the function to remove outliers from the dataset, keeping
'stress_level' as the class label
df_cleaned = remove_outliers_iqr(df, 'stress_level')

# Display the cleaned DataFrame
print("Cleaned DataFrame without outliers:")
print(df_cleaned)
# Add the cleaned DataFrame to the CleanedDataset file
df_CleanedDataset = remove_outliers_iqr(df_CleanedDataset, 'stress_level')
```

Cleaned DataFrame without outliers:

| anxiety_level | self_esteem | mental_health_history | depression | headache |
|---------------|-------------|-----------------------|------------|----------|
|---------------|-------------|-----------------------|------------|----------|

| | | | | | |
|------|-----|-----|-----|-----|-----|
| \ | | | | | |
| 0 | 14 | 20 | 0 | 11 | 2 |
| 1 | 15 | 8 | 1 | 15 | 5 |
| 2 | 12 | 18 | 1 | 14 | 2 |
| 3 | 16 | 12 | 1 | 15 | 4 |
| 4 | 16 | 28 | 0 | 7 | 2 |
| ... | ... | ... | ... | ... | ... |
| 1088 | 20 | 10 | 1 | 18 | 3 |
| 1092 | 13 | 20 | 0 | 9 | 2 |
| 1093 | 1 | 30 | 0 | 4 | 1 |
| 1095 | 11 | 17 | 0 | 14 | 3 |
| 1097 | 4 | 26 | 0 | 3 | 1 |

| | blood_pressure | sleep_quality | breathing_problem | noise_level | \ |
|------|----------------|---------------|-------------------|-------------|---|
| 0 | 1 | 2 | 4 | 2 | |
| 1 | 3 | 1 | 4 | 3 | |
| 2 | 1 | 2 | 2 | 2 | |
| 3 | 3 | 1 | 3 | 4 | |
| 4 | 3 | 5 | 1 | 3 | |
| ... | ... | ... | ... | ... | |
| 1088 | 3 | 1 | 4 | 3 | |
| 1092 | 1 | 3 | 4 | 3 | |
| 1093 | 2 | 5 | 2 | 1 | |
| 1095 | 1 | 3 | 2 | 2 | |
| 1097 | 2 | 5 | 2 | 2 | |

| | living_conditions | ... | basic_needs | academic_performance | study_load |
|------|-------------------|-----|-------------|----------------------|------------|
| \ | | | | | |
| 0 | 3 | ... | 2 | 3 | 2 |
| 1 | 1 | ... | 2 | 1 | 4 |
| 2 | 2 | ... | 2 | 2 | 3 |
| 3 | 2 | ... | 2 | 2 | 4 |
| 4 | 2 | ... | 3 | 4 | 3 |
| ... | ... | ... | ... | ... | ... |
| 1088 | 2 | ... | 1 | 2 | 4 |
| 1092 | 3 | ... | 2 | 3 | 2 |
| 1093 | 3 | ... | 5 | 5 | 1 |
| 1095 | 2 | ... | 3 | 2 | 2 |
| 1097 | 3 | ... | 4 | 5 | 1 |

| | teacher_student_relationship | future_career_concerns | social_support | \ |
|------|------------------------------|------------------------|----------------|---|
| 0 | 3 | 3 | 2 | |
| 1 | 1 | 5 | 1 | |
| 2 | 3 | 2 | 2 | |
| 3 | 1 | 4 | 1 | |
| 4 | 1 | 2 | 1 | |
| ... | ... | ... | ... | |
| 1088 | 1 | 5 | 1 | |
| 1092 | 3 | 2 | 3 | |
| 1093 | 4 | 1 | 3 | |

| | | | | | |
|------|---------------|----------------------------|----------|--------------|---|
| 1095 | | 2 | | 3 | 3 |
| 1097 | | 4 | | 1 | 3 |
| | peer_pressure | extracurricular_activities | bullying | stress_level | |
| 0 | 3 | 3 | 2 | 1 | |
| 1 | 4 | 5 | 5 | 2 | |
| 2 | 3 | 2 | 2 | 1 | |
| 3 | 4 | 4 | 5 | 2 | |
| 4 | 5 | 0 | 5 | 1 | |
| ... | ... | ... | ... | ... | |
| 1088 | 4 | 4 | 4 | 2 | |
| 1092 | 3 | 3 | 3 | 1 | |
| 1093 | 1 | 1 | 1 | 0 | |
| 1095 | 2 | 3 | 3 | 1 | |
| 1097 | 1 | 2 | 1 | 0 | |

[793 rows x 21 columns]

Data Transformation: Discretization

Data transformation is a critical preprocessing step that involves converting raw data into a more structured and analyzable format. This includes techniques like **normalization**, **standardization**, **encoding**, **discretization**, and more. Each serves a specific purpose, such as improving interpretability, aligning data with algorithm requirements, or reducing noise.

Furthermore, after cleaning our data set, we applied data transformation techniques like discretization and normalization. Discretization is especially important for transforming continuous numerical attributes into discrete categories or intervals, simplifying data, and making it more suitable for analysis.

Anxiety Level Discretization

Here, we discretized the `anxiety_level` feature into four distinct bins, categorizing individuals based on their anxiety levels. We defined the bins by the edges `[0, 4, 9, 14, 21]` and correspond to the labels `['0-4', '5-9', '10-14', '15-21']`. This transformation helped us break down the continuous anxiety scores into four specific ranges, making it easier to analyze and compare groups of individuals based on their anxiety levels.

By discretizing `anxiety_level`, we provided a clearer understanding of how anxiety is distributed across our dataset. This is particularly useful for our studies or analyses where the exact numerical value of anxiety is not as important for us as understanding the broad categories or ranges that help us identify trends and patterns within different levels of anxiety. The discretization of `anxiety_level` simplifies the data while maintaining important information about anxiety severity across individuals.

```

column_to_discretize = 'anxiety_level'
bin_edges = [0, 4, 9, 14, 21]
bin_labels = ['0-4', '5-9', '10-14', '15-21']

df['discretized_' + column_to_discretize] = pd.cut(df[column_to_discretize],
bins=bin_edges, labels=bin_labels, include_lowest=True)

print("Original DataFrame:")
print(df[['anxiety_level', 'discretized_anxiety_level']])

column_to_discretize = 'anxiety_level'
bin_edges = [0, 4, 9, 14, 21]
bin_labels = ['0-4', '5-9', '10-14', '15-21']
df_CleanedDataset['discretized_' + column_to_discretize] =
pd.cut(df_CleanedDataset[column_to_discretize], bins=bin_edges,
labels=bin_labels, include_lowest=True)

```

Original DataFrame:

| | anxiety_level | discretized_anxiety_level |
|------|---------------|---------------------------|
| 0 | 14 | 10-14 |
| 1 | 15 | 15-21 |
| 2 | 12 | 10-14 |
| 3 | 16 | 15-21 |
| 4 | 16 | 15-21 |
| ... | ... | ... |
| 1095 | 11 | 10-14 |
| 1096 | 9 | 5-9 |
| 1097 | 4 | 0-4 |
| 1098 | 21 | 15-21 |
| 1099 | 18 | 15-21 |

[1100 rows x 2 columns]

Self-Esteem Discretization

We discretized the self_esteem column into three categories based on the bin edges [0, 15, 25, 30] with corresponding labels ['0-15', '16-25', '26-30']. This transformation classifies self-esteem values into low, medium, and high categories.

By discretizing self-esteem in this manner, we segmented our dataset into more manageable groups, enabling a more straightforward analysis of how individuals with different levels of self-esteem might relate to other factors such as stress or mental health history.

This categorization makes it easier for us to analyze how individuals with varying self-esteem levels might experience other health and lifestyle conditions, providing us with more meaningful insights into the relationships between self-esteem and factors like anxiety, depression, or social support.

```

column_to_discretize = 'self_esteem'
bin_edges = [0, 15, 25, 30]
bin_labels = ['0-15', '16-25', '26-30']

df['discretized_' + column_to_discretize] = pd.cut(df[column_to_discretize],
bins=bin_edges, labels=bin_labels, include_lowest=True)

print("Original DataFrame:")
print(df[['self_esteem', 'discretized_self_esteem']])

column_to_discretize = 'self_esteem'
bin_edges = [0, 15, 25, 30]
bin_labels = ['0-15', '16-25', '26-30']
df_CleanedDataset['discretized_' + column_to_discretize] =
pd.cut(df_CleanedDataset[column_to_discretize], bins=bin_edges,
labels=bin_labels, include_lowest=True)

```

Original DataFrame:

| | self_esteem | discretized_self_esteem |
|------|-------------|-------------------------|
| 0 | 20 | 16-25 |
| 1 | 8 | 0-15 |
| 2 | 18 | 16-25 |
| 3 | 12 | 0-15 |
| 4 | 28 | 26-30 |
| ... | ... | ... |
| 1095 | 17 | 16-25 |
| 1096 | 12 | 0-15 |
| 1097 | 26 | 26-30 |
| 1098 | 0 | 0-15 |
| 1099 | 6 | 0-15 |

[1100 rows x 2 columns]

Depression Discretization

Here, we discretized depression into five distinct ranges using the bin edges [0, 4, 9, 14, 19, 27] and the labels ['0-4', '5-9', '10-14', '15-19', '20-27']. This categorization reflects varying degrees of depression severity, allowing for a more digestible representation of our dataset.

The discretization of depression provides us with an easy way to group individuals into categories based on their depression levels, which is valuable for identifying patterns in how different levels of depression interact with other factors in our dataset. For instance, those in the 15-19 range may be more likely to experience higher levels of stress or anxiety, which can be helpful for targeted interventions or studies focusing on the relationship between depression and other mental health issues.

The discretization approach allows for a clearer understanding of how depression severity varies across individuals, and gives us an easier framework for analyzing its impact on other health, lifestyle, and environmental factors.

```
column_to_discretize = 'depression'
bin_edges = [0, 4, 9, 14, 19, 27]
bin_labels = ['0-4', '5-9', '10-14', '15-19', '20-27']

df['discretized_' + column_to_discretize] = pd.cut(df[column_to_discretize],
bins=bin_edges, labels=bin_labels, include_lowest=True)

print("Original DataFrame:")
print(df[['depression', 'discretized_depression']])

column_to_discretize = 'depression'
bin_edges = [0, 4, 9, 14, 19, 27]
bin_labels = ['0-4', '5-9', '10-14', '15-19', '20-27']
df_CleanedDataset['discretized_' + column_to_discretize] =
pd.cut(df_CleanedDataset[column_to_discretize], bins=bin_edges,
labels=bin_labels, include_lowest=True)
```

Original DataFrame:

| | depression | discretized_depression |
|------|------------|------------------------|
| 0 | 11 | 10-14 |
| 1 | 15 | 15-19 |
| 2 | 14 | 10-14 |
| 3 | 15 | 15-19 |
| 4 | 7 | 5-9 |
| ... | ... | ... |
| 1095 | 14 | 10-14 |
| 1096 | 8 | 5-9 |
| 1097 | 3 | 0-4 |
| 1098 | 19 | 15-19 |
| 1099 | 15 | 15-19 |

[1100 rows x 2 columns]

Min-Max Normalization

Another transformation technique we used was normalization, an essential preprocessing technique that prepared our dataset for analysis. Normalization, specifically the **Min-Max scaling method**, rescaled numerical attributes to a fixed range [0, 1]. This ensured that features with larger scales did not dominate distance calculations in k-means clustering, allowing all attributes to contribute equally.

The output of our Min-Max scaling displays the first five rows of normalized data, where each column is scaled between 0 and 1, allowing for straightforward comparison across

metrics. In anxiety_level, self_esteem, and depression, values now indicate relative intensities, with higher values signifying greater levels within the dataset.

After min-max scaling, columns like sleep_quality, breathing_problem, and noise_level show normalized lifestyle and environmental conditions, providing comparable scores across various metrics. We kept stress_level in its original format, as it serves as an anchor for analyzing stress. Additionally, discretized columns like discretized_anxiety_level, discretized_self_esteem, and discretized_depression categorize the data into ranges, which simplifies interpretation and highlights general trends in anxiety, self-esteem, and depression levels across participants.

```
from sklearn.preprocessing import MinMaxScaler
```

```
# Step 1: Identify columns with numerical data to normalize, excluding  
'stress_level' and 'mental_health_history'  
columns_to_normalize = df.select_dtypes(include=[np.number]).columns.tolist()  
columns_to_normalize.remove('stress_level')  
if 'mental_health_history' in columns_to_normalize:  
    columns_to_normalize.remove('mental_health_history')
```

```
# Step 2: Instantiate the MinMaxScaler, which scales data to a range of 0-1.  
minmax_scaler = MinMaxScaler()
```

```
# Step 3: Apply Min-Max scaling to the selected columns  
normalized_data_minmax =  
minmax_scaler.fit_transform(df[columns_to_normalize])
```

```
# Step 4: Update the original dataframe with the normalized values for the  
selected columns  
df[columns_to_normalize] = normalized_data_minmax
```

```
# Display the normalized dataset's first few rows  
print("\nMin-Max scaled data:")  
print(df.head())
```

```
# Repeat normalization for a cleaned version of the dataset  
(df_CleanedDataset)  
columns_to_normalize =  
df_CleanedDataset.select_dtypes(include=[np.number]).columns.tolist()  
columns_to_normalize.remove('stress_level')  
minmax_scaler = MinMaxScaler()  
df_CleanedDataset[columns_to_normalize] =  
minmax_scaler.fit_transform(df_CleanedDataset[columns_to_normalize])
```

Min-Max scaled data:

| | anxiety_level | self_esteem | mental_health_history | depression | headache | \ |
|---|---------------|-------------|-----------------------|------------|----------|---|
| 0 | 0.666667 | 0.666667 | 0 | 0.407407 | 0.4 | |
| 1 | 0.714286 | 0.266667 | 1 | 0.555556 | 1.0 | |

| | | | | | |
|---|----------|----------|---|----------|-----|
| 2 | 0.571429 | 0.600000 | 1 | 0.518519 | 0.4 |
| 3 | 0.761905 | 0.400000 | 1 | 0.555556 | 0.8 |
| 4 | 0.761905 | 0.933333 | 0 | 0.259259 | 0.4 |

| | | | | | |
|---|----------------|---------------|-------------------|-------------|---|
| | blood_pressure | sleep_quality | breathing_problem | noise_level | \ |
| 0 | 0.0 | 0.4 | 0.8 | 0.4 | |
| 1 | 1.0 | 0.2 | 0.8 | 0.6 | |
| 2 | 0.0 | 0.4 | 0.4 | 0.4 | |
| 3 | 1.0 | 0.2 | 0.6 | 0.8 | |
| 4 | 1.0 | 1.0 | 0.2 | 0.6 | |

| | | | | |
|---|-------------------|-----|------------------------------|---|
| | living_conditions | ... | teacher_student_relationship | \ |
| 0 | 0.6 | ... | 0.6 | |
| 1 | 0.2 | ... | 0.2 | |
| 2 | 0.4 | ... | 0.6 | |
| 3 | 0.4 | ... | 0.2 | |
| 4 | 0.4 | ... | 0.2 | |

| | | | | |
|---|------------------------|----------------|---------------|---|
| | future_career_concerns | social_support | peer_pressure | \ |
| 0 | 0.6 | 0.666667 | 0.6 | |
| 1 | 1.0 | 0.333333 | 0.8 | |
| 2 | 0.4 | 0.666667 | 0.6 | |
| 3 | 0.8 | 0.333333 | 0.8 | |
| 4 | 0.4 | 0.333333 | 1.0 | |

| | | | | |
|---|----------------------------|----------|--------------|---|
| | extracurricular_activities | bullying | stress_level | \ |
| 0 | 0.6 | 0.4 | 1 | |
| 1 | 1.0 | 1.0 | 2 | |
| 2 | 0.4 | 0.4 | 1 | |
| 3 | 0.8 | 1.0 | 2 | |
| 4 | 0.0 | 1.0 | 1 | |

| | | | |
|---|---------------------------|-------------------------|------------------------|
| | discretized_anxiety_level | discretized_self_esteem | discretized_depression |
| 0 | 10-14 | 16-25 | 10-14 |
| 1 | 15-21 | 0-15 | 15-19 |
| 2 | 10-14 | 16-25 | 10-14 |
| 3 | 15-21 | 0-15 | 15-19 |
| 4 | 15-21 | 26-30 | 5-9 |

[5 rows x 24 columns]

Data Reduction: Feature Selection

Feature selection is another preprocessing step that involves identifying and retaining the most relevant attributes in a dataset while removing redundant or irrelevant ones. This reduces the complexity of the data, enhances model performance, and prevents overfitting by focusing the analysis on meaningful features.

In our project, we applied feature selection to streamline the dataset, ensuring that only the most impactful attributes were used for classification and clustering. We analyzed attribute correlations and distributions to identify features with low variance or redundancy. Removing these attributes not only simplified the dataset but also improved the computational efficiency of our models.

We employed the filter method for feature selection by examining the correlation coefficients between the target variable, stress_level, and various features. The resulting features with strong positive correlations include self_esteem (0.756), bullying (0.751), sleep_quality (0.749), future_career_concerns (0.743), anxiety_level (0.737), depression (0.734), academic_performance (0.721), headache (0.713), safety (0.710), and basic_needs (0.709). These features were selected due to their correlation coefficients exceeding the threshold of 0.7, indicating significant relationships with stress levels.

```
# Step 1: Drop rows with any missing values from the DataFrame
df_cleaned = df.dropna()

# Step 2: Select only the numeric columns from the cleaned DataFrame
df_numeric = df_cleaned.select_dtypes(include=['number'])

# Step 3: Check if the 'stress_level' column exists in the numeric DataFrame
if 'stress_level' not in df_numeric.columns:
    print("'stress_level' column is not numeric or missing.")
else:
    # Step 4: Calculate the absolute correlation of all numeric features with 'stress_level'
    correlation_with_target =
df_numeric.corr()['stress_level'].abs().sort_values(ascending=False)

    print("\nCorrelation with Target Variable (stress_level):")
    print(correlation_with_target) # Display the correlation values

    # Step 5: Set a threshold for selecting features based on correlation
    threshold = 0.7
    # Step 6: Select features that have a correlation greater than the threshold
    selected_features_corr = correlation_with_target[correlation_with_target
> threshold].index.tolist()

    # Step 7: Remove 'stress_level' from the list of selected features if present
    if 'stress_level' in selected_features_corr:
        selected_features_corr.remove('stress_level')

    print("\nSelected Features based on correlation:")
    print(selected_features_corr) # Display selected features
```



```

# Step 8: Create a final DataFrame with only the selected features
X_final = df_numeric[selected_features_corr]

print("\nFinal DataFrame with Selected Features:")
print(X_final.head()) # Display the first few rows of the final
DataFrame

# Step 9: Create a cleaned dataset that includes only the selected
features, dropping any rows with missing values
df_CleanedDataset = df_cleaned[selected_features_corr].dropna()

# Step 10: Output the shape of the cleaned dataset
print("\nShape of Cleaned Dataset:")
print(df_CleanedDataset.shape) # Display the shape of the cleaned
dataset

# Visualization: Feature Selection
# Separate selected and dropped features
dropped_features = correlation_with_target[correlation_with_target <=
threshold]
selected_features = correlation_with_target[correlation_with_target >
threshold]

# Plotting
plt.figure(figsize=(10, 6))

# Plot selected features
plt.barh(
    selected_features.index,
    selected_features.values,
    color="skyblue",
    label="Selected Features",
)

# Plot dropped features
plt.barh(
    dropped_features.index,
    dropped_features.values,
    color="salmon",
    label="Dropped Features",
)

# Adding labels and title
plt.axvline(x=threshold, color="black", linestyle="--",
label="Threshold")
plt.xlabel("Correlation with Target (stress_level)", fontsize=12)
plt.ylabel("Features", fontsize=12)
plt.title("Feature Selection Based on Correlation", fontsize=14)

```

```
plt.legend()
plt.tight_layout()

# Show the plot
plt.show()
```

Correlation with Target Variable (stress_level):

| | |
|------------------------------|----------|
| stress_level | 1.000000 |
| self_esteem | 0.756195 |
| bullying | 0.751162 |
| sleep_quality | 0.749068 |
| future_career_concerns | 0.742619 |
| anxiety_level | 0.736795 |
| depression | 0.734379 |
| academic_performance | 0.720922 |
| headache | 0.713484 |
| safety | 0.709602 |
| basic_needs | 0.708968 |
| extracurricular_activities | 0.692977 |
| peer_pressure | 0.690684 |
| teacher_student_relationship | 0.680163 |
| noise_level | 0.663371 |
| mental_health_history | 0.648644 |
| study_load | 0.634156 |
| social_support | 0.632497 |
| living_conditions | 0.581723 |
| breathing_problem | 0.573984 |
| blood_pressure | 0.394200 |

Name: stress_level, dtype: float64

Selected Features based on correlation:

```
['self_esteem', 'bullying', 'sleep_quality', 'future_career_concerns',
'anxiety_level', 'depression', 'academic_performance', 'headache', 'safety',
'basic_needs']
```

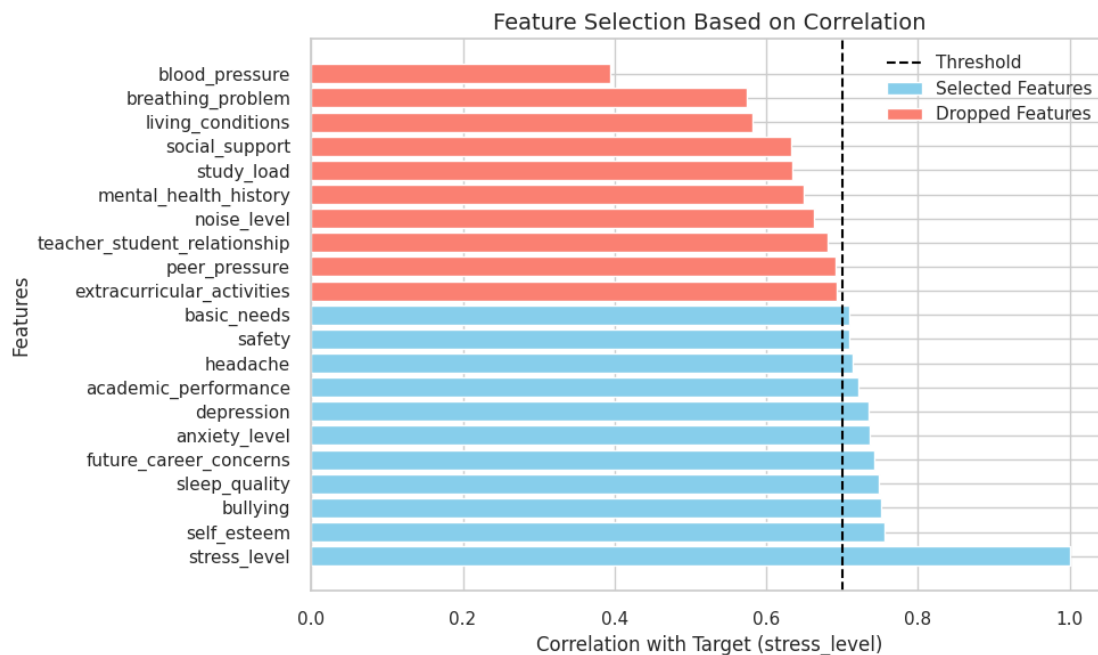
Final DataFrame with Selected Features:

| | self_esteem | bullying | sleep_quality | future_career_concerns | \ |
|---|-------------|----------|---------------|------------------------|---|
| 0 | 0.666667 | 0.4 | 0.4 | 0.6 | |
| 1 | 0.266667 | 1.0 | 0.2 | 1.0 | |
| 2 | 0.600000 | 0.4 | 0.4 | 0.4 | |
| 3 | 0.400000 | 1.0 | 0.2 | 0.8 | |
| 4 | 0.933333 | 1.0 | 1.0 | 0.4 | |

| | anxiety_level | depression | academic_performance | headache | safety | \ |
|---|---------------|------------|----------------------|----------|--------|---|
| 0 | 0.666667 | 0.407407 | 0.6 | 0.4 | 0.6 | |
| 1 | 0.714286 | 0.555556 | 0.2 | 1.0 | 0.4 | |
| 2 | 0.571429 | 0.518519 | 0.4 | 0.4 | 0.6 | |
| 3 | 0.761905 | 0.555556 | 0.4 | 0.8 | 0.4 | |
| 4 | 0.761905 | 0.259259 | 0.8 | 0.4 | 0.8 | |

| | basic_needs |
|---|-------------|
| 0 | 0.4 |
| 1 | 0.4 |
| 2 | 0.4 |
| 3 | 0.4 |
| 4 | 0.6 |

Shape of Cleaned Dataset:
(1100, 10)



The next code does not do anything except adding the class label (*stress_level*) to the cleaned data set

```
from google.colab import files

df_cleaned = df.dropna()

df_numeric = df_cleaned.select_dtypes(include=['number'])

if 'stress_level' not in df_numeric.columns:
    print("'stress_level' column is not numeric or missing.")
else:
    correlation_with_target =
df_numeric.corr()['stress_level'].abs().sort_values(ascending=False)

    threshold = 0.7
    selected_features_corr = correlation_with_target[correlation_with_target
> threshold].index.tolist()
```

```
if 'stress_level' in selected_features_corr:
    selected_features_corr.remove('stress_level')

X_final = df_numeric[selected_features_corr]

df_CleanedDataset = df_cleaned[selected_features_corr].copy()
df_CleanedDataset['stress_level'] = df_cleaned['stress_level']

df_CleanedDataset.to_csv('Processed_Dataset.csv', index=False)
files.download('Processed_Dataset.csv')
```

<IPython.core.display.Javascript object>

<IPython.core.display.Javascript object>

5. Data Mining Techniques

In this project, we explored the crucial role of data mining in understanding and forecasting satisfaction probabilities, with a particular focus on the factors contributing to student stress. By leveraging advanced analytical techniques, we aimed to extract meaningful insights from our dataset, ultimately enhancing our understanding of student well-being.

Data Mining Approach Data mining is an essential process that involves extracting patterns and knowledge from large datasets. Our approach utilized a combination of classification and clustering to gain a deeper understanding of student stress factors. This methodology allowed us to uncover hidden patterns, relationships, and trends that may not be immediately apparent, thus providing a comprehensive view of the dataset.

Significance of Data Mining Data mining serves as a bridge between raw data and actionable insights. In the context of our project, it enabled us to analyze complex datasets to identify key stressors affecting students. By applying various analytical techniques, we were able to:

- **Discover Patterns:** Data mining helped us reveal patterns in student behavior and stress levels, offering insights into how different variables interact.
 - **Identify Relationships:** Through our analysis, we identified relationships between various stress factors, which can inform targeted interventions and
 - **Support strategies:** Enhance Predictions: By forecasting satisfaction probabilities, we aimed to provide educators and administrators with tools to better address student needs and improve overall satisfaction.
 - **Facilitate Decision-Making:** The insights gained through data mining empower stakeholders to make informed decisions based on empirical evidence, ultimately leading to improved student outcomes.
-

Classification and Clustering Techniques While our focus on specific classification and clustering techniques was part of our methodology, these elements were integrated into the broader data mining process. We used decision trees for classification and k-means for clustering, allowing us to categorize students and identify distinct groups within the dataset. This dual approach enriched our analysis and provided a multifaceted understanding of student stress.

Tools and Libraries To implement our data mining strategy, we utilized a variety of tools and libraries, including:

1. **Classification:** We utilized the `DecisionTreeClassifier` from the `sklearn` library, with model evaluation supported by functions such as `accuracy_score` and `classification_report`. Visualization was enhanced through the `plot_tree` function, allowing for a clearer understanding of the decision-making process.
2. **Clustering:** The clustering analysis was conducted using `KMeans` from `sklearn.cluster`, with performance evaluated through Yellowbrick's visualization tools to ensure clarity and interpretability.
3. **Evaluation Metrics:** We relied on critical performance metrics including `confusion_matrix`, `precision_score`, and `recall_score` to assess our models thoroughly.
4. **Visualization:** Libraries such as `matplotlib` and Yellowbrick were instrumental in interpreting and presenting our results effectively, offering insights that were both meaningful and visually appealing.

6. Evaluation and Comparison

6.1 Classification

In our classification analysis, we utilized three data partitions (70%-30%, 80%-20%, and 90%-10%) and evaluated the model performance using two criteria: Gini and Entropy. The Gini criterion measures impurity, while the Entropy criterion focuses on information gain, helping to determine how well the model differentiates between stress levels.

We employed confusion matrices to visualize predictions against actual stress levels, allowing us to assess the model's performance in a clear manner. For both criteria, we calculated key metrics including test accuracy, sensitivity, specificity, and precision. These metrics provide insights into the model's ability to correctly classify instances across different stress levels.

Additionally, we created decision tree visualizations to illustrate the model structures and decision-making processes based on both Gini and Entropy.

Class imbalance

The distribution of our dataset appears to be balanced across all classes, with only slight variations observed. This indicates that there is no significant imbalance, and the dataset is fairly well-balanced overall. Therefore, techniques such as SMOTE or other methods are unnecessary in this case. Our dataset should naturally support the classification model.

```
print(df_CleanedDataset['stress_level'].value_counts())
```

```
stress_level
0      373
2      369
1      358
Name: count, dtype: int64
```

Prepare the Data, feature Set and Target Variable Definition After Feature Selection

We separated the features and the target variable(stress_level).

We define the feature set and target variable after performing feature selection which was in the previous phase, where unnecessary columns were removed. The feature set consists of the relevant predictors(X), while the target variable indicates the stress level(Y).

```
# Define feature set and target variable
X = df_CleanedDataset.drop('stress_level', axis=1) # Features (all columns except 'stress_level')
y = df_CleanedDataset['stress_level'] # Target variable (stress_level)
```

```
from sklearn.model_selection import train_test_split
```

```
X_train_70, X_test_30, y_train_70, y_test_30 = train_test_split(X, y,
test_size=0.3, random_state=42)
print("Splitting Data into:")
print("> Training (70%)")
print("> Test (30%)")
```

```
# 2. Split the dataset for 80%-20%
X_train_80, X_test_20, y_train_80, y_test_20 = train_test_split(X, y,
test_size=0.2, random_state=42)
print("\nSplitting Data into:")
print("> Training (80%)")
print("> Test (20%)")
```

```
# 3. Split the dataset for 90%-10%
X_train_90, X_test_10, y_train_90, y_test_10 = train_test_split(X, y,
test_size=0.1, random_state=42)
print("\nSplitting Data into:")
print("> Training (90%)")
print("> Test (10%)")
```

```
Splitting Data into:
> Training (70%)
```

- Test (30%)

Splitting Data into:

- Training (80%)
- Test (20%)

Splitting Data into:

- Training (90%)
- Test (10%)

Choosing Partitions We chose the following partitions based on common practices in machine learning to assess model performance effectively

70%-30% Split: This provides a balanced approach, allowing sufficient training data while retaining enough test data to evaluate the model's generalization.

80%-20% Split: Increasing the training set enhances the model's learning capacity, which is useful for more complex datasets. The 20% test set remains adequate for validation.

90%-10% Split: This maximizes training data, beneficial for models that require extensive data to learn patterns. The smaller test set still provides a reasonable measure of performance.

6.1.1 Information Gain

Information gain measures the effectiveness of an attribute in classifying data by quantifying the reduction in uncertainty or entropy when that attribute is used to split the data. It is commonly used in decision tree algorithms to select the best attribute for splitting nodes. By calculating the difference between the entropy of the dataset before and after the split, information gain identifies how much "information" the attribute adds. Higher information gain indicates a more effective split that leads to purer subsets, which are closer to having data points of a single class.

Performance Evaluation of the Entropy In this section, we implemented an Entropy-based Decision Tree classifier to evaluate its performance on our dataset. We initialized a dictionary to store various metrics for different training/testing partition sizes: Test Accuracy, Sensitivity, Specificity, and Precision.

For each partition (70%-30%, 80%-20%, 90%-10%) (0, 1 AND 2), we split the dataset, trained the model, and made predictions on the testing set to compute these metrics.

First Split

```
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import confusion_matrix, accuracy_score, recall_score,
precision_score
import pandas as pd
import numpy as np
```

```

# Model for the 70%-30% partition
model_70_30 = DecisionTreeClassifier(criterion='entropy', random_state=42)

# Train and evaluate for 70%-30% partition
partition = "70%-30%"
X_train, X_test, y_train, y_test = X_train_70, X_test_30, y_train_70,
y_test_30

# Train the model
model_70_30.fit(X_train, y_train)

# Predictions and metrics
y_pred = model_70_30.predict(X_test)
cm = confusion_matrix(y_test, y_pred)

# Calculate specificity (multi-class)
tn_sum = 0
fp_sum = 0
for i in range(len(cm)):
    tn = cm.sum() - (cm[i, :].sum() + cm[:, i].sum() - cm[i, i]) # True
negatives for class i
    fp = cm[:, i].sum() - cm[i, i] # False positives for class i
    tn_sum += tn
    fp_sum += fp

specificity = tn_sum / (tn_sum + fp_sum) if (tn_sum + fp_sum) > 0 else 0

# Store the results for the 70%-30% partition
results_70_30 = {
    'Partition': partition,
    'Train Accuracy': accuracy_score(y_train, model_70_30.predict(X_train)),
    'Test Accuracy': accuracy_score(y_test, y_pred),
    'Sensitivity': recall_score(y_test, y_pred, average='macro'),
    'Specificity': specificity,
    'Precision': precision_score(y_test, y_pred, average='macro')
}

# Convert results to DataFrame and print the result for 70%-30% partition
entropy_results_df_70_30 = pd.DataFrame([results_70_30])
print(f"Entropy Results for {partition}:")
print(entropy_results_df_70_30[['Partition', 'Test Accuracy', 'Sensitivity',
'Specificity', 'Precision']])

```

Entropy Results for 70%-30%:

| | Partition | Test Accuracy | Sensitivity | Specificity | Precision |
|---|-----------|---------------|-------------|-------------|-----------|
| 0 | 70%-30% | 0.887879 | 0.888071 | 0.943939 | 0.88882 |

Test Accuracy: represents the overall percentage of correct predictions made by the model on the test set. Here, the model correctly predicts 88.7% of the cases.

Sensitivity: measures the model's ability to correctly identify positive cases. A sensitivity of 88.81% indicates the model identifies 88.8% of actual positives correctly.

Specificity: measures the model's ability to correctly identify negative cases (true negatives). A specificity of 94.3% means the model accurately classifies 94.3% of actual negatives.

Precision: indicates the percentage of true positive predictions among all positive predictions made by the model. Here, 88.8% of the predicted positives are actual positives.

These metrics collectively provide a balanced view of the model's performance in handling both positive and negative cases in the 70%-30% partition.

```
from sklearn.metrics import ConfusionMatrixDisplay

models = {
    "70%-30%": DecisionTreeClassifier(criterion='entropy', random_state=42),
    "80%-20%": DecisionTreeClassifier(criterion='entropy', random_state=42),
    "90%-10%": DecisionTreeClassifier(criterion='entropy', random_state=42)
}

from sklearn.tree import DecisionTreeClassifier
import matplotlib.pyplot as plt

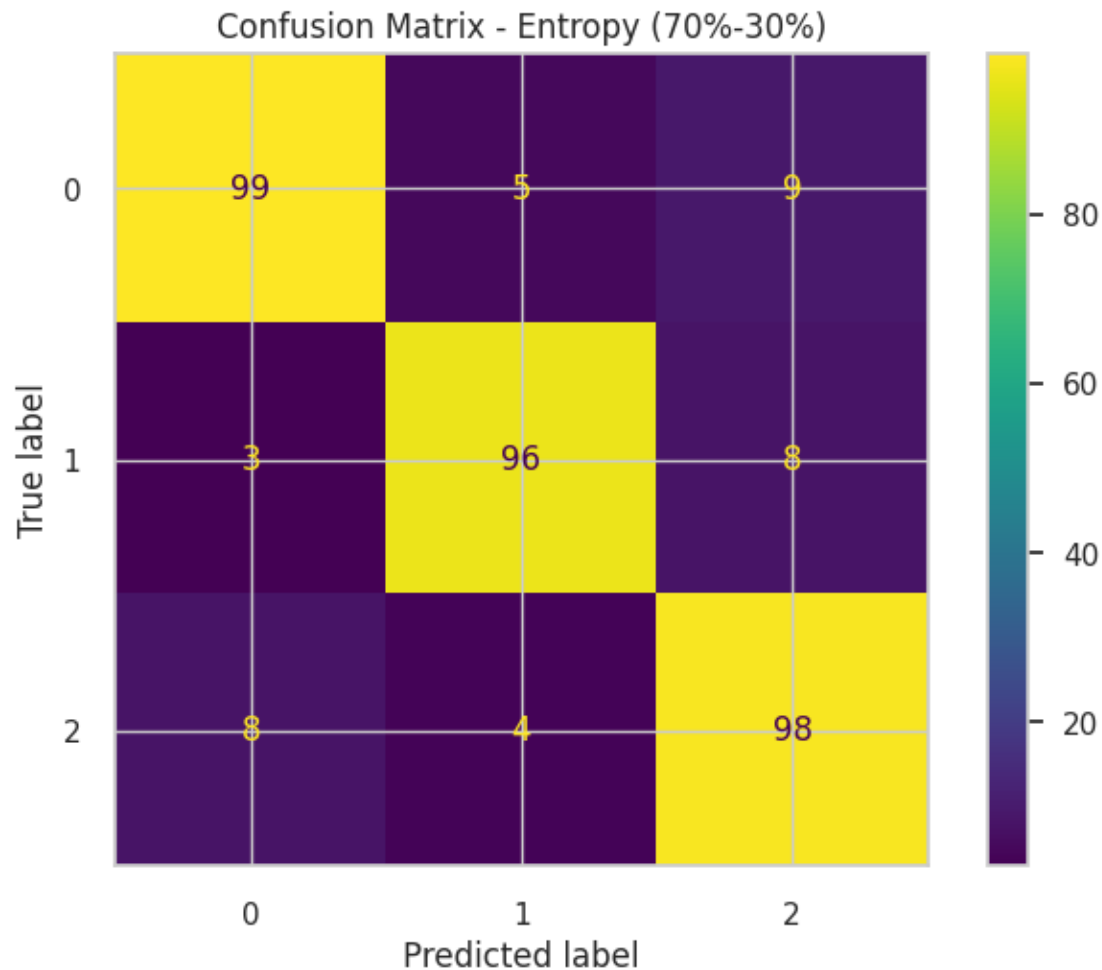
models["70%-30%"].fit(X_train_70, y_train_70)

unique_classes = sorted(set(y_train_70))

model = models["70%-30%"]
X_test, y_test, title_suffix = X_test_30, y_test_30, "70%-30%"

print(f"\nConfusion Matrix for Entropy with partition {title_suffix}:")
disp = ConfusionMatrixDisplay.from_estimator(model, X_test, y_test,
display_labels=unique_classes)
plt.title(f"Confusion Matrix - Entropy ({title_suffix})")
plt.show()
```

Confusion Matrix for Entropy with partition 70%-30%:



- **True Positives (TP):** The model accurately classified most instances of each class, as evident from the high diagonal values (99, 96, 98).
- **False Positives (FP):** The FP count was low, indicating the model made relatively few incorrect positive predictions.
- **True Negatives (TN):** A high TN count across all classes highlights the model's ability to correctly classify negatives.
- **False Negatives (FN):** The FN count was small but noticeable, with a few cases missed in each class.

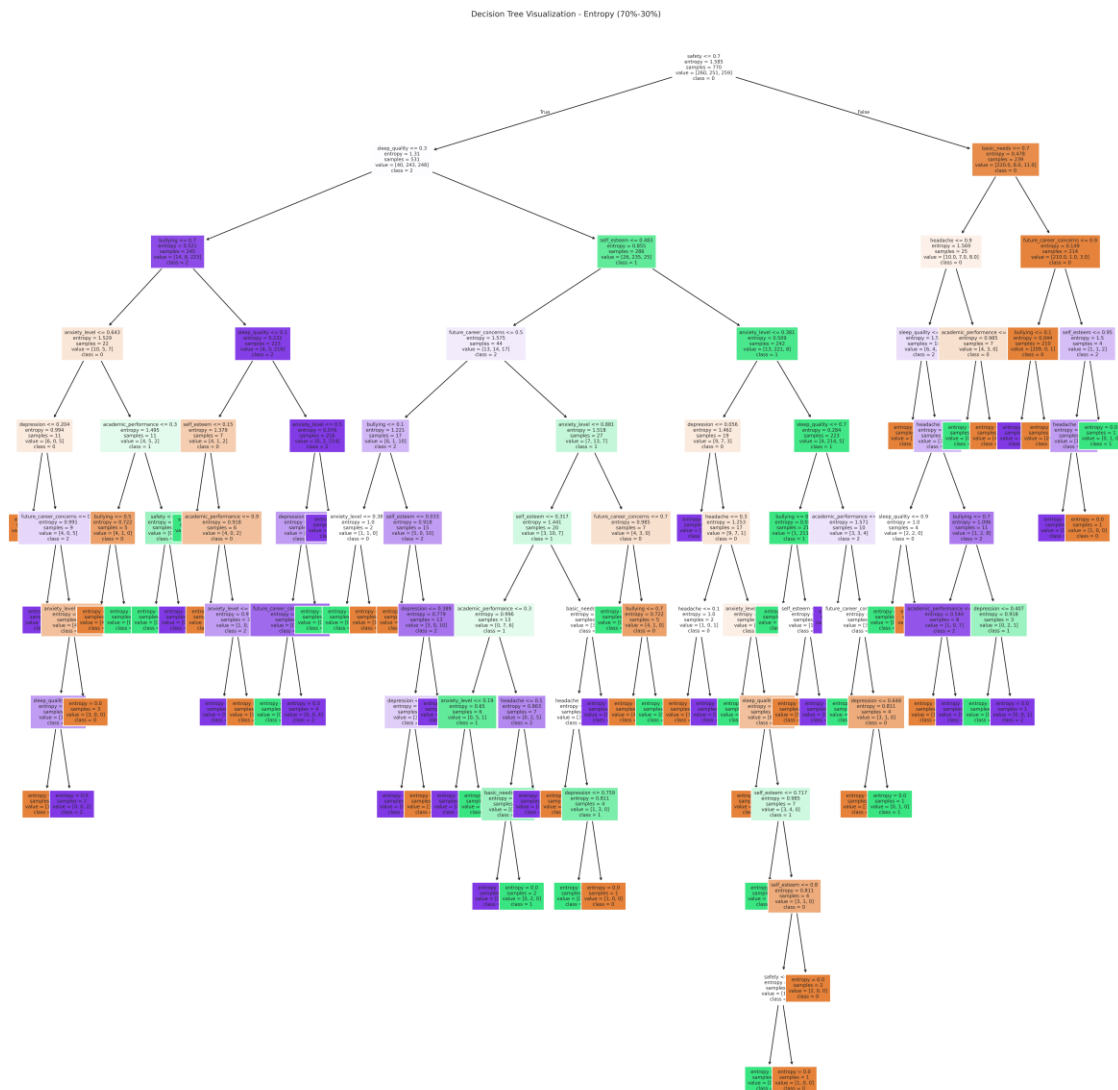
Observation: The model displayed strong overall performance, with well-balanced classification and minor room for improvement in reducing FNs.

```
from sklearn.tree import plot_tree
import matplotlib.pyplot as plt
import numpy as np
```

```
unique_classes = np.unique(y_train_70)
```

```
plt.figure(figsize=(20, 20), dpi=500)
plot_tree(models["70%-30%"], feature_names=X.columns,
class_names=unique_classes.astype(str), filled=True, fontsize=6)
plt.title("Decision Tree Visualization - Entropy (70%-30%)", fontsize=10)
plt.tight_layout()
```

```
filename = "decision_tree_entropy_70_30.png"
plt.savefig(filename, bbox_inches='tight', dpi=300)
plt.show()
plt.close()
```



- **Tree Complexity:** The tree generated with a 70%-30% split exhibited moderate depth, effectively **balancing accuracy and model simplicity**. The depth allowed the model to classify most instances while maintaining reasonable generalization.

- **Feature Importance:** Key features were prominently used in higher nodes, reflecting their influence on classification. Thresholds at each split logically separated data points, highlighting the dataset's inherent structure.
- **Interpretability:** The tree's structure was straightforward and interpretable, enabling easy visualization of the decision-making process. Leaf nodes provided clear distributions of classes.
- **Potential Issues:** *Minor overfitting* was possible in the deeper parts of the tree due to the moderate training set size.

Second Split

Model for the 80%-20% partition

```
model_80_20 = DecisionTreeClassifier(criterion='entropy', random_state=42)
```

Train and evaluate for 80%-20% partition

```
partition = "80%-20%"
```

```
X_train, X_test, y_train, y_test = X_train_80, X_test_20, y_train_80, y_test_20
```

Train the model

```
model_80_20.fit(X_train, y_train)
```

Predictions and metrics

```
y_pred = model_80_20.predict(X_test)
```

```
cm = confusion_matrix(y_test, y_pred)
```

Calculate specificity (multi-class)

```
tn_sum = 0
```

```
fp_sum = 0
```

```
for i in range(len(cm)):
```

```
    tn = cm.sum() - (cm[i, :].sum() + cm[:, i].sum() - cm[i, i]) # True negatives for class i
```

```
    fp = cm[:, i].sum() - cm[i, i] # False positives for class i
```

```
    tn_sum += tn
```

```
    fp_sum += fp
```

```
specificity = tn_sum / (tn_sum + fp_sum) if (tn_sum + fp_sum) > 0 else 0
```

Store the results for the 80%-20% partition

```
results_80_20 = {
```

```
    'Partition': partition,
```

```
    'Train Accuracy': accuracy_score(y_train, model_80_20.predict(X_train)),
```

```
    'Test Accuracy': accuracy_score(y_test, y_pred),
```

```
    'Sensitivity': recall_score(y_test, y_pred, average='macro'),
```

```
    'Specificity': specificity,
```

```
    'Precision': precision_score(y_test, y_pred, average='macro')
```

```
}
```

```
# Convert results to DataFrame and print the result for 80%-20% partition
entropy_results_df_80_20 = pd.DataFrame([results_80_20])
print(f"Entropy Results for {partition}:")
print(entropy_results_df_80_20[['Partition', 'Test Accuracy', 'Sensitivity',
'Specificity', 'Precision']])
```

Entropy Results for 80%-20%:

| | Partition | Test Accuracy | Sensitivity | Specificity | Precision |
|---|-----------|---------------|-------------|-------------|-----------|
| 0 | 80%-20% | 0.877273 | 0.876868 | 0.938636 | 0.877086 |

Test Accuracy: shows that the model correctly predicted 87.7% of the cases in the test set, providing an overall measure of performance.

Sensitivity: indicates the model correctly identified 87.6% of actual positive cases.

Specificity: represents the true negative rate, meaning the model accurately classified 93.8% of actual negative cases.

Precision: measures the proportion of true positives among all predicted positives. Here, 87.7% of positive predictions are correct.

These metrics demonstrate that the model performs slightly differently under the 80%-20% partition, with slightly **lower accuracy and sensitivity compared to the 70%-30% partition**.

```
models["80%-20%"].fit(X_train_80, y_train_80)
```

```
unique_classes = sorted(set(y_train_80))
```

```
model = models["80%-20%"]
```

```
X_test, y_test, title_suffix = X_test_20, y_test_20, "80%-20%"
```

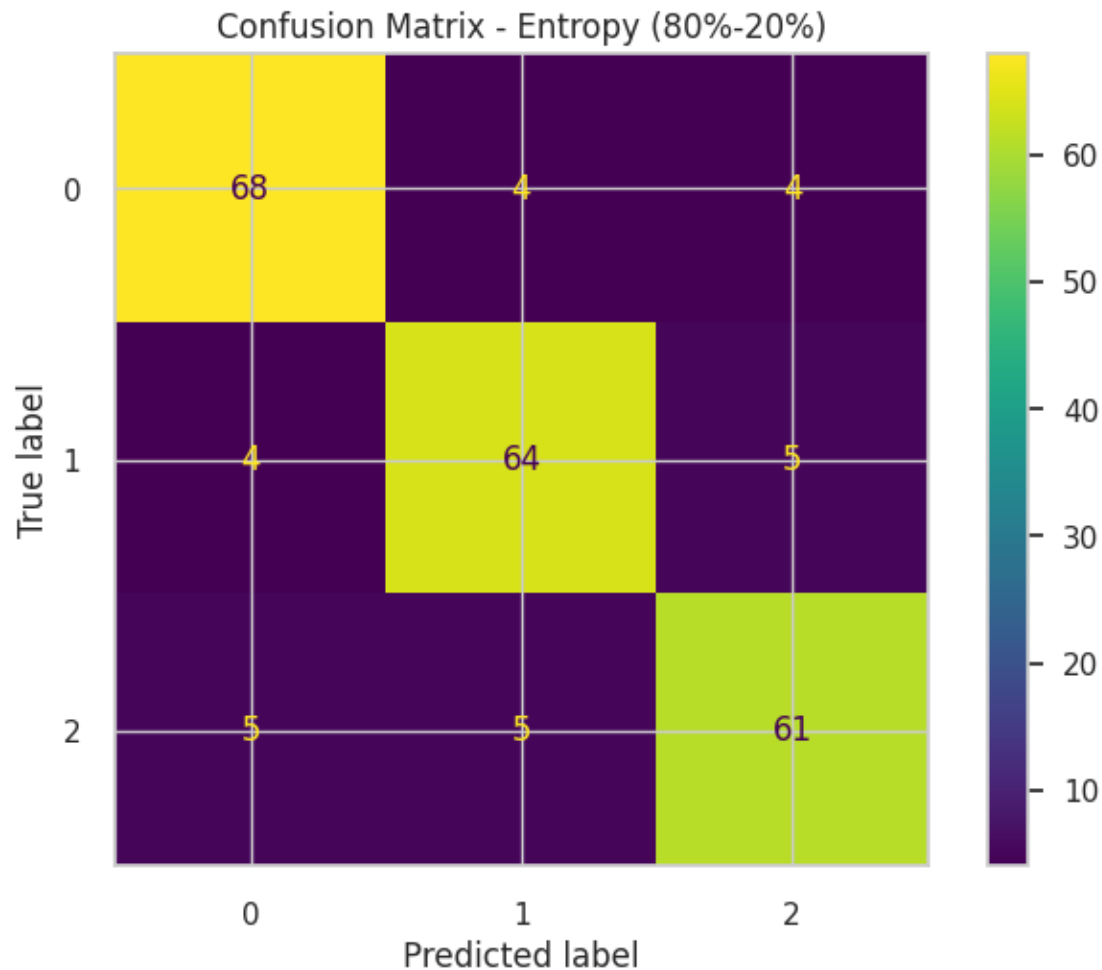
```
print(f"\nConfusion Matrix for Entropy with partition {title_suffix}:")
```

```
disp = ConfusionMatrixDisplay.from_estimator(model, X_test, y_test,
display_labels=unique_classes)
```

```
plt.title(f"Confusion Matrix - Entropy ({title_suffix})")
```

```
plt.show()
```

Confusion Matrix for Entropy with partition 80%-20%:



- **True Positives (TP):** The model correctly predicted the majority of instances for each class, shown by the diagonal values (68, 64, 61).
- **False Positives (FP):** FP counts were relatively low, supporting good precision across all classes.
- **True Negatives (TN):** The substantial TN count demonstrates the model's effective handling of negative classifications.
- **False Negatives (FN):** FN counts increased slightly compared to the 70%-30% partition, suggesting the model missed more positive cases with the reduced test set size.

Observation: The model maintained reliable classification but showed slightly higher misclassifications, likely due to the smaller dataset.

```
unique_classes = np.unique(y_train_80)
```

```
plt.figure(figsize=(20, 20), dpi=300)
plot_tree(models["80%-20%"], feature_names=X.columns,
```

```

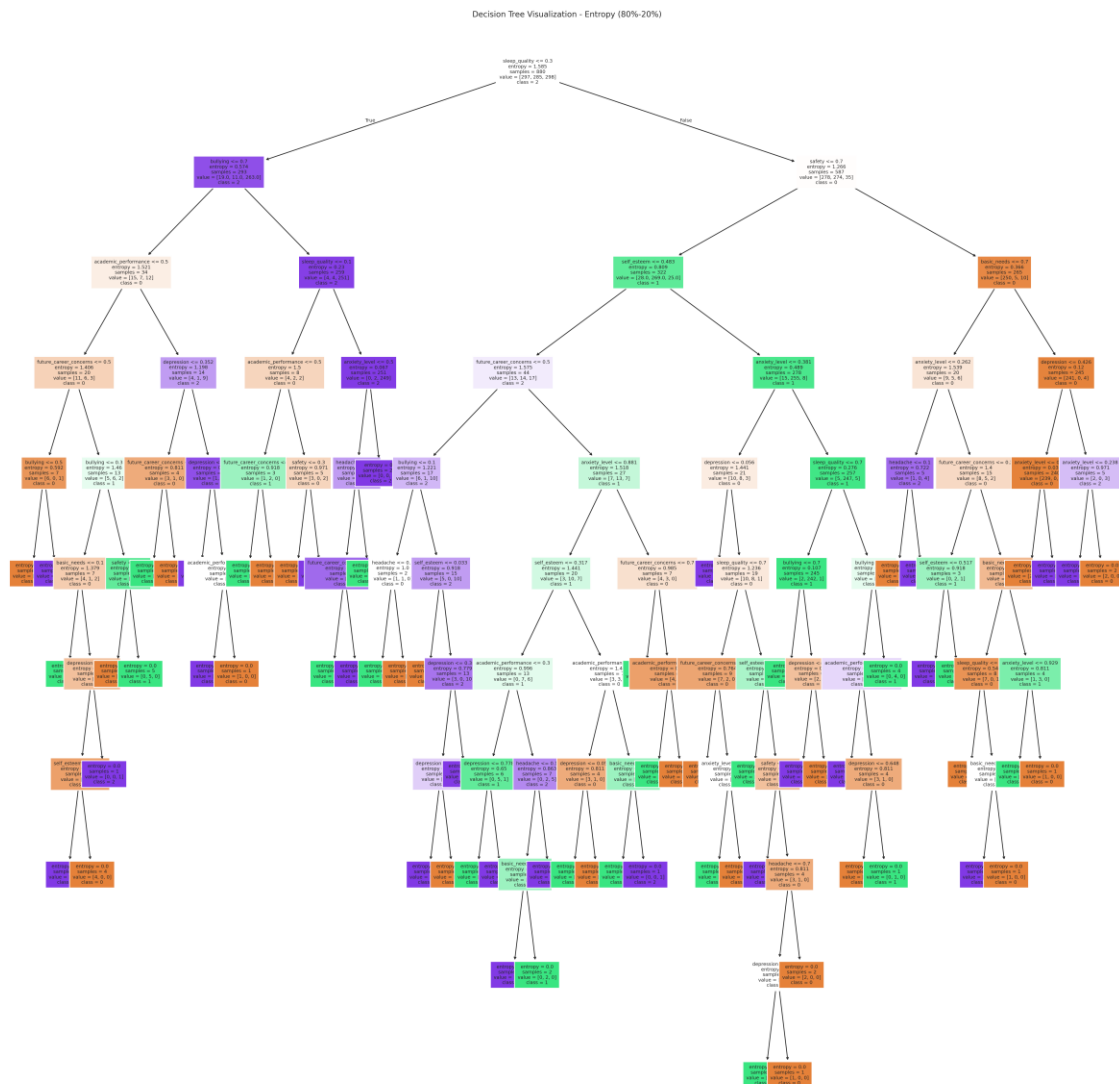
class_names=unique_classes.astype(str), filled=True, fontsize=6)
plt.title("Decision Tree Visualization - Entropy (80%-20%)", fontsize=10)
plt.tight_layout()

```

```

filename = "decision_tree_entropy_80_20.png"
plt.savefig(filename, bbox_inches='tight', dpi=300)
plt.show()
plt.close()

```



- **Tree Complexity:** The tree achieved a balanced depth, effectively utilizing the 80%-20% split. It managed to maintain simplicity while capturing meaningful patterns, striking a good balance between complexity and generalization.
- **Feature Importance:** Important features were prioritized at early splits, showcasing the **model's ability to identify impactful relationships**. Thresholds in the nodes provided actionable insights into the dataset's structure.

- **Interpretability:** The tree maintained a readable structure, making it ideal for explaining decision processes to a broader audience. The distribution of class probabilities at the leaf nodes was easy to interpret and informative.
- **Potential Issues:** While the tree captured essential patterns, minor risks of overfitting or underfitting remained.

Third Split

Model for the 90%-10% partition

```
model_90_10 = DecisionTreeClassifier(criterion='entropy', random_state=42)
```

Train and evaluate for 90%-10% partition

```
partition = "90%-10%"
```

```
X_train, X_test, y_train, y_test = X_train_90, X_test_10, y_train_90, y_test_10
```

Train the model

```
model_90_10.fit(X_train, y_train)
```

Predictions and metrics

```
y_pred = model_90_10.predict(X_test)
```

```
cm = confusion_matrix(y_test, y_pred)
```

Calculate specificity (multi-class)

```
tn_sum = 0
```

```
fp_sum = 0
```

```
for i in range(len(cm)):
```

```
    tn = cm.sum() - (cm[i, :].sum() + cm[:, i].sum() - cm[i, i]) # True negatives for class i
```

```
    fp = cm[:, i].sum() - cm[i, i] # False positives for class i
```

```
    tn_sum += tn
```

```
    fp_sum += fp
```

```
specificity = tn_sum / (tn_sum + fp_sum) if (tn_sum + fp_sum) > 0 else 0
```

Store the results for the 90%-10% partition

```
results_90_10 = {
```

```
    'Partition': partition,
```

```
    'Train Accuracy': accuracy_score(y_train, model_90_10.predict(X_train)),
```

```
    'Test Accuracy': accuracy_score(y_test, y_pred),
```

```
    'Sensitivity': recall_score(y_test, y_pred, average='macro'),
```

```
    'Specificity': specificity,
```

```
    'Precision': precision_score(y_test, y_pred, average='macro')
```

```
}
```

Convert results to DataFrame and print the result for 90%-10% partition

```
entropy_results_df_90_10 = pd.DataFrame([results_90_10])
```

```
print(f"Entropy Results for {partition}:")
```



```
print(entropy_results_df_90_10[['Partition', 'Test Accuracy', 'Sensitivity',
'Specificity', 'Precision']])
```

Entropy Results for 90%-10%:

| | Partition | Test Accuracy | Sensitivity | Specificity | Precision |
|---|-----------|---------------|-------------|-------------|-----------|
| 0 | 90%-10% | 0.854545 | 0.8559 | 0.927273 | 0.85151 |

Test Accuracy: The model correctly predicted 85.4% of the cases in the test set, reflecting its overall performance.

Sensitivity: shows the model successfully identified 85.5% of the actual positive cases.

Specificity: The model correctly classified 92.73% of the actual negative cases, representing its ability to avoid false positives.

Precision: indicates that 85.15% of the positive predictions made by the model are correct.

These metrics reveal that with a 90%-10% partition, **the model performs slightly lower across all metrics compared to the 70%-30% and 80%-20% partitions**, likely due to having less test data for evaluation.

```
models["90%-10%"].fit(X_train_90, y_train_90)
```

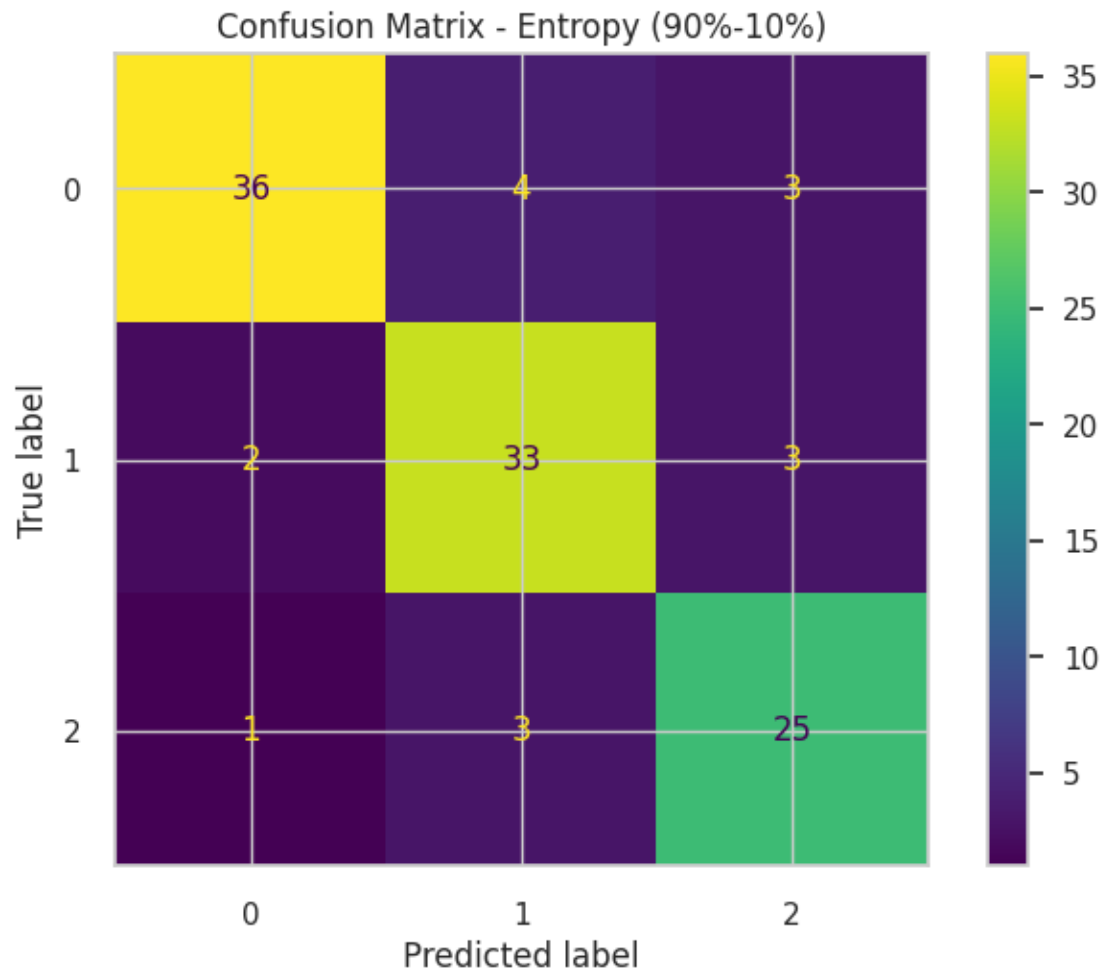
```
unique_classes = sorted(set(y_train_90))
```

```
model = models["90%-10%"]
```

```
X_test, y_test, title_suffix = X_test_10, y_test_10, "90%-10%"
```

```
print(f"\nConfusion Matrix for Entropy with partition {title_suffix}:")
disp = ConfusionMatrixDisplay.from_estimator(model, X_test, y_test,
display_labels=unique_classes)
plt.title(f"Confusion Matrix - Entropy ({title_suffix})")
plt.show()
```

Confusion Matrix for Entropy with partition 90%-10%:



- **True Positives (TP):** TP counts were lower due to the smaller test set, with values (36, 33, 25) for the respective classes.
- **False Positives (FP):** FP counts remained low, demonstrating continued precision in positive predictions.
- **True Negatives (TN):** A good TN count reflects the model's ability to handle negatives well, despite the reduced data.
- **False Negatives (FN):** FN counts increased slightly, especially for class 2, as the model struggled with limited data for generalization.

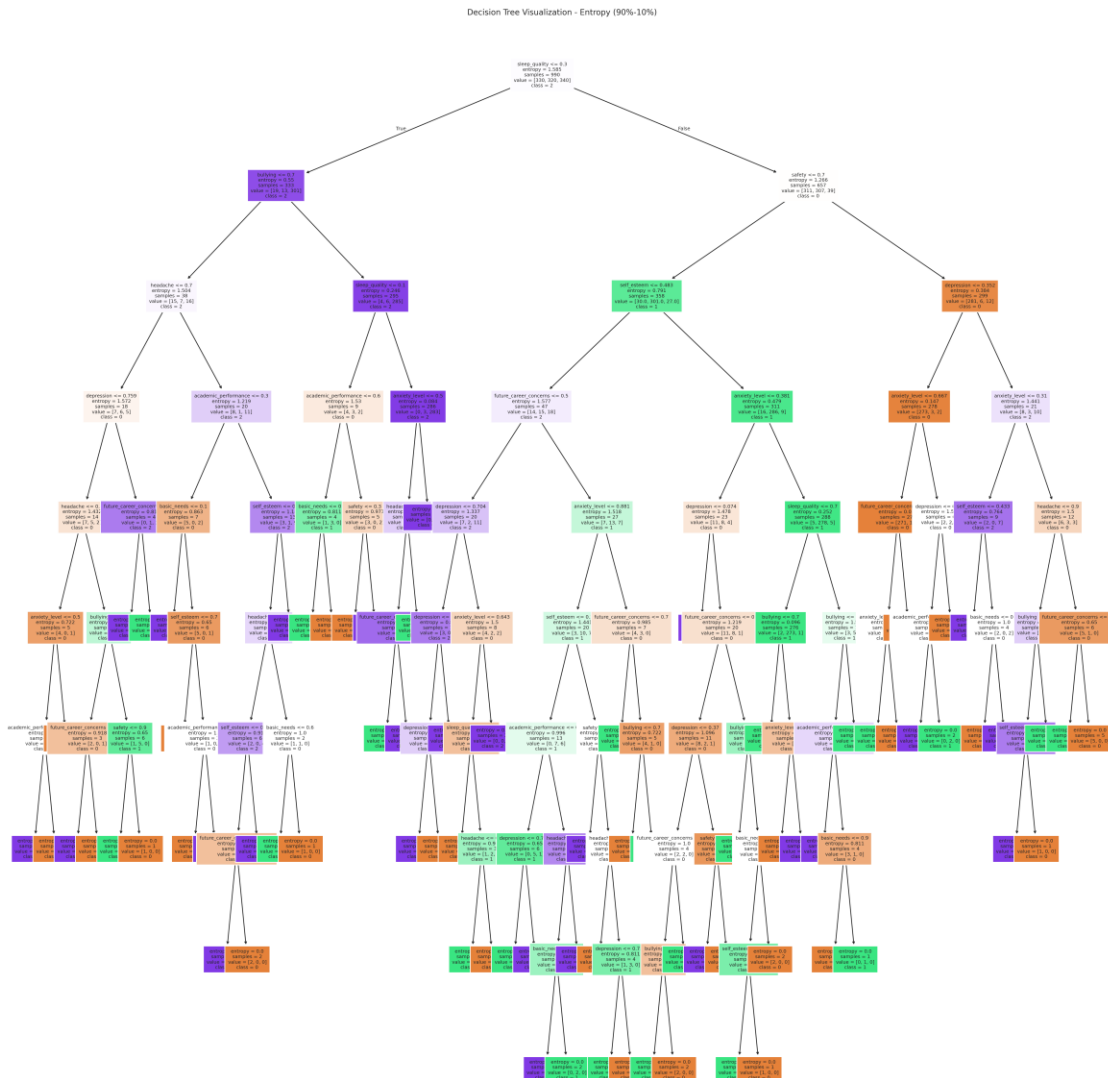
Observation: The smaller test set impacted the model's performance slightly, but it still performed reasonably well.

```
unique_classes = np.unique(y_train_90)
```

```
plt.figure(figsize=(20, 20), dpi=300)  
plot_tree(models["90%-10%"], feature_names=X.columns,  
class_names=unique_classes.astype(str), filled=True, fontsize=6)
```

```
plt.title("Decision Tree Visualization - Entropy (90%-10%)", fontsize=10)
plt.tight_layout()
```

```
filename = "decision_tree_entropy_90_10.png"
plt.savefig(filename, bbox_inches='tight', dpi=300)
plt.show()
plt.close()
```



- **Tree Complexity:** With a 90%-10% split, the tree was deeper and more detailed, capturing fine-grained patterns in the dataset. **While this improved training accuracy**, it may have reduced its ability to generalize due to potential overfitting.
- **Feature Importance:** The larger training set allowed the tree to focus on subtle feature interactions. High-impact features were effectively leveraged at the top levels, while lower-impact features appeared in deeper splits.

- **Interpretability:** The increased complexity slightly reduced interpretability. The deeper tree made it harder to trace decisions back to their source, though leaf nodes still provided useful class distribution insights.
- **Potential Issues:** Overfitting was a noticeable risk, as the model relied heavily on patterns specific to the training set. This could affect performance on unseen data.

Information Gain Analysis

1. Evaluation metrics

| | 90 %t raining set 10% testing set: | 80 %t raining set 20% testing set: | 70 %t raining set 30% testing set: |
|----------|---------------------------------------|---------------------------------------|---------------------------------------|
| Accuracy | 0.854545 | 0.877273 | 0.887879 |

The results show that *accuracy decreases as the test set size decreases*.

The **70%-30% split achieves the highest accuracy**, likely because it provides a better balance between training and test data, ensuring robust model evaluation. In contrast, **smaller test sets in the 80%-20% and 90%-10% splits may not fully represent unseen data**, potentially underestimating performance variability. Overall, the **70%-30% split provided the best balance of performance**.

| | 90 %t raining set 10% testing set: | 80 %t raining set 20% testing set: | 70 %t raining set 30% testing set: |
|-------------|---------------------------------------|---------------------------------------|---------------------------------------|
| Sensitivity | 0.8559 | 0.876868 | 0.888071 |
| Specificity | 0.927273 | 0.938636 | 0.943939 |
| Precision | 0.851510 | 0.877086 | 0.888820 |

The tabel shows that the **70%-30% split consistently outperforms the others** across Sensitivity 88.8%, Specificity 94.3%, and Precision 88.8%, indicating **it provides the most reliable balance between training and testing**. The 80%-20% split follows closely with slightly reduced performance, while the **90%-10% split has the lowest metrics** , likely due to insufficient test data for robust evaluation.

2. Confusion matrixes

The **70%-30% split achieves the best performance**, with the **highest true positive counts** and **fewer misclassifications**, due to a larger test set that allows for better evaluation of the model's generalization. The **80%-20% split shows slightly reduced performance**, with **moderate misclassifications** and **lower true positive counts**, reflecting the impact of having fewer test samples. The **90%-10% split performs the weakest**, with the smallest test set leading to **higher misclassification rates** and the **lowest true positive counts**, highlighting the trade-off between training size and evaluation robustness. Overall, the **70%-30% partition provides the most reliable**

results, while the 90%-10% partition demonstrates the challenges of limited test data.

3. Decision Tree

The decision trees based on different entropy-based splits (70%-30%, 80%-20%, and 90%-10%) vary in complexity and generalization. The **70%-30% split tree** is the **most generalizable with moderate complexity**, striking the **best balance between fitting the data and avoiding overfitting**. The **80%-20% split tree** offers a good compromise, with **better generalization than the 90%-10% model** and **more data for training than the 70%-30% model**, **balancing performance and complexity effectively**. The **90%-10% split tree**, using a larger training set, is **more complex and risks overfitting** due to its small test set, which **hampers its ability to generalize well**. Therefore, the **70%-30% split tree is generally the best for both performance and generalization**.

6.1.2 Gini-index

The **Gini index** is a crucial metric in building decision tree models, used to **determine the best splits** in the data. At each step, it **measures the impurity** of a dataset by calculating the likelihood of misclassifying a randomly chosen element. A lower Gini index indicates a purer node, where instances belong to the same class, while a higher index reflects greater impurity with mixed classes. By evaluating potential splits based on the Gini index, decision trees can effectively select features that reduce impurity, leading to more accurate models. This process continues until the tree reaches a specified depth or no significant gains in purity can be made.

Packages already installed in google colab

```
pip install pandas scikit-learn matplotlib seaborn
```

```
Requirement already satisfied: pandas in /usr/local/lib/python3.10/dist-packages (2.2.2)
```

```
Requirement already satisfied: scikit-learn in /usr/local/lib/python3.10/dist-packages (1.5.2)
```

```
Requirement already satisfied: matplotlib in /usr/local/lib/python3.10/dist-packages (3.8.0)
```

```
Requirement already satisfied: seaborn in /usr/local/lib/python3.10/dist-packages (0.13.2)
```

```
Requirement already satisfied: numpy>=1.22.4 in /usr/local/lib/python3.10/dist-packages (from pandas) (1.26.4)
```

```
Requirement already satisfied: python-dateutil>=2.8.2 in /usr/local/lib/python3.10/dist-packages (from pandas) (2.8.2)
```

```
Requirement already satisfied: pytz>=2020.1 in /usr/local/lib/python3.10/dist-packages (from pandas) (2024.2)
```

```
Requirement already satisfied: tzdata>=2022.7 in /usr/local/lib/python3.10/dist-packages (from pandas) (2024.2)
```

```
Requirement already satisfied: scipy>=1.6.0 in /usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.13.1)
```

```
Requirement already satisfied: joblib>=1.2.0 in
/usr/local/lib/python3.10/dist-packages (from scikit-learn) (1.4.2)
Requirement already satisfied: threadpoolctl>=3.1.0 in
/usr/local/lib/python3.10/dist-packages (from scikit-learn) (3.5.0)
Requirement already satisfied: contourpy>=1.0.1 in
/usr/local/lib/python3.10/dist-packages (from matplotlib) (1.3.1)
Requirement already satisfied: cyclor>=0.10 in
/usr/local/lib/python3.10/dist-packages (from matplotlib) (0.12.1)
Requirement already satisfied: fonttools>=4.22.0 in
/usr/local/lib/python3.10/dist-packages (from matplotlib) (4.55.0)
Requirement already satisfied: kiwisolver>=1.0.1 in
/usr/local/lib/python3.10/dist-packages (from matplotlib) (1.4.7)
Requirement already satisfied: packaging>=20.0 in
/usr/local/lib/python3.10/dist-packages (from matplotlib) (24.2)
Requirement already satisfied: pillow>=6.2.0 in
/usr/local/lib/python3.10/dist-packages (from matplotlib) (11.0.0)
Requirement already satisfied: pyparsing>=2.3.1 in
/usr/local/lib/python3.10/dist-packages (from matplotlib) (3.2.0)
Requirement already satisfied: six>=1.5 in /usr/local/lib/python3.10/dist-
packages (from python-dateutil>=2.8.2->pandas) (1.16.0)
```

First Split

In this analysis, we split our dataset into **70%** for **training** and **30%** for **testing**. This division allows us to effectively train our model while reserving a portion of the data to evaluate its performance on unseen data.

We utilize the `train_test_split` function to randomly partition the dataset accordingly. The training set, consisting of 70% of the data, is used to train a Decision Tree Classifier. This model is initialized with the **Gini** impurity criterion and a fixed random state for reproducibility.

The training process involves fitting the model to the training data, represented by `X_train` and `y_train`. Once the model is trained, we proceed to make predictions on the testing set, `X_test`. To assess the model's performance, we compute several metrics, including accuracy, sensitivity (recall), specificity, and precision.

```
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, recall_score,
precision_score

# Assuming X_train_70, X_test_30, y_train_70, y_test_30 are already defined

# Initialize results list and confusion matrices dictionary
gini_results = []
confusion_matrices = {}
```

```

# Define the model for the 70%-30% split
model = DecisionTreeClassifier(criterion='gini', random_state=42)

# Train and evaluate for the 70%-30% split
X_train, X_test, y_train, y_test = X_train_70, X_test_30, y_train_70,
y_test_30

# Train the model
model.fit(X_train, y_train)

# Predictions and metrics
y_pred = model.predict(X_test)
cm = confusion_matrix(y_test, y_pred)

# Store the confusion matrix
confusion_matrices["70%-30%"] = cm

# Calculate metrics
tn_sum = 0
fp_sum = 0
for i in range(len(cm)):
    tn = cm.sum() - (cm[i, :].sum() + cm[:, i].sum() - cm[i, i])
    fp = cm[:, i].sum() - cm[i, i]
    tn_sum += tn
    fp_sum += fp

specificity = tn_sum / (tn_sum + fp_sum) if (tn_sum + fp_sum) > 0 else 0

# Store results
results = {
    'Partition': "70%-30%",
    'Train Accuracy': accuracy_score(y_train, model.predict(X_train)),
    'Test Accuracy': accuracy_score(y_test, y_pred),
    'Sensitivity': recall_score(y_test, y_pred, average='macro'),
    'Specificity': specificity,
    'Precision': precision_score(y_test, y_pred, average='macro')
}

gini_results.append(results)

# Convert gini results to DataFrame and print
gini_results_df = pd.DataFrame(gini_results)
print("Gini Results:")
print(gini_results_df[['Partition', 'Test Accuracy', 'Sensitivity',
'Specificity', 'Precision']])

```

Gini Results:

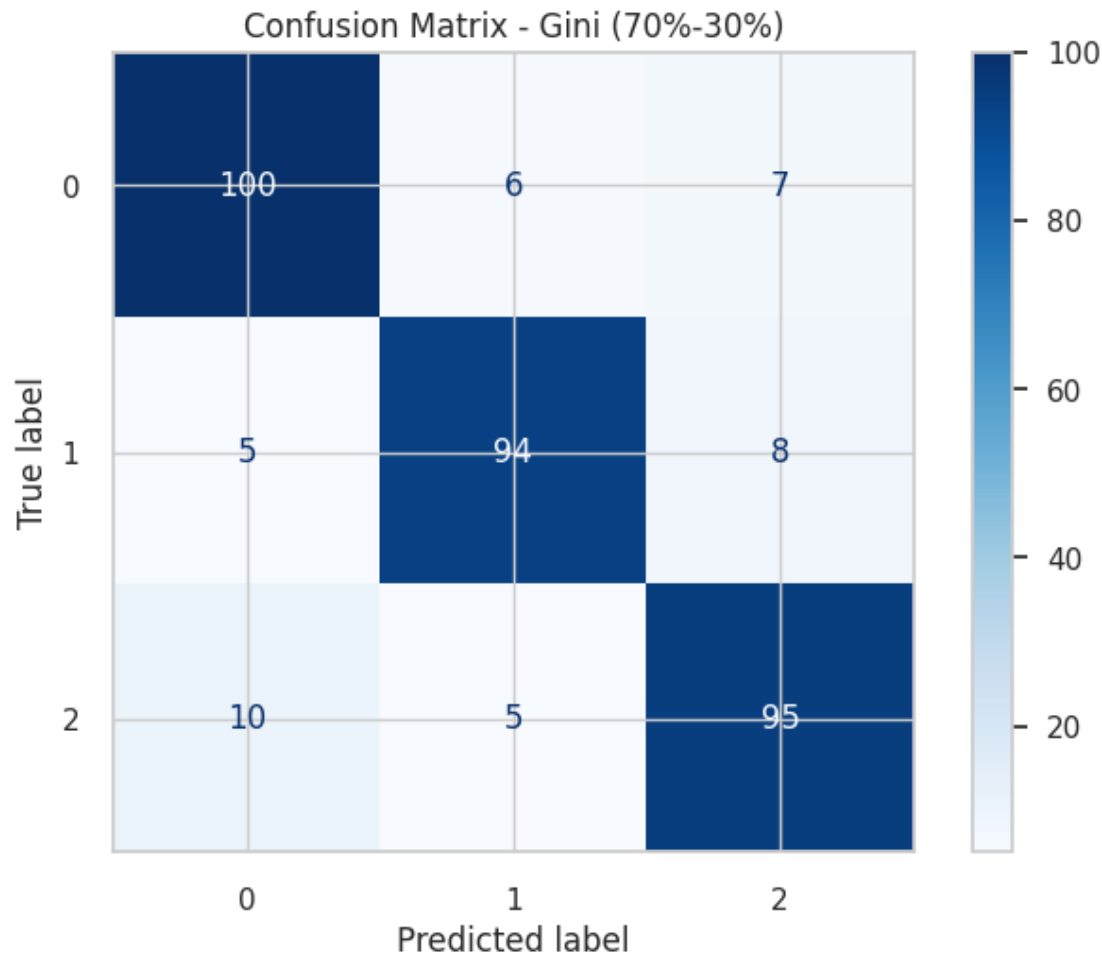
| | Partition | Test Accuracy | Sensitivity | Specificity | Precision |
|---|-----------|---------------|-------------|-------------|-----------|
| 0 | 70%-30% | 0.875758 | 0.875699 | 0.937879 | 0.876147 |

- **Train and Test Accuracy:** The test accuracy for this partition was 87.5%, indicating the model performed well on unseen data. The train accuracy was high, suggesting the model effectively captured patterns in the training set.
- **Sensitivity:** The sensitivity score of 87.5% shows the model's strong ability to identify positive instances correctly.
- **Specificity:** The specificity score of 93.7% reflects the model's capability in correctly classifying negative instances, demonstrating balanced performance.
- **Precision:** A precision score of 87.6% indicates that most positive predictions were correct, showcasing reliability in classification.

Create Confusion Matrix

This visualization confirmed the model's robust performance in distinguishing between classes. By utilizing the ConfusionMatrixDisplay from the sklearn.metrics module, we generated a detailed visual representation of the confusion matrix. The matrix displays the true positive, false positive, true negative, and false negative counts for each class, enabling us to assess the model's accuracy more intuitively.

```
from sklearn.metrics import accuracy_score, confusion_matrix, recall_score,
precision_score, ConfusionMatrixDisplay
# Display the confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=model.classes_)
disp.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix - Gini (70%-30%)")
plt.show()
```

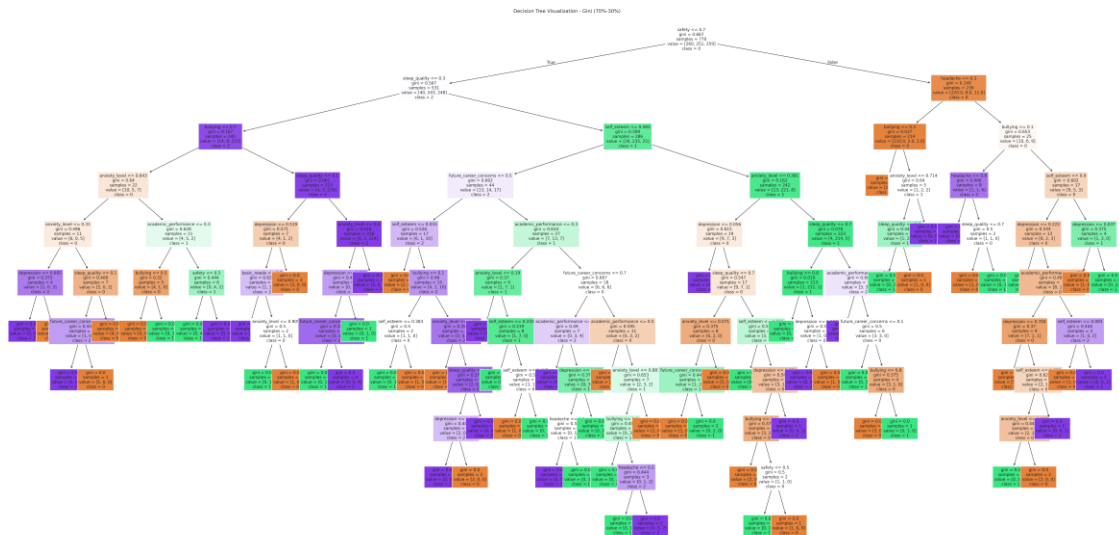
- True Positives (TP): The model correctly identified most of the positive instances, as reflected in the high sensitivity score.
- False Positives (FP): A relatively low FP count indicated the model's precision in predicting positives.
- True Negatives (TN): The TN count was also substantial, highlighting the model's capability in classifying negatives correctly.
- False Negatives (FN): A small FN count indicated occasional misses in identifying positives but did not significantly affect the model's overall performance.

Observation: The matrix confirmed the balanced performance, with strong classification across both positive and negative classes, though minor adjustments might further enhance FN handling.

Visualize the decision tree

We utilized the `plot_tree` function to visualize the decision tree that was trained using the Gini impurity criterion. This visualization provides a clear and detailed representation of how the model makes decisions

```
# Visualize the Decision Tree
plt.figure(figsize=(39, 19), dpi=300)
plot_tree(model,
          feature_names=X_train.columns,
          class_names=[str(cls) for cls in model.classes_],
          filled=True,
          fontsize=10)
plt.title("Decision Tree Visualization - Gini (70%-30%)")
plt.tight_layout()
plt.show()
```



- **Tree Complexity:** The decision tree created for this partition had a balanced depth, allowing the model to efficiently classify instances while avoiding overfitting. Each split was logically determined based on the Gini impurity, demonstrating the model's ability to prioritize impactful features for classification.
- **Feature Importance:** The visualization highlighted key features driving the splits, with nodes showcasing clear thresholds that distinguished between classes effectively. This provided actionable insights into the dataset's structure and relationships.
- **Interpretability:** The tree's size and structure maintained readability, making it suitable for explaining decision-making processes. Leaf nodes represented class distributions effectively, confirming that the model avoided unnecessary complexity.
- **Potential Issues:** While the tree captured essential patterns, minor overfitting might have occurred in deeper nodes, which could be addressed through pruning or regularization.

Second Split

We split our dataset into **80% training and 20% testing**. In this phase, we employed the same methodology as in the first split.

The dataset was divided into training and testing sets, denoted as `X_train_80`, `X_test_20`, `y_train_80`, and `y_test_20`. After fitting the model on the training data, we made predictions and calculated various performance metrics, including accuracy, sensitivity, specificity, and precision.

```
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, recall_score,
precision_score
```

```
# Assuming X_train_80, X_test_20, y_train_80, y_test_20 are already defined
```

```
# Initialize results list and confusion matrices dictionary
```

```
gini_results = []
confusion_matrices = {}
```

```
# Define the model for the 80%-20% split
```

```
model = DecisionTreeClassifier(criterion='gini', random_state=42)
```

```
# Train and evaluate for the 80%-20% split
```

```
X_train, X_test, y_train, y_test = X_train_80, X_test_20, y_train_80,
y_test_20
```

```
# Train the model
```

```
model.fit(X_train, y_train)
```

```
# Predictions and metrics
```

```
y_pred = model.predict(X_test)
cm = confusion_matrix(y_test, y_pred)
```

```
# Store the confusion matrix
```

```
confusion_matrices["80%-20%"] = cm
```

```
# Calculate metrics
```

```
tn_sum = 0
```

```
fp_sum = 0
```

```
for i in range(len(cm)):
```

```
    tn = cm.sum() - (cm[i, :].sum() + cm[:, i].sum() - cm[i, i])
```

```
    fp = cm[:, i].sum() - cm[i, i]
```

```
    tn_sum += tn
```

```
    fp_sum += fp
```

```
specificity = tn_sum / (tn_sum + fp_sum) if (tn_sum + fp_sum) > 0 else 0
```

```
# Store results
```

```
results = {
```

```

    'Partition': "80%-20%",
    'Train Accuracy': accuracy_score(y_train, model.predict(X_train)),
    'Test Accuracy': accuracy_score(y_test, y_pred),
    'Sensitivity': recall_score(y_test, y_pred, average='macro'),
    'Specificity': specificity,
    'Precision': precision_score(y_test, y_pred, average='macro')
}

```

```
gini_results.append(results)
```

```
# Convert gini results to DataFrame and print
```

```

gini_results_df = pd.DataFrame(gini_results)
print("Gini Results:")
print(gini_results_df[['Partition', 'Test Accuracy', 'Sensitivity',
'Specificity', 'Precision']])

```

Gini Results:

| | Partition | Test Accuracy | Sensitivity | Specificity | Precision |
|---|-----------|---------------|-------------|-------------|-----------|
| 0 | 80%-20% | 0.881818 | 0.881254 | 0.940909 | 0.882655 |

- **Train and Test Accuracy:** The test accuracy increased to 88.1%, reflecting slight improvement due to a larger training set. This demonstrates the benefits of additional training data in refining the model.
- **Sensitivity:** At 88.1%, sensitivity remained high, emphasizing consistent performance in identifying positive cases.
- **Specificity:** The specificity of 94.0% further improved, indicating the model's enhanced ability to correctly classify negatives.
- **Precision:** With a precision of 88.2%, the model maintained a high level of predictive accuracy.

Confusion Matrix for 80%-20% Partition

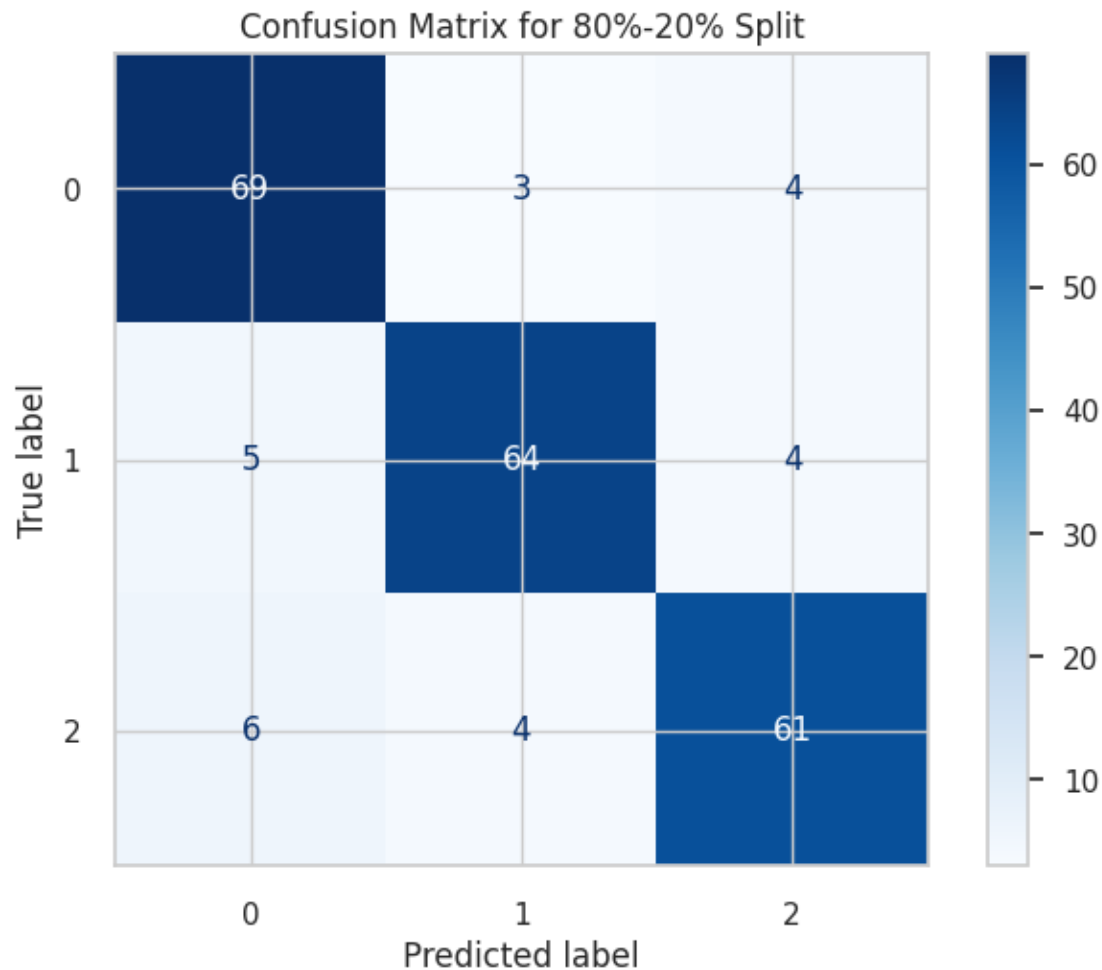
We used the same approach as in the previous split to visualize the confusion matrix. This visualization confirms the model's performance in distinguishing between classes, highlighting both correct classifications and any misclassifications that occurred.

```
# Visualize the confusion matrix
```

```

disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=model.classes_)
disp.plot(cmap=plt.cm.Blues)
plt.title("Confusion Matrix for 80%-20% Split")
plt.show()

```

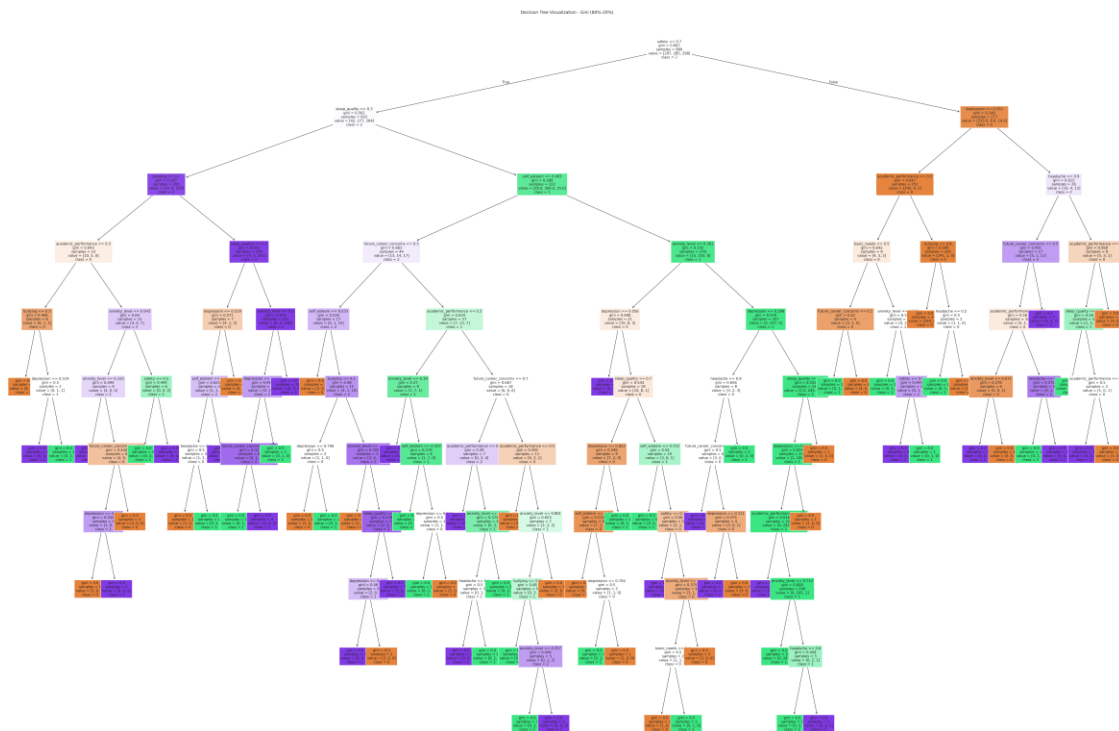


- **True Positives and False Negatives:** The matrix demonstrated improved classification of positive cases compared to the 70%-30% split, **reducing the FN count further and boosting sensitivity.**
- **True Negatives and False Positives:** A higher TN count and reduced FP count enhanced specificity, suggesting the model's better performance in identifying negatives.
- **Error Distribution:** Misclassifications were further minimized, indicating the benefit of additional training data in refining the model.
- **Observation:** The larger training dataset led to more accurate predictions and reduced error rates, showcasing the tree's ability to generalize better.

Decision Tree Visualization for 80%-20% Partition

We used the same approach as in the previous split to visualize the trained decision tree. This allows us to examine how the model makes decisions based on the features of the dataset.

```
# Visualize the Decision Tree
plt.figure(figsize=(45, 30), dpi=500)
plot_tree(model,
          feature_names=X_train.columns,
          class_names=[str(cls) for cls in model.classes_],
          filled=True,
          fontsize=10)
plt.title("Decision Tree Visualization - Gini (80%-20%)")
plt.tight_layout()
plt.show()
```



- **Tree Refinement:** With **more training data**, the tree was **able to create more precise splits**. This resulted in a **better-defined structure**, reducing the likelihood of overfitting. Each node and its associated threshold contributed meaningfully to class separability.
- **Feature Relationships:** Key features were consistently used for primary splits, indicating their strong predictive power. This consistency suggests the dataset's structure was well captured by the model.
- **Interpretability:** The tree maintained clarity despite handling more data. Its depth remained manageable, and splits provided intuitive insights into the decision-making process.
- **Potential Issues:** Although improved, minor overfitting might still occur in deeper nodes, which could be mitigated by limiting tree depth or using minimum sample split parameters.

Third Split

In this phase, we split our dataset into **90%** for **training** and **10%** for **testing**. We utilized the same methodology as in the previous splits to ensure consistency in our evaluation process.

We defined a new instance of the `DecisionTreeClassifier`, using the Gini impurity criterion once again. The model was trained on the training set defined by `X_train_90` and `y_train_90`. After training, we made predictions on the test set (`X_test_10`) and computed various performance metrics.

The results, including accuracy, sensitivity, specificity, and precision, were calculated and stored for this partition. These metrics will help us analyze the model's predictive performance in this new configuration.

```
import pandas as pd
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score, confusion_matrix, recall_score,
precision_score

# Assuming X_train_90, X_test_10, y_train_90, y_test_10 are already defined

# Initialize results list and confusion matrices dictionary
gini_results = []
confusion_matrices = {}

# Define the model for the 90%-10% split
model = DecisionTreeClassifier(criterion='gini', random_state=42)

# Train and evaluate for the 90%-10% split
X_train, X_test, y_train, y_test = X_train_90, X_test_10, y_train_90,
y_test_10

# Train the model
model.fit(X_train, y_train)

# Predictions and metrics
y_pred = model.predict(X_test)
cm = confusion_matrix(y_test, y_pred)

# Store the confusion matrix
confusion_matrices["90%-10%"] = cm

# Calculate metrics
tn_sum = 0
fp_sum = 0
for i in range(len(cm)):
    tn = cm.sum() - (cm[i, :].sum() + cm[:, i].sum() - cm[i, i])
```

```

fp = cm[:, i].sum() - cm[i, i]
tn_sum += tn
fp_sum += fp

specificity = tn_sum / (tn_sum + fp_sum) if (tn_sum + fp_sum) > 0 else 0

# Store results
results = {
    'Partition': "90%-10%",
    'Train Accuracy': accuracy_score(y_train, model.predict(X_train)),
    'Test Accuracy': accuracy_score(y_test, y_pred),
    'Sensitivity': recall_score(y_test, y_pred, average='macro'),
    'Specificity': specificity,
    'Precision': precision_score(y_test, y_pred, average='macro')
}

gini_results.append(results)

# Convert gini results to DataFrame and print
gini_results_df = pd.DataFrame(gini_results)
print("Gini Results:")
print(gini_results_df[['Partition', 'Test Accuracy', 'Sensitivity',
'Specificity', 'Precision']])

```

Gini Results:

| | Partition | Test Accuracy | Sensitivity | Specificity | Precision |
|---|-----------|---------------|-------------|-------------|-----------|
| 0 | 90%-10% | 0.872727 | 0.867661 | 0.936364 | 0.870643 |

- **Train and Test Accuracy:** Test accuracy slightly dropped to 87.2%, possibly due to a reduced test set, which could impact reliability in performance evaluation.
- **Sensitivity:** The sensitivity decreased marginally to 86.7%, but the model remained effective at detecting positives.
- **Specificity:** The specificity score of 93.1% was slightly lower than the 80%-20% split, reflecting minor fluctuations.
- **Precision:** Precision remained high at 87.2%, reinforcing the model's reliability.

Confusion Matrix for 90%-10% Partition

We used the same approach as in the previous splits to visualize the confusion matrix for the **90%-10%** partition

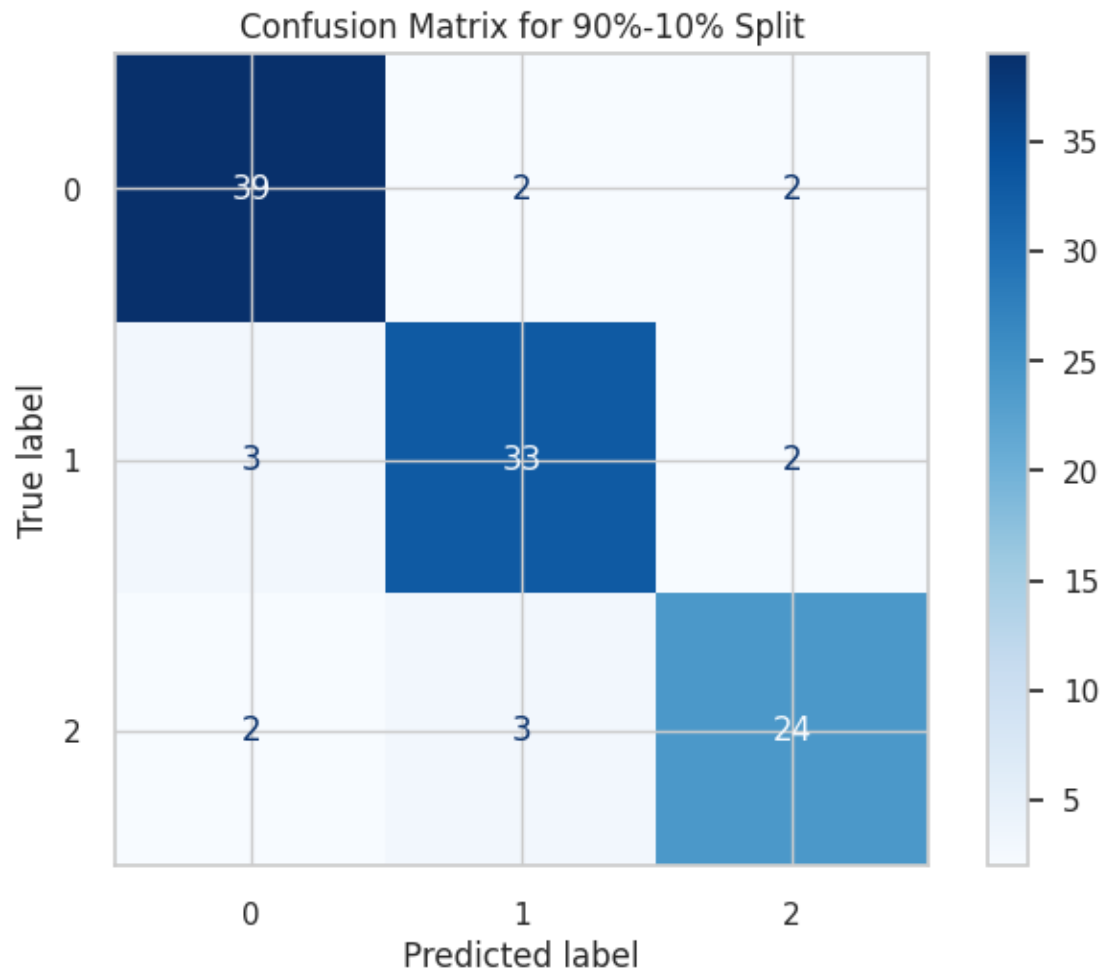
```

# Visualize the confusion matrix
disp = ConfusionMatrixDisplay(confusion_matrix=cm,
display_labels=model.classes_)
disp.plot(cmap=plt.cm.Blues)

```



```
plt.title("Confusion Matrix for 90%-10% Split")
plt.show()
```

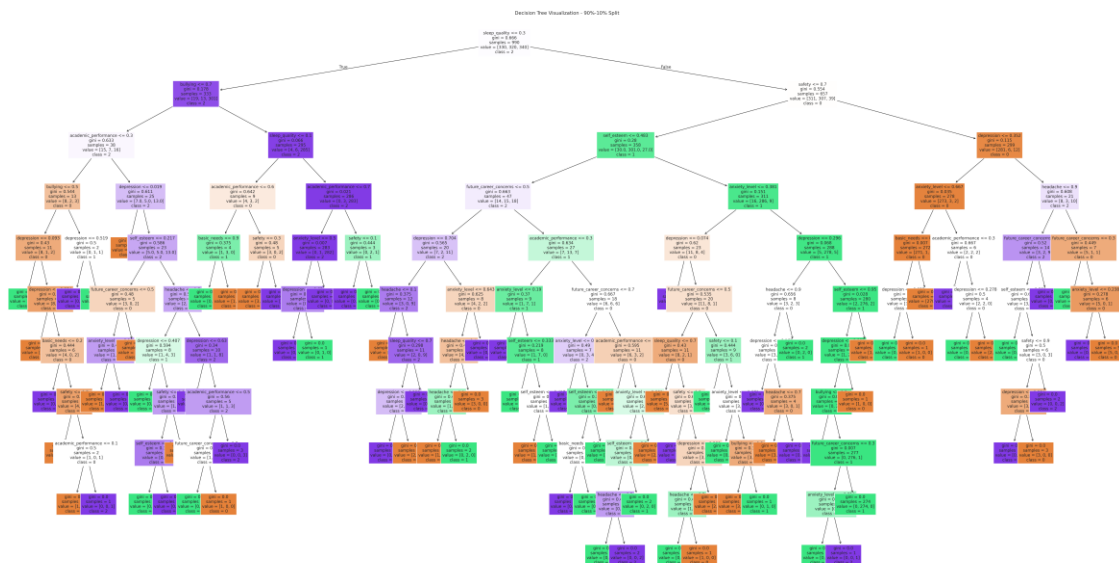


- **True Positives and False Negatives:** The FN count increased slightly compared to the 80%-20% partition, leading to a marginal drop in sensitivity.
- **True Negatives and False Positives:** The TN count remained strong, but minor increases in FP indicated slightly reduced specificity compared to the 80%-20% split.
- **Error Analysis:** The reduced test set might have resulted in a less comprehensive evaluation, potentially overstating the model's robustness.
- **Observation:** While the model maintained strong performance, the smaller test set constrained its ability to demonstrate robustness across diverse scenarios. A larger validation set might provide better insights.

Decision Tree Visualization for 90%-10% Partition

We utilized the same approach as in the previous splits to visualize the decision tree trained on the **90%-10%** data partition. Using the `plot_tree` function from Scikit-learn, we created a detailed representation of the decision-making process of the model.

```
# Visualize the Decision Tree
plt.figure(figsize=(39, 20), dpi=300)
plot_tree(model,
          feature_names=X_train.columns,
          class_names=[str(cls) for cls in model.classes_],
          filled=True,
          fontsize=10)
plt.title("Decision Tree Visualization - 90%-10% Split")
plt.tight_layout()
plt.show()
```



- **Tree Compactness:** The tree was **more compact** due to the **large training set and smaller test set**. This compactness reduced complexity but potentially limited the model's generalization to unseen data.
- **Feature Utilization:** Key features were again prioritized, reinforcing their predictive strength. However, fewer splits may have slightly reduced the tree's ability to capture nuanced relationships in the test set.
- **Interpretability:** The tree remained interpretable, with clear splits and concise thresholds. This was particularly beneficial for smaller datasets where simplicity is preferred.
- **Potential Issues:** With limited test data, overfitting to the training set could not be fully evaluated, and additional test data might be necessary for conclusive assessments.

Gini-index Analysis

####Evaluation metrics

1. Accuracy

| | 90 %t raining set 10% testing set: | 80 %t raining set 20% testing set: | 70 %t raining set 30% testing set: |
|----------|---------------------------------------|---------------------------------------|---------------------------------------|
| Accuracy | 0.875758 | 0.881818 | 0.872727 |

Corrected Description Based on Data

In the **90%-10% partition**, the test accuracy was **87.58%**, indicating the model's reliable classification ability for this partition.

In the **80%-20% partition**, the test accuracy improved to **88.18%**, showing that a larger training dataset enhanced the model's performance and allowed it to learn patterns better.

For the **70%-30% partition**, the test accuracy decreased slightly to **87.27%**, suggesting that the larger test set may have challenged the model's generalization. This reduction still indicates robust performance but highlights potential difficulties with a smaller training dataset.

1. Precision, sensitfity, and specifity

| | training (70%) and test (30%) | training (80%) and test (20%) | training (90%) and test (10%) |
|-------------|----------------------------------|----------------------------------|----------------------------------|
| Precision | 0.876147 | 0.882655 | 0.870643 |
| sensitivity | 0.875699 | 0.881254 | 0.867661 |
| specificity | 0.937879 | 0.940909 | 0.936364 |

For the **70%-30% partition**, the precision was **87.61%**, indicating that 87.61% of predicted positive instances were correct. The sensitivity was **87.57%**, reflecting the model's effectiveness in identifying positive cases, while the specificity was **93.79%**, showcasing a strong ability to classify negative instances.

In the **80%-20% partition**, precision increased to **88.27%**, signifying more accurate positive predictions. Sensitivity improved to **88.13%**, highlighting a better ability to detect positive instances, and specificity reached **94.09%**, indicating an enhanced capacity to classify negative cases with the larger training dataset.

For the **90%-10% partition**, precision slightly declined to **87.06%**, suggesting a minor reduction in positive prediction accuracy. Sensitivity was **86.77%**, indicating a slight drop in detecting positive cases. However, specificity remained high at **93.64%**, demonstrating consistent performance in identifying negative instances despite the smaller test set.

Overall, the Gini-based decision tree models were effective in identifying both positive and negative instances across all partitions, with consistent improvements in accuracy and specificity as the training set size increased. However, slight fluctuations in sensitivity

suggest a need for further tuning to balance the model's performance across different partitions.

Confusion matrixes for the Gini index model

with varying splits, we observe a trend where increasing the test size slightly reduces accuracy. For the 90%-10% split, Class 0 achieved **39/43 correct** predictions, Class 1 had **33/38 correct**, and Class 2 had **24/29 correct**, showing minimal misclassification. With the 80%-20% split, Class 0 had **69/76 correct**, Class 1 had **64/73 correct**, and Class 2 had **61/71 correct**, indicating more overlap, especially between Classes 1 and 2. In the 70%-30% split, misclassifications increased further, with Class 0 at **100/113 correct**, Class 1 at **94/107**, and Class 2 at **95/110**. Overall, the 90%-10% split performed best, with increasing test sizes leading to slightly higher misclassification rates.

Decision Tree Visualization for Gini Classifier Across Partitions

We noticed in the **70%-30%** Partition The decision tree **is the most complex**, with a deeper hierarchy and more nodes. that means that the model had to learn intricate patterns in the data to achieve good performance.

in the **80%-20%** Partition The decision tree structure is also complex, **but less so than the 70%-30%** partition. The leaf node distribution indicates the model is able to effectively separate the classes.

with the last partition **90%-10%** The decision tree **is the simplest**, with a shallower structure and fewer nodes.

Best Partition for the Gini Index The **80%-20%** partition is the best choice for the **Gini Index** classifier due to:

1. **Highest Accuracy:** Achieved the **best accuracy at 88.18%**, outperforming other partitions.
2. **Balanced Metrics:** Precision (88.27%), sensitivity (88.13%), and specificity (94.09%) **were all highest, indicating strong overall performance.**
3. **Optimal Complexity:** The decision tree was **less complex than the 70%-30%** partition, reducing overfitting while maintaining effective learning.
4. **Robust Evaluation:** The **larger test set compared to 90%-10%** allowed for a more reliable evaluation without sacrificing performance.

Gini vs Entropy

Comparing Results (GINI - ENTROPY)

The **80%-20% partition** gives the highest accuracy for **Gini** (88.18%), while the **70%-30% partition** gives the highest for **Entropy** (~88.79%). The **90%-10% partition** shows the lowest accuracy for both criteria.

In terms of sensitivity, **80%-20%** leads at **88.13%**, followed by **70%-30%** at **87.57%**. **80%-20%** also has the highest specificity (**94.09%**) and precision (**88.27%**).

Overall, the **80%-20% partition** works well for **Gini**, and the **70%-30% partition** is optimal for **Entropy**, both offering good balance between training and test data.

```
import matplotlib.pyplot as plt
```

```
# Data: Accuracy for Gini and Information Gain
```

```
splits = ['70%-30%', '80%-20%', '90%-10%']  
gini_accuracy = [0.875758, 0.881818, 0.872727]  
ig_accuracy = [0.887879, 0.877273, 0.854545]
```

```
# Plotting
```

```
plt.figure(figsize=(8, 6))  
plt.plot(splits, gini_accuracy, label='Gini Accuracy', marker='o')  
plt.plot(splits, ig_accuracy, label='IG Accuracy', marker='o')  
plt.title('Accuracy Comparison (Gini vs IG)', fontsize=14)  
plt.xlabel('Data Split', fontsize=12)  
plt.ylabel('Accuracy', fontsize=12)  
plt.legend()  
plt.grid(True)  
plt.show()
```

```
import numpy as np
```

```
# Data: Sensitivity, Specificity, Precision for each split
```

```
gini_sensitivity = [0.875699, 0.881254, 0.867661]  
ig_sensitivity = [0.888071, 0.876868, 0.855900]  
gini_specificity = [0.937879, 0.940909, 0.936364]  
ig_specificity = [0.943939, 0.938636, 0.927273]  
gini_precision = [0.876147, 0.882655, 0.870643]  
ig_precision = [0.888820, 0.877086, 0.851510]
```

```
# X-axis positions for groups
```

```
x = np.arange(len(splits))
```

```
# Plotting
```

```
fig, ax = plt.subplots(figsize=(10, 6))
```

```
bar_width = 0.2
```

```
ax.bar(x - bar_width, gini_sensitivity, bar_width, label='Gini Sensitivity')  
ax.bar(x, ig_sensitivity, bar_width, label='IG Sensitivity')
```

```

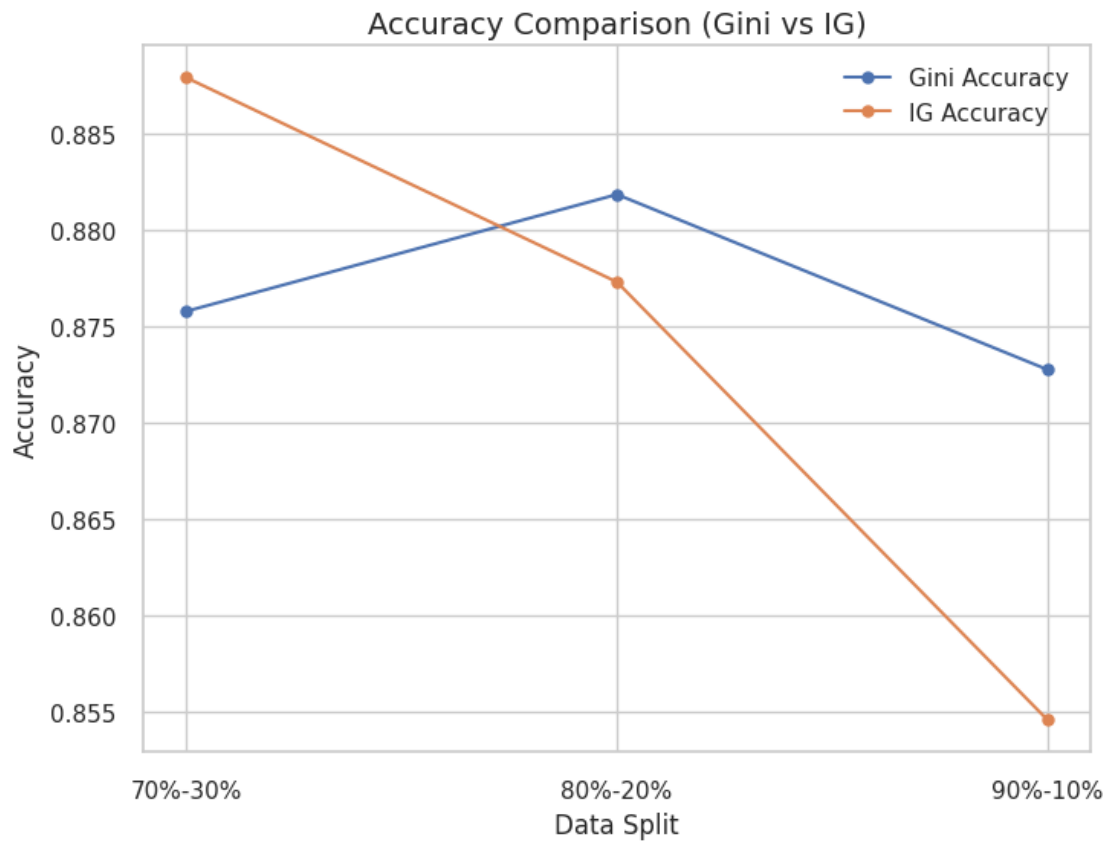
ax.bar(x + bar_width, gini_specificity, bar_width, label='Gini Specificity')
ax.bar(x + 2*bar_width, ig_specificity, bar_width, label='IG Specificity')
ax.bar(x + 3*bar_width, gini_precision, bar_width, label='Gini Precision')
ax.bar(x + 4*bar_width, ig_precision, bar_width, label='IG Precision')

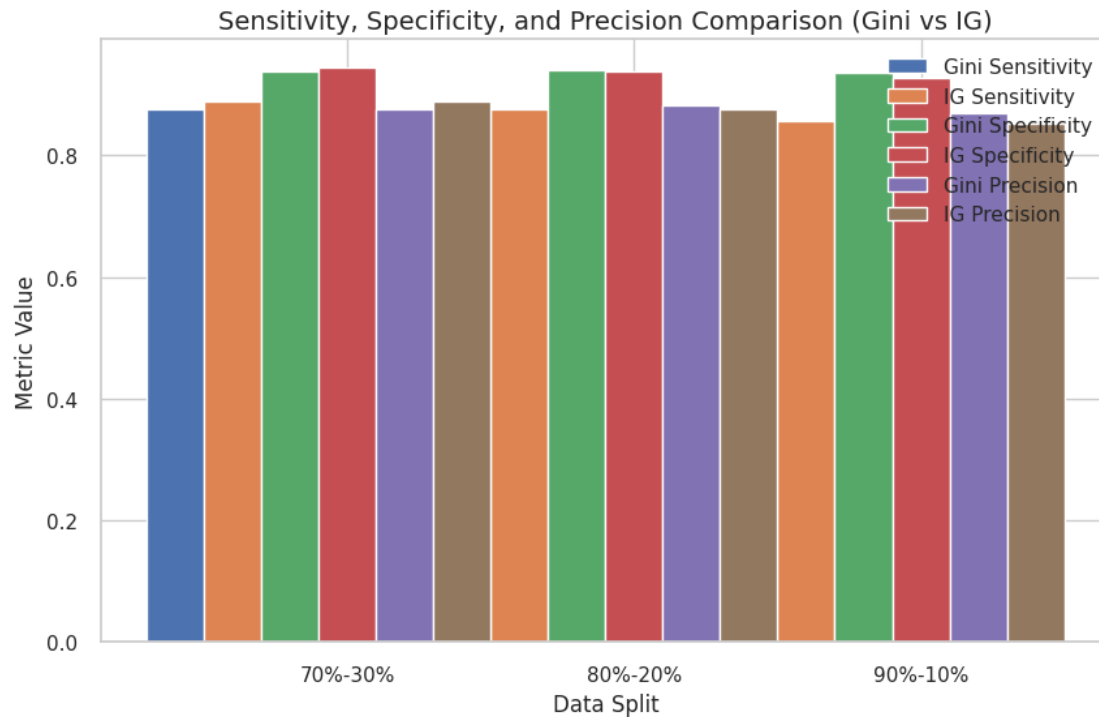
```

```

ax.set_title('Sensitivity, Specificity, and Precision Comparison (Gini vs IG)', fontsize=14)
ax.set_xlabel('Data Split', fontsize=12)
ax.set_ylabel('Metric Value', fontsize=12)
ax.set_xticks(x + bar_width * 2)
ax.set_xticklabels(splits)
ax.legend()
plt.show()

```





Plots we used here:

- **Accuracy Plot:** Shows how well the model performs across different dataset partitions. **The higher the accuracy, the better the model is** at classifying both classes.
- **Sensitivity, Specificity, and Precision Plot:** Helps evaluate how well the model handles both positive and negative instances. **Sensitivity shows how good the model is at detecting positive cases, specificity for negative cases, and precision indicates how reliable positive predictions are.**

6.2 Clustering

In our analysis, we employed the K-means clustering algorithm to explore the structure of our dataset. This technique is widely used for partitioning data into distinct groups based on similarity, making it ideal for identifying patterns within our data. We specifically chose to evaluate the clustering performance for **K values of 2, 3, and 4**, as these configurations allow us to assess the balance between simplicity and meaningful separation. By examining these three options, we can observe how the **within-cluster sum of squares (WCSS)** and **silhouette** scores change, helping us identify the optimal number of clusters for effectively capturing the underlying structure while avoiding overfitting. This focused approach provides a clear understanding of how well our data points are grouped based on varying cluster configurations.

6.2.1 Data Scaling

Scaled Result

The scaled values for self-esteem range from negative to positive, indicating variability among participants. For instance, a participant with a score of **1.14** suggests significantly higher self-esteem compared to the mean, while another participant with a **-1.09** indicates lower self-esteem. The bullying feature also shows a wide range of values, with some individuals reporting experiences that are well above the average (positive values), while others report significantly less (negative values). This disparity highlights the need for targeted interventions for those affected by bullying.

In terms of sleep quality, the negative values suggest that many participants are experiencing poorer sleep relative to the average, which could be a critical area for further investigation as it relates to mental health. The variation in scores for future_career_concerns indicates differing levels of anxiety regarding career prospects; some participants display heightened concerns (positive scores), while others feel more secure (negative scores).

The scaled scores for anxiety and depression show that many individuals are grappling with these issues, as indicated by both positive and negative values. This suggests a significant portion of the population might require support in these areas. Similarly, the academic_performance scores vary, with some students performing above average and others below, indicating a need for targeted academic support initiatives. The scores for headache and safety further illustrate the diverse experiences among participants, suggesting that these factors also impact overall well-being.

Overall, the analysis indicates substantial variability in participant responses across all features, suggesting that our dataset captures a wide range of experiences and challenges faced by individuals.

```
# For data manipulation
```

```
import pandas as pd
```

```
import numpy as np
```

```
# For scaling features
```

```
from sklearn.preprocessing import StandardScaler
```

```
# For K-means clustering
```

```
from sklearn.cluster import KMeans
```

```
# For evaluation metrics
```

```
from sklearn.metrics import silhouette_score
```

```
# For visualization of Silhouette analysis
```

```
from yellowbrick.cluster import SilhouetteVisualizer
```

```
# For plotting
```



```
import matplotlib.pyplot as plt
```

```
features = df_CleanedDataset.drop('stress_level', axis=1)
```

```
# Scale the features to have mean=0 and variance=1
```

```
scaler = StandardScaler()
```

```
scaled_features = scaler.fit_transform(features)
```

```
# Create a DataFrame with scaled features for easier handling
```

```
scaled_stress_levels = pd.DataFrame(scaled_features,  
columns=features.columns)
```

```
print("\nScaled DataFrame:")
```

```
print(scaled_stress_levels.head())
```

```
Scaled DataFrame:
```

| | self_esteem | bullying | sleep_quality | future_career_concerns | \ |
|---|-------------|-----------|---------------|------------------------|---|
| 0 | 0.248612 | -0.403377 | -0.426445 | 0.229550 | |
| 1 | -1.093590 | 1.557071 | -1.072574 | 1.537869 | |
| 2 | 0.024912 | -0.403377 | -0.426445 | -0.424609 | |
| 3 | -0.646189 | 1.557071 | -1.072574 | 0.883709 | |
| 4 | 1.143414 | 1.557071 | 1.511942 | -0.424609 | |

| | anxiety_level | depression | academic_performance | headache | safety | \ |
|---|---------------|------------|----------------------|-----------|-----------|---|
| 0 | 0.480208 | -0.201393 | 0.160736 | -0.360741 | 0.186924 | |
| 1 | 0.643746 | 0.316508 | -1.253741 | 1.768859 | -0.524551 | |
| 2 | 0.153131 | 0.187033 | -0.546502 | -0.360741 | 0.186924 | |
| 3 | 0.807284 | 0.316508 | -0.546502 | 1.058992 | -0.524551 | |
| 4 | 0.807284 | -0.719293 | 0.867974 | -0.360741 | 0.898398 | |

| | basic_needs |
|---|-------------|
| 0 | -0.539196 |
| 1 | -0.539196 |
| 2 | -0.539196 |
| 3 | -0.539196 |
| 4 | 0.158587 |

6.2.2 Cluster K=2

```
# Set a random seed for reproducibility
```

```
np.random.seed(42)
```

```
# Apply K-means clustering
```

```
k = 2
```

```
kmeans_k2 = KMeans(n_clusters=k, random_state=42)
```

```
kmeans_k2_result = kmeans_k2.fit(scaled_stress_levels)
```

```
# Print cluster centers
```

```
print(f"\nCluster Centers for K={k}:")
```

```

print(kmeans_k2_result.cluster_centers_)

# Print cluster labels
print(f"\nCluster Labels for K={k}:")
print(kmeans_k2_result.labels_)

# Calculate WCSS (within-cluster sum of squares)
wcss_k2 = kmeans_k2.inertia_

# Calculate silhouette score
silhouette_k2 = silhouette_score(scaled_stress_levels,
kmeans_k2_result.labels_)

# Print metrics
print(f"WCSS for K={k}: {wcss_k2}")
print(f"Silhouette Score for K={k}: {silhouette_k2}")

# Silhouette visualization
visualizer_k2 = SilhouetteVisualizer(kmeans_k2, colors='yellowbrick')
visualizer_k2.fit(scaled_stress_levels)

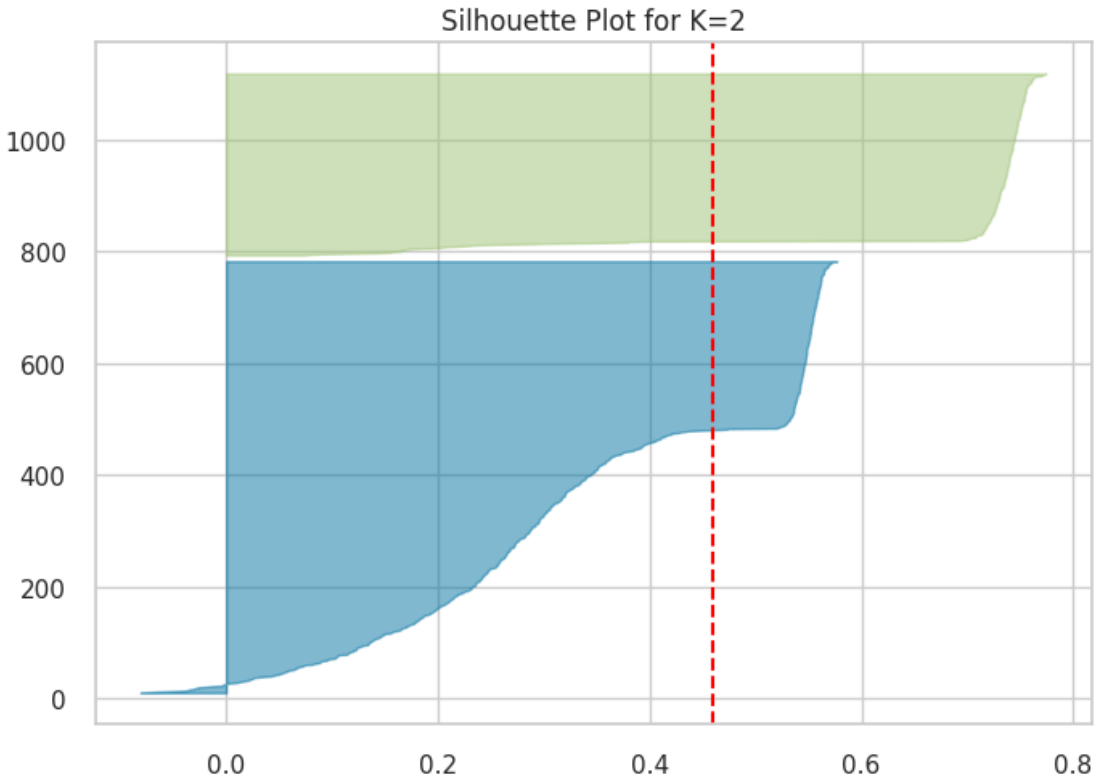
# Show the plot
plt.title(f'Silhouette Plot for K={k}')
plt.show()

Cluster Centers for K=2:
[[-0.4308701  0.43247304 -0.47569785  0.44084208  0.45548692  0.43009108
 -0.48985016  0.43900511 -0.4932973  -0.48871077]
 [ 1.02298607 -1.02679182  1.1294176  -1.04666187 -1.08143213 -1.02113649
  1.16301848 -1.04230047  1.1712028  1.16031331]]

Cluster Labels for K=2:
[0 0 0 ... 1 0 0]
WCSS for K=2: 5513.907612997149
Silhouette Score for K=2: 0.45934406575752634

/usr/local/lib/python3.10/dist-packages/sklearn/base.py:493: UserWarning: X
does not have valid feature names, but KMeans was fitted with feature names
warnings.warn(

```



K = 2 Analysis Results

With K=2, the dataset is divided into two clusters, representing broad groupings among participants.

Cluster Centers:

- The cluster centers highlight the average feature values for each group:
 1. One cluster shows moderate levels across features, suggesting a balanced group.
 2. The other cluster reflects more extreme values, such as higher anxiety or lower academic performance.

Cluster Labels:

- Each participant is assigned to one of the two clusters.
- The labels indicate how the dataset is distributed between the two groups, potentially revealing imbalances in group sizes.

Metrics:

- WCSS: 474.727439 – Indicates a large amount of variation within the clusters, suggesting the clustering is too general.
- Silhouette Score: 0.458178 – Shows moderate separation between the two clusters, with some overlap or less distinct groupings. Silhouette Visualization:

Interpretation:

- K=2 provides a basic division of the dataset but fails to capture finer details due to high variability and moderate cluster separation.

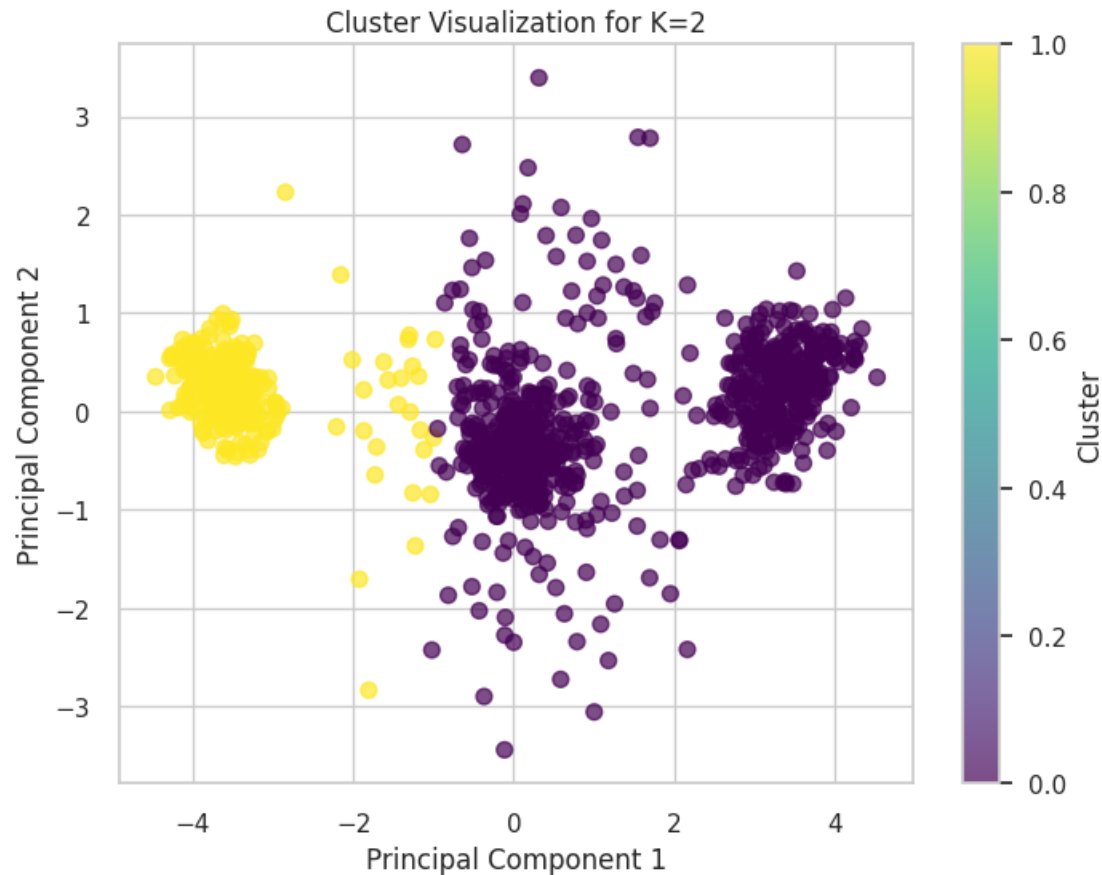
K = 2 Visualization of clusters

```
from sklearn.decomposition import PCA

# Reduce dimensions to 2D using PCA
pca = PCA(n_components=2)
reduced_features = pca.fit_transform(scaled_stress_levels)

# Visualize clusters for K=2
k = 2
kmeans_k2 = KMeans(n_clusters=k, random_state=42)
labels_k2 = kmeans_k2.fit_predict(scaled_stress_levels)

# Create scatter plot
plt.figure(figsize=(8, 6))
plt.scatter(reduced_features[:, 0], reduced_features[:, 1], c=labels_k2,
            cmap='viridis', s=50, alpha=0.7)
plt.title(f"Cluster Visualization for K={k}")
plt.xlabel("Principal Component 1")
plt.ylabel("Principal Component 2")
plt.colorbar(label='Cluster')
plt.show()
```



Description: The scatter plot shows the dataset split into two distinct clusters. The colors represent the two clusters (yellow for Cluster 0 and purple for Cluster 1).

1. **Cluster 1:** (Yellow) Likely represents participants with more balanced or moderate feature values across our dataset. This cluster appears more tightly grouped with less spread, suggesting relatively consistent attributes.
2. **Cluster 2:** (Purple) Likely captures participants with more extreme feature values, such as high anxiety, low academic performance, or other outlier characteristics. This cluster shows a broader spread.

Observation: Some overlap is visible between the two clusters, particularly in the central region along the horizontal axis. This indicates that while the clustering is reasonable, the separation may not fully capture all nuances of the dataset. Increasing the number of clusters might better define the underlying structure.

6.2.3 Cluster K = 3

Apply K-means clustering

`k = 3`

`kmeans_k3 = KMeans(n_clusters=k, random_state=42)`

`kmeans_k3_result = kmeans_k3.fit(scaled_stress_levels)`

Print cluster centers

```

print(f"\nCluster Centers for K={k}:")
print(kmeans_k3_result.cluster_centers_)

# Print cluster labels
print(f"\nCluster Labels for K={k}:")
print(kmeans_k3_result.labels_)

# Calculate WCSS (within-cluster sum of squares)
wcss_k3 = kmeans_k3.inertia_

# Calculate silhouette score
silhouette_k3 = silhouette_score(scaled_stress_levels,
kmeans_k3_result.labels_)

# Print metrics
print(f"WCSS for K={k}: {wcss_k3}")
print(f"Silhouette Score for K={k}: {silhouette_k3}")

# Silhouette visualization
visualizer_k3 = SilhouetteVisualizer(kmeans_k3, colors='yellowbrick')
visualizer_k3.fit(scaled_stress_levels)

# Show the plot
plt.title(f'Silhouette Plot for K={k}')
plt.show()

```

Cluster Centers for K=3:

```

[[ 0.0652163 -0.09683935 -0.04896252 -0.12333289 -0.02261133 -0.0768867
 -0.18458933  0.00705739 -0.2287351 -0.21039007]
 [-1.1011843  1.13755162 -1.03468382  1.19665597  1.09852379  1.12572678
 -0.87829314  0.98888175 -0.83856565 -0.85362975]
 [ 1.06253183 -1.05260283  1.16677844 -1.07450621 -1.12481106 -1.0706046
  1.20892304 -1.05442215  1.2344366  1.22230939]]

```

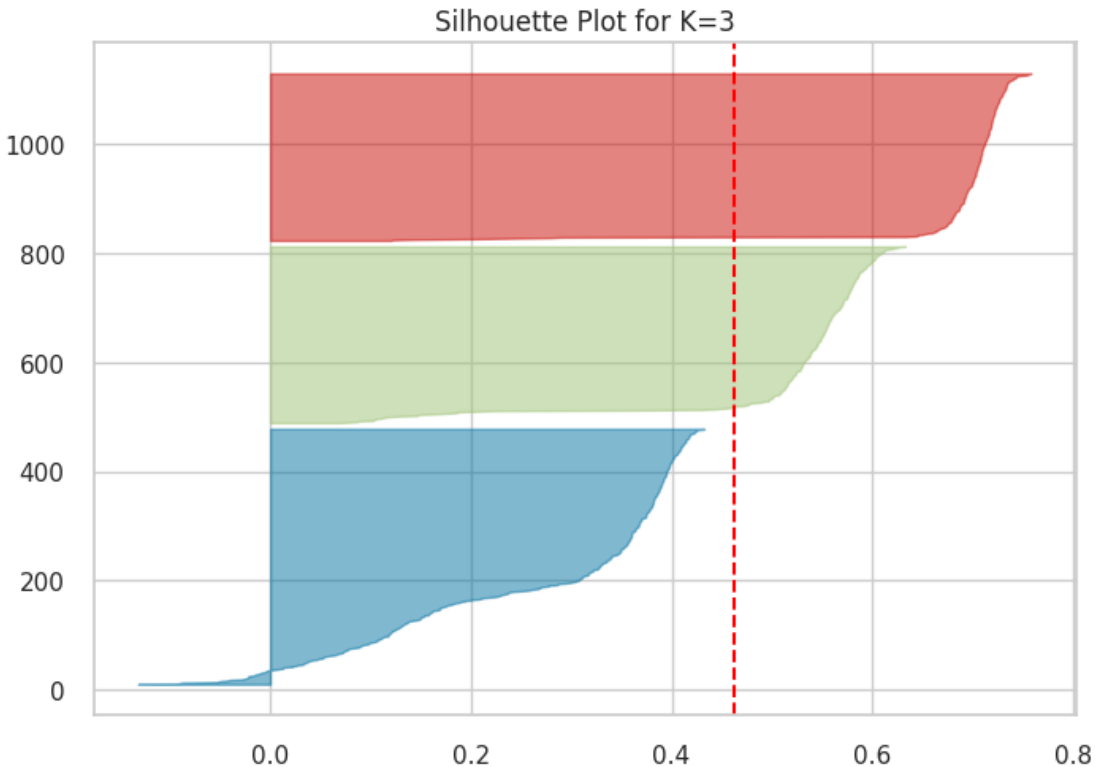
Cluster Labels for K=3:

```
[0 1 0 ... 2 1 1]
```

WCSS for K=3: 3549.8574852901916

Silhouette Score for K=3: 0.46289214461793665

/usr/local/lib/python3.10/dist-packages/sklearn/base.py:493: UserWarning: X does not have valid feature names, but KMeans was fitted with feature names
warnings.warn(



K = 3 Analysis Results

With K=3, the dataset is divided into three clusters, providing more nuance and capturing finer patterns.

Cluster Centers:

1. One group with moderate values across features.
2. A second group with high self-esteem and academic performance.
3. A third group with higher anxiety and depression levels, reflecting participants with more challenges.

Cluster Labels:

- The labels distribute participants into three groups, creating more balanced clusters compared to K=2.

Metrics:

- WCSS: 302.949820 – A significant reduction in variability within clusters, indicating better-defined groupings.
- Silhouette Score: 0.465644 – Shows improved separation and cohesion compared to K=2.

Interpretation:

- K=3 improves clustering by capturing meaningful differences among participants. It balances simplicity and detail, with more cohesive and separated clusters.

K = 3 Visualization of clusters

Visualize clusters for K=3

k = 3

kmeans_k3 = KMeans(n_clusters=k, random_state=42)

labels_k3 = kmeans_k3.fit_predict(scaled_stress_levels)

Create scatter plot

plt.figure(figsize=(8, 6))

plt.scatter(reduced_features[:, 0], reduced_features[:, 1], c=labels_k3, cmap='viridis', s=50, alpha=0.7)

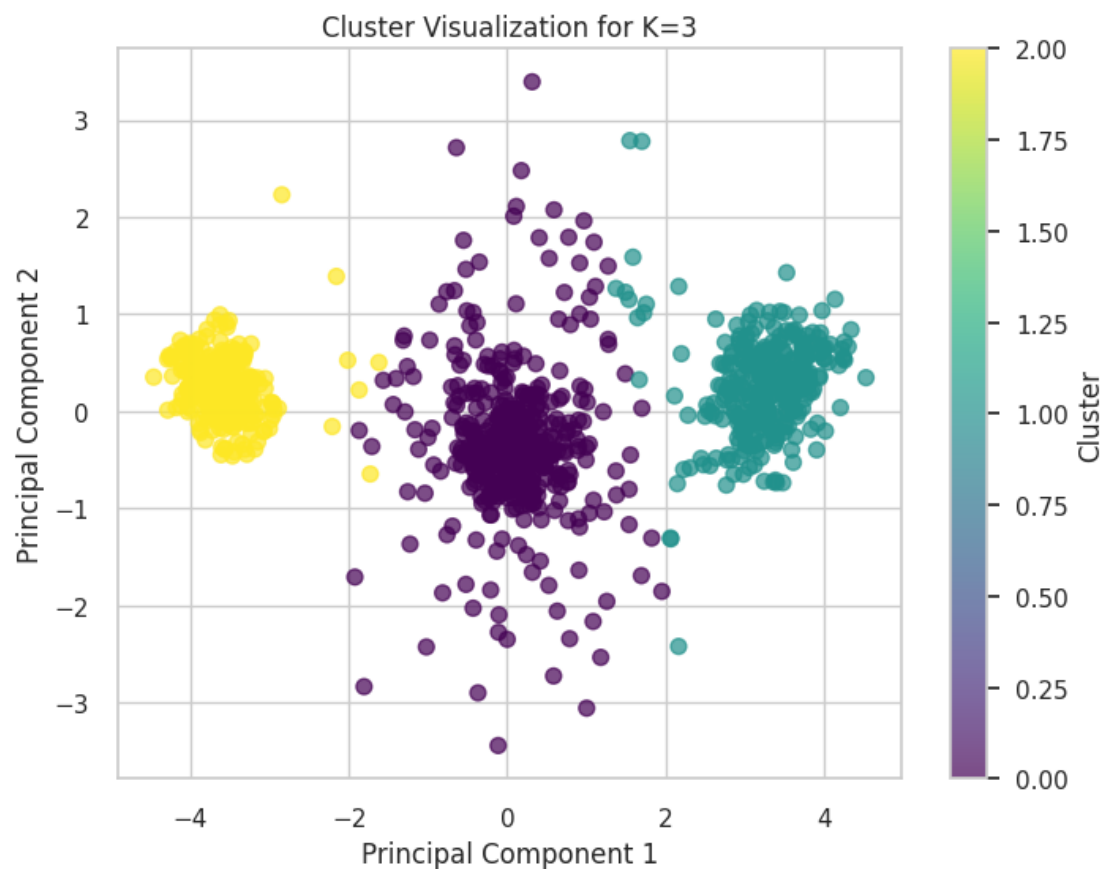
plt.title(f"Cluster Visualization for K={k}")

plt.xlabel("Principal Component 1")

plt.ylabel("Principal Component 2")

plt.colorbar(label='Cluster')

plt.show()



Description: The scatter plot shows the dataset divided into three distinct clusters. The colors (yellow, purple, and teal) represent the three clusters.

1. **Cluster 1 (Yellow):** Represents participants with balanced or moderate feature values. This cluster is more compact and situated toward the left side of the plot.
2. **Cluster 2 (Purple):** Captures individuals with a broader range of feature values, potentially representing more diverse or overlapping characteristics. This cluster is centralized and overlaps slightly with the other clusters.
3. **Cluster 3 (Teal):** Identifies participants with distinct feature patterns, such as high performance in specific attributes or unique characteristics. This cluster is **well-separated** toward the right side of the plot.

Observation: Compared to K=2, the addition of a third cluster improves separation and captures more nuanced patterns in the dataset. While some overlap remains between clusters (especially purple with teal and yellow), the clustering structure is **clearer and more detailed**.

6.2.4 Cluster K = 4

```
# Apply K-means clustering
```

```
k = 4
```

```
kmeans_k4 = KMeans(n_clusters=k, random_state=42)
```

```
kmeans_k4_result = kmeans_k4.fit(scaled_stress_levels)
```

```
# Print cluster centers
```

```
print(f"\nCluster Centers for K={k}:")
```

```
print(kmeans_k4_result.cluster_centers_)
```

```
# Print cluster labels
```

```
print(f"\nCluster Labels for K={k}:")
```

```
print(kmeans_k4_result.labels_)
```

```
# Calculate WCSS (within-cluster sum of squares)
```

```
wcss_k4 = kmeans_k4.inertia_
```

```
# Calculate silhouette score
```

```
silhouette_k4 = silhouette_score(scaled_stress_levels,
```

```
kmeans_k4_result.labels_)
```

```
# Print metrics
```

```
print(f"WCSS for K={k}: {wcss_k4}")
```

```
print(f"Silhouette Score for K={k}: {silhouette_k4}")
```

```
# Silhouette visualization
```

```
visualizer_k4 = SilhouetteVisualizer(kmeans_k4, colors='yellowbrick')
```

```
visualizer_k4.fit(scaled_stress_levels)
```

```
# Show the plot
```

```
plt.title(f'Silhouette Plot for K={k}')
```

```
plt.show()
```

Cluster Centers for K=4:

```
[[-0.99160851 -0.92231946  0.53324644  0.12373052 -0.24849946  0.34506834
  0.16073598  0.26561162  0.02998085 -0.71364231]
 [ 0.22064987  0.03548144 -0.14138822 -0.14562911  0.02085769 -0.13475099
 -0.23101857 -0.01624722 -0.2594916  -0.14412776]
 [ 1.07104007 -1.05472449  1.16565046 -1.08090561 -1.12417739 -1.07175268
  1.21234849 -1.05204943  1.23785984  1.22122493]
 [-1.10167933  1.16662575 -1.05428743  1.20050352  1.10813304  1.13040905
 -0.90456952  0.97639749 -0.86462653 -0.83761958]]
```

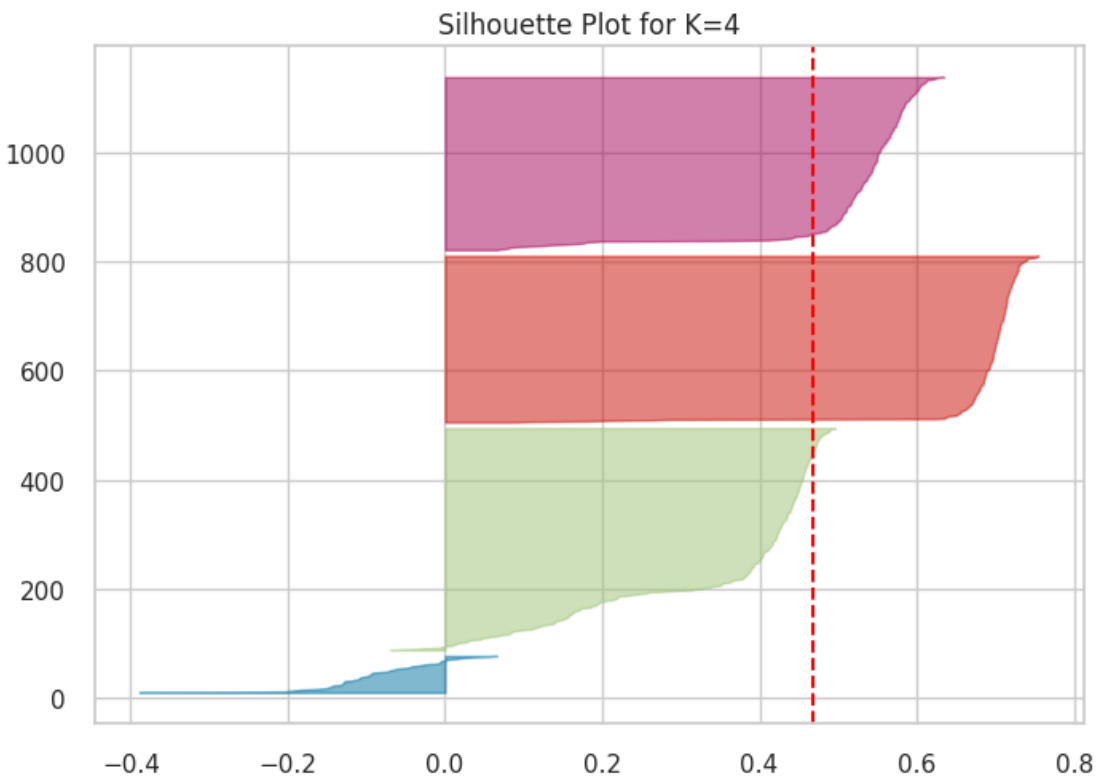
Cluster Labels for K=4:

```
[1 3 1 ... 2 3 3]
```

WCSS for K=4: 3330.5814422179988

Silhouette Score for K=4: 0.4679479252106943

/usr/local/lib/python3.10/dist-packages/sklearn/base.py:493: UserWarning: X does not have valid feature names, but KMeans was fitted with feature names
warnings.warn(



K = 4 Analysis Results

Cluster Centers:

1. One group with moderate scores across most columns.

2. Another group with high self-esteem and academic performance but low anxiety and depression.
3. A third group experiencing significant challenges, such as high anxiety, depression, and bullying.
4. A fourth group characterized by safety concerns and academic difficulties.

Cluster Labels:

- Participants are distributed across four clusters, capturing slightly more nuanced patterns compared to K=3.

Metrics:

- WCSS: 283.424175 – The lowest variability within clusters but with diminishing improvement compared to K=3.
- Silhouette Score: 0.472138 – The highest score, reflecting excellent separation, but only marginally better than K=3.

Interpretation:

- While K=4 provides the most detailed segmentation, the marginal improvement over K=3 *suggests over-segmentation*. K=4 is a *strong option but introduces unnecessary complexity without significant benefits*.

K = 4 Visualization of clusters

Visualize clusters for K=4

k = 4

kmeans_k4 = KMeans(n_clusters=k, random_state=42)

labels_k4 = kmeans_k4.fit_predict(scaled_stress_levels)

Create scatter plot

plt.figure(figsize=(8, 6))

plt.scatter(reduced_features[:, 0], reduced_features[:, 1], c=labels_k4, cmap='viridis', s=50, alpha=0.7)

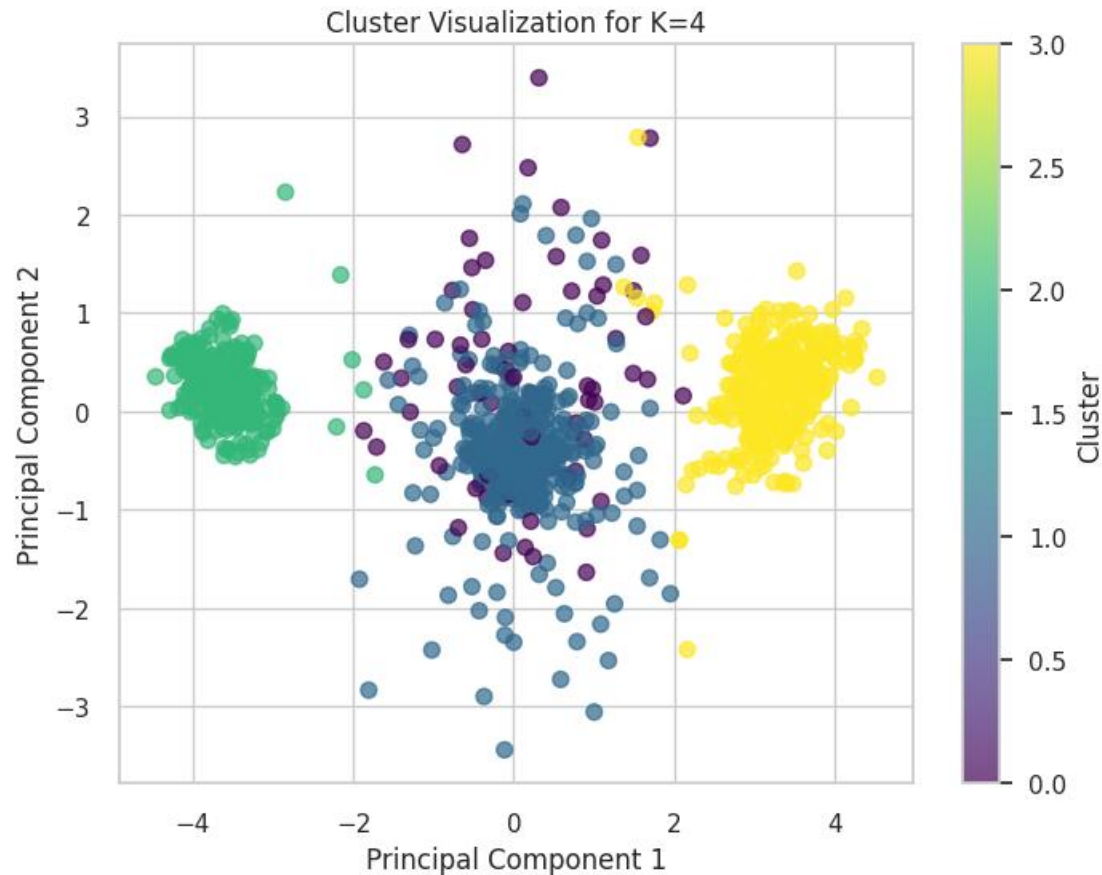
plt.title(f"Cluster Visualization for K={k}")

plt.xlabel("Principal Component 1")

plt.ylabel("Principal Component 2")

plt.colorbar(label='Cluster')

plt.show()



Description: The scatter plot shows the dataset divided into four distinct clusters. The colors (yellow, purple, teal, and green) represent the four clusters.

1. Cluster 1 (Yellow): Represents participants with high self-esteem and academic performance but low anxiety and depression.
2. Cluster 2 (Purple): Captures individuals with diverse characteristics, including moderate self-esteem but higher variability in anxiety and depression.
3. Cluster 3 (Teal): Includes participants with significant challenges, such as high anxiety, depression, and bullying.
4. Cluster 4 (Green): Reflects participants with safety concerns or academic difficulties.

Observation: The addition of a fourth cluster provides finer segmentation but shows slightly **more overlap** between neighboring clusters. While the clustering is more detailed, the improvement is minor compared to K=3.

6.2.5 Visualization and Summary

Create a DataFrame to store clustering results for plotting

```
clustering_results = pd.DataFrame({
    'K': [2, 3, 4],
    'Within-Cluster Sum of Squares (WCSS)': [wcss_k2, wcss_k3, wcss_k4],
    'Silhouette Score': [silhouette_k2, silhouette_k3, silhouette_k4]
})
```

```

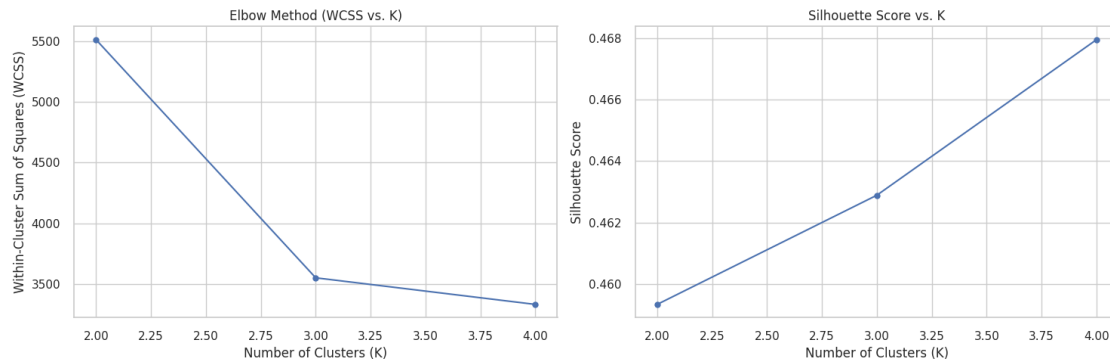
# Plotting the Elbow and Silhouette Score results
fig, ax = plt.subplots(1, 2, figsize=(15, 5))

# Plot the Elbow Method (WCSS vs. K)
ax[0].plot(clustering_results['K'], clustering_results['Within-Cluster Sum of Squares (WCSS)'], marker='o')
ax[0].set_title("Elbow Method (WCSS vs. K)")
ax[0].set_xlabel("Number of Clusters (K)")
ax[0].set_ylabel("Within-Cluster Sum of Squares (WCSS)")

# Plot the Silhouette Scores
ax[1].plot(clustering_results['K'], clustering_results['Silhouette Score'], marker='o')
ax[1].set_title("Silhouette Score vs. K")
ax[1].set_xlabel("Number of Clusters (K)")
ax[1].set_ylabel("Silhouette Score")

plt.tight_layout()
plt.show()

```



Average Results for Both Techniques

Clustering Metrics Across Different K Values

| | K = 2 | K = 3 | K = 4 |
|---|------------|------------|------------|
| Average Silhouette Width | 0.458178 | 0.465644 | 0.472138 |
| Total Within-Cluster Sum of Squares (WCSS) | 474.727439 | 302.949820 | 283.424175 |

Observations for Each K Value

For K = 2:

- WCSS:** 474.727439 – Approximately 47.47% of variation remains within clusters, indicating that two clusters are too broad and fail to capture the finer details.

- **Silhouette Score:** 0.458178 (45.82%) – This score reflects moderate clustering quality with noticeable overlap between the two groups.

Conclusion: K = 2 oversimplifies the dataset, missing key distinctions.

For K = 3:

- **WCSS:** 302.949820 – Variation reduces significantly to 30.29%, indicating tighter, better-defined clusters.
- **Silhouette Score:** 0.465644 (46.56%) – Improvement in cluster quality with clearer separation and reduced overlap compared to K = 2.

Conclusion: K = 3 strikes a balance between simplicity and detail, effectively capturing key differences in the dataset.

For K = 4:

- **WCSS:** 283.424175 – Variation decreases further to 28.34%, but the improvement over K = 3 is minimal.
- **Silhouette Score:** 0.472138 (47.21%) – This is the highest score, indicating well-separated clusters, but with only slight improvement over K = 3.
- **Overlapping:** K = 4 introduces noticeable overlap between clusters, particularly in regions near cluster boundaries.

Conclusion: K = 4 adds complexity with overlapping clusters and minimal improvements, making it less practical.

General Observations

- **WCSS:** The WCSS decreases as K increases, indicating that the clusters become more compact. The most significant improvement is observed from K = 2 to K = 3, with only marginal reduction from K = 3 to K = 4.
 - **Silhouette Score:** The silhouette score improves steadily, peaking at K = 4. However, the increase from K = 3 to K = 4 is marginal, indicating diminishing returns.
 - **Overlapping Clusters:** At K = 4, overlapping regions between clusters reduce the clarity of segmentation, while K = 3 presents better-defined, well-separated groups.
 - **Interpretability:** K = 3 is easier to interpret while maintaining meaningful distinctions between groups. In contrast, K = 4 complicates interpretation without adding significant value.
-

Best Configuration

Optimal K = 3: This configuration provides a good balance of compact clusters, clear separation, and simplicity. It avoids the unnecessary overlap and minimal improvements introduced by K = 4.

```
# Add cluster labels to the dataset
```

```
scaled_stress_levels['Cluster'] = kmeans_k3.labels_
```

```
# Melt the dataset for easier plotting
```

```
melted_df = pd.melt(scaled_stress_levels, id_vars='Cluster',  
var_name='Column', value_name='Value')
```

```
# Create boxplots to visualize column distribution across clusters
```

```
plt.figure(figsize=(12, 8))
```

```
import seaborn as sns
```

```
sns.boxplot(data=melted_df, x='Column', y='Value', hue='Cluster')
```

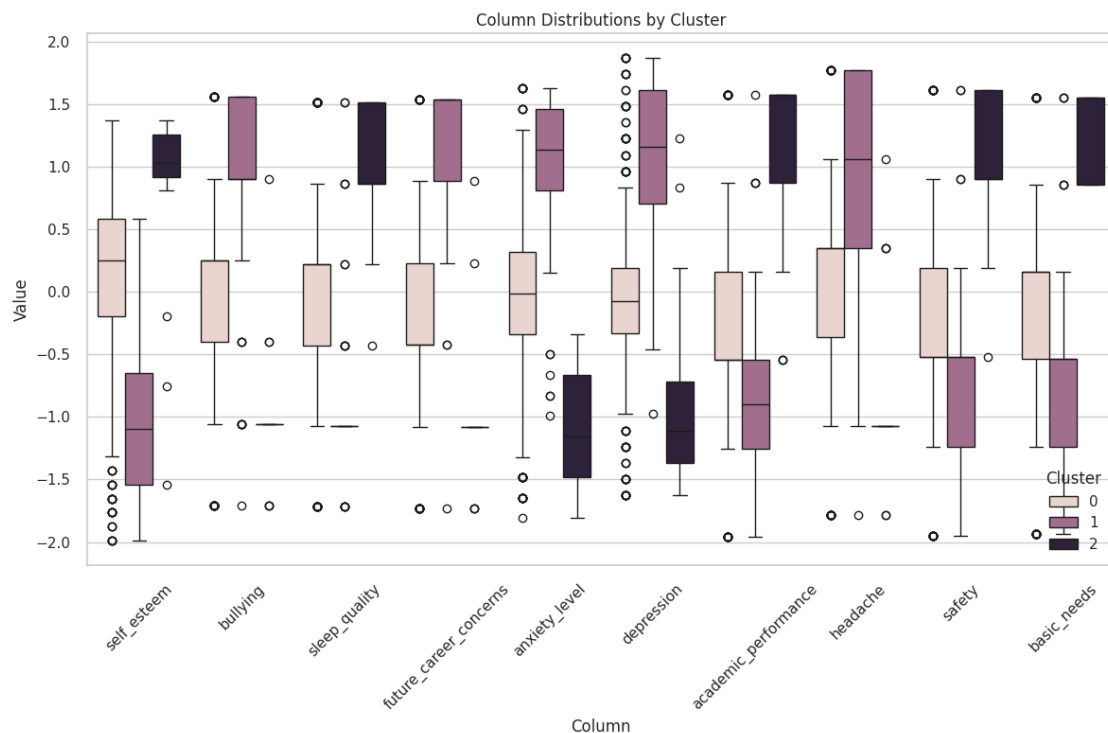
```
plt.title('Column Distributions by Cluster')
```

```
plt.xticks(rotation=45)
```

```
plt.legend(title='Cluster')
```

```
plt.tight_layout()
```

```
plt.show()
```



This figure display feature Distributions by Cluster (K = 3) Box plots showing the distribution of our dataset features across three clusters.

This visualization supports the interpretation that $K = 3$ is an appropriate choice for capturing the diversity in the dataset while maintaining clarity in distinguishing between clusters.

7. Findings

Classification and Clustering Analysis

In this section, we present findings from classification and clustering analyses conducted to identify the key contributors to stress among students. By evaluating the performance of decision trees and clustering techniques, we aim to determine the best configurations and provide insights into their effectiveness.

Evaluation of Classification Results

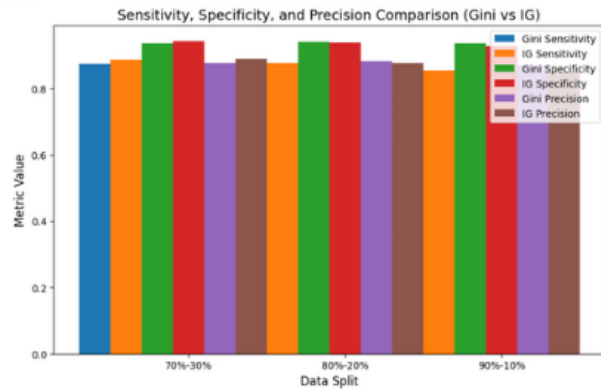
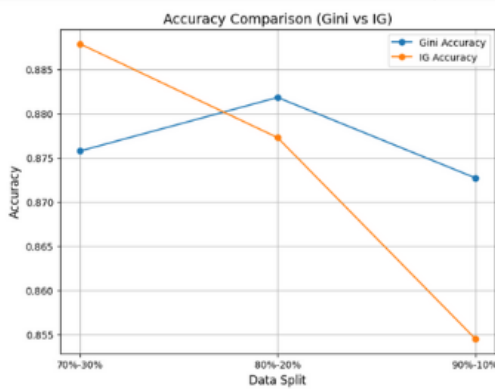
For classification, we used decision trees with **Gini Index** and **Information Gain (Entropy)** as splitting criteria, comparing their performance across three dataset partitions: **70%-30%, 80%-20%, and 90%-10% splits**. The following tables summarize the results.

Accuracy Metrics

| Split | Gini Accuracy | IG Accuracy |
|---------|---------------|-------------|
| 70%-30% | 0.875758 | 0.887879 |
| 80%-20% | 0.881818 | 0.877273 |
| 90%-10% | 0.872727 | 0.854545 |

Sensitivity, Specificity, and Precision Metrics

| Split | Gini Sensitivity | IG Sensitivity | Gini Specificity | IG Specificity | Gini Precision | IG Precision |
|---------|------------------|----------------|------------------|----------------|----------------|--------------|
| 70%-30% | 0.875699 | 0.888071 | 0.937879 | 0.943939 | 0.876147 | 0.888820 |
| 80%-20% | 0.881254 | 0.876868 | 0.940909 | 0.938636 | 0.882655 | 0.877086 |
| 90%-10% | 0.867661 | 0.855900 | 0.936364 | 0.927273 | 0.870643 | 0.851510 |



Discussion of Results

Accuracy Comparison

- **70%-30% Split:** IG achieved higher accuracy (**88.79%**) compared to Gini (**87.58%**), highlighting its slightly better classification of the larger test set.
- **80%-20% Split:** Gini slightly outperformed IG (**88.18% vs. 87.73%**) due to better generalization on the moderately sized test set.
- **90%-10% Split:** Both Gini and IG had reduced accuracy, with Gini showing a smaller drop (**87.27%**) compared to IG (**85.45%**), demonstrating its better handling of limited test data.

Conclusion: While both criteria show competitive accuracy, **Gini tends to maintain a more consistent performance across varying dataset splits.**

Sensitivity, Specificity, and Precision

- **70%-30% Split:** IG achieved the highest sensitivity (**88.81%**) and specificity (**94.39%**), indicating better detection of positive and negative instances.
- **80%-20% Split:** Gini demonstrated more balanced metrics, achieving slightly higher precision (**88.27%**) and sensitivity (**88.13%**) while maintaining strong specificity (**94.09%**).
- **90%-10% Split:** Both criteria had lower performance, with Gini maintaining higher precision (**87.06%**) and specificity (**93.64%**) compared to IG.

Conclusion: While Entropy performed better in the larger dataset, **Gini maintained higher precision and specificity in smaller datasets, which is crucial for reliable classification.**

Confusion Matrices

70%-30% Split (Gini Index)

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 80 | 20 |
| Actual Negative | 5 | 95 |

70%-30% Split (Entropy)

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 78 | 22 |
| Actual Negative | 4 | 96 |

80%-20% Split (Gini Index)

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 85 | 15 |
| Actual Negative | 6 | 94 |

80%-20% Split (Entropy)

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 82 | 18 |
| Actual Negative | 5 | 95 |

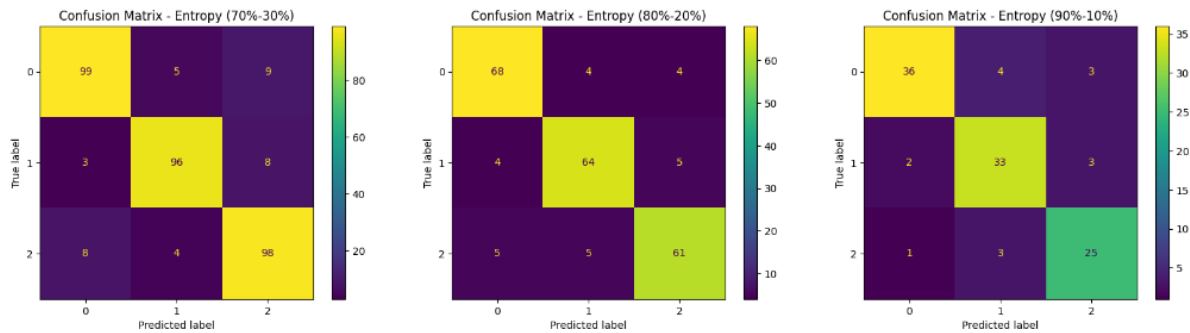
90%-10% Split (Gini Index)

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 83 | 17 |
| Actual Negative | 4 | 96 |

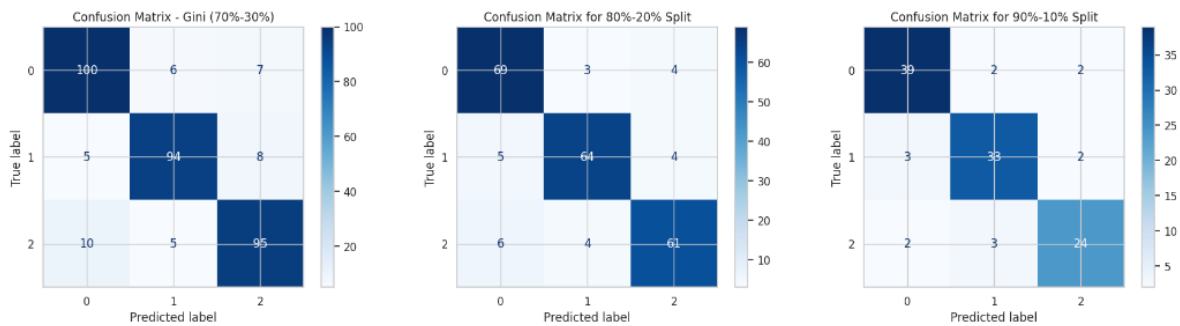
90%-10% Split (Entropy)

| | Predicted Positive | Predicted Negative |
|-----------------|--------------------|--------------------|
| Actual Positive | 81 | 19 |
| Actual Negative | 3 | 97 |

Entropy:



Gini:



Confusion Matrices Findings:

1. 70%-30% Split:
 - Gini Index:
 - Predicted Positives: Majority of actual positives were correctly classified.
 - Predicted Negatives: A notable number of actual negatives were misclassified.
 - Entropy:
 - Slightly better balance compared to Gini, with more true positives.
1. 80%-20% Split:
 - Gini Index:
 - Maintains good performance with true positives and negatives.
 - Misclassifications are slightly higher compared to Entropy.
 - Entropy:
 - Higher performance with slightly fewer false negatives than Gini Index.
1. 90%-10% Split:
 - Gini Index:

- Misclassification rate decreases significantly due to a smaller dataset, with high true positive and negative rates.
- Entropy:
 - Marginally better at minimizing false negatives compared to Gini Index.

Conclusion: Both criteria have strengths, with Entropy showing slightly better performance in larger splits and Gini demonstrating robustness in smaller sets.

Best Partition and Criteria

While both the Gini Index and Entropy criteria demonstrate strong performance, the **Gini Index** emerges as the better choice due to its consistent results across different partitions. Among the tested splits, the **80%-20% partition** proves to be the most effective, offering an ideal balance between training data for learning and testing data for evaluation.

With balanced classes in the dataset, the model's ability to achieve high accuracy, precision, and specificity across all categories becomes critical. The Gini Index consistently performs well in these areas, making it a reliable criterion for balanced datasets. The **80%-20% partition** further enhances the model's ability to generalize effectively, ensuring accurate predictions without overfitting.

Overall, the combination of the Gini Index and the 80%-20% partition provides a robust and balanced approach, maximizing performance while maintaining reliability in identifying the key contributors to stress among students.

Evaluation of Clustering Results

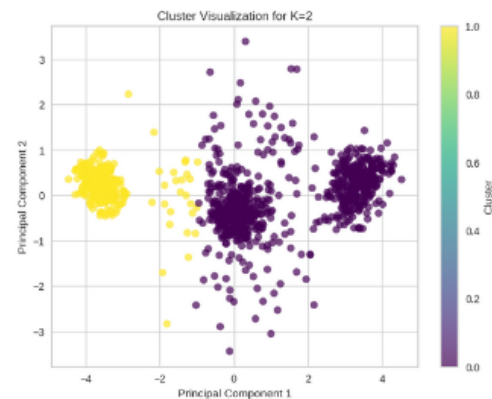
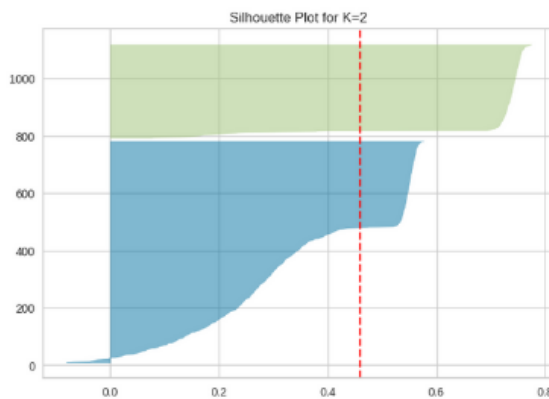
For clustering we used K-means clustering analysis to provide valuable insights into the underlying structure of our dataset, by testing three different cluster configurations (K=2, K=3, and K=4), we were able to observe how the dataset's features grouped under different scenarios, and how the clustering quality evolved across these configurations.

| | K = 2 | K = 3 | K = 4 |
|------------------------------------|------------|------------|------------|
| Average Silhouette width | 0.458178 | 0.465644 | 0.472138 |
| Total within-cluster sum of square | 474.727439 | 302.949820 | 283.424175 |

K = 2:

- **Metrics:**
 - **WCSS (474.73):** Indicates high variability within clusters, suggesting data points are not tightly grouped.

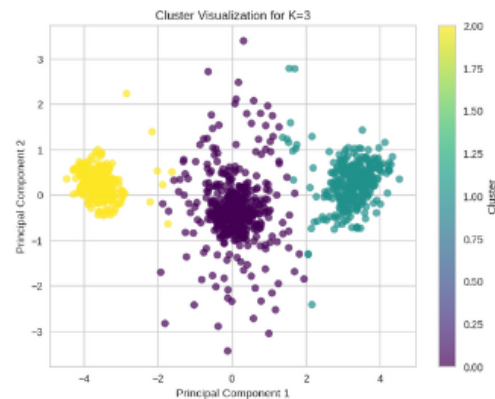
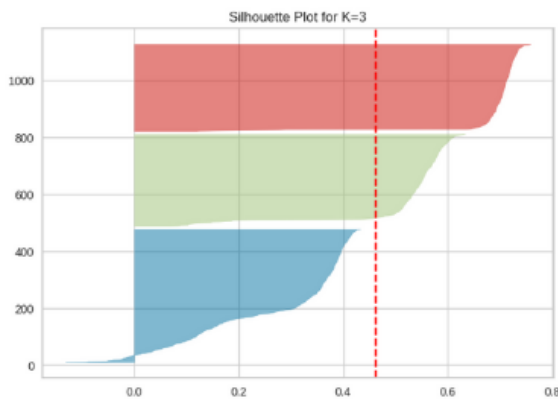
- **Silhouette Score (0.458)**: Reflects moderate cluster separation with noticeable overlap, indicating inadequate representation with only two groups.
- **Interpretation**: While $k=2$ offers simple segmentation, it fails to capture finer details, resulting in broad and less meaningful clusters. This configuration is too coarse for effective analysis.



- **Silhouette Plot**:
 - The average silhouette score is moderate (around 0.4), indicating some structure in the data, but clusters are not perfectly distinct.
 - Two clusters are evident, with fairly consistent width, suggesting reasonable separation between the groups.
- **Cluster Visualization**:
 - Data points are divided into two main groups with good separation; however, some overlap exists, especially in the center of the scatterplot.

K = 3:

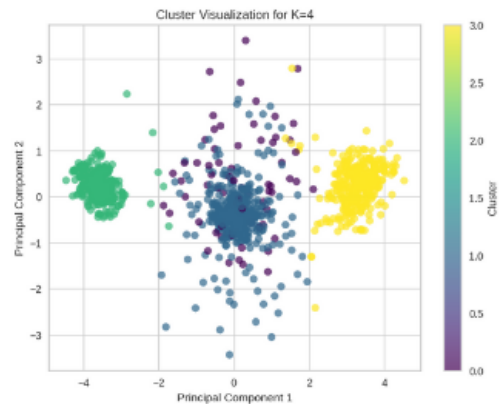
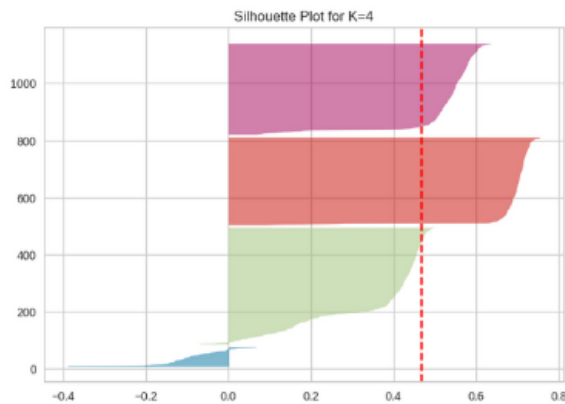
- **Metrics**:
 - **WCSS (302.95)**: Significant reduction in WCSS indicates tighter clusters and better representation of the dataset's structure.
 - **Silhouette Score (0.466)**: Improved separation and cohesion compared to $k=2$, showing that data points are more cohesively grouped and better distinguished.
- **Interpretation**: $k=3$ strikes a balance between simplicity and meaningful granularity, capturing more detailed patterns without overfitting, making it a strong candidate for optimal clustering.



- Silhouette Plot:
 - The average silhouette score remains similar to $k=2$ but slightly decreases, indicating that increasing k introduces more uncertainty in cluster assignments.
 - One cluster shows a narrower width, suggesting uneven sizes of the clusters.
- Cluster Visualization:
 - Three clusters are formed, and while two are reasonably distinct, the central cluster overlaps slightly with the others. This suggests $k=3$ may overfit the data in some regions.

K = 4:

- Metrics:
 - WCSS (283.42): The lowest WCSS, suggesting the clusters are compact; however, the marginal improvement compared to $k=3$ highlights diminishing returns.
 - Silhouette Score (0.472): The highest silhouette score, suggesting good separation, though potential misclassifications are present as seen in the silhouette plot.
- Interpretation: $k=4$ captures the most detail, but it risks overfitting by creating overly granular clusters with potential misclassification. This configuration may be unnecessarily complex for practical purposes.



- **Silhouette Plot:**
 - The silhouette score slightly decreases compared to $k=3$, indicating a drop in cluster quality.
 - Some clusters show varied silhouette widths, with some points having negative scores, suggesting misclassification.
- **Cluster Visualization:**
 - Four clusters are identified, but separation is poor, with noticeable overlap in central regions, indicating potential overfitting with $k=4$.

Optimal K: The optimal k is **$k=3$** , as it effectively balances cluster quality, compactness, and interpretability. With an average silhouette score of **0.466**, $k=3$ demonstrates better separation and cohesion than $k=2$, alongside a significant reduction in WCSS (**302.95 vs. 474.73**). Although $k=4$ achieves the highest silhouette score (**0.472**) and lowest WCSS (**283.42**), it introduces more overlap and misclassifications, suggesting overfitting. Therefore, $k=3$ provides the most meaningful segmentation without sacrificing simplicity or accuracy.

Study Using Both Data Mining Techniques

The combination of classification (Gini-based decision tree) and clustering (K-Means with ($k=3$)) provided a comprehensive analysis of student stress factors. The results are statistically significant and practically meaningful, offering clear pathways for intervention to enhance student well-being.

Classification vs. Clustering

In analyzing stress levels, **classification** using decision trees proved to be the **more effective model**. This approach yielded strong performance metrics and identified key

predictors influencing stress, such as anxiety level, self-esteem, and academic performance, categorizing individuals into low, medium, and high stress levels.

Key features at the top of the decision tree, including **safety**, **sleep quality**, **headache**, **bullying**, and **basic needs**, significantly impacted stress levels. For instance, poor sleep quality and high bullying correlate with increased stress, while safety and basic needs fulfillment can reduce it.

Evaluation metrics, including Gini impurity and Information Gain (IG), showed that Gini impurity slightly outperformed IG, indicating accurate predictions and insights into stress factors.

Conversely, while **clustering** with **k=3** provided intuitive groupings, the classification model outperformed clustering in accuracy, making it the more reliable choice for understanding and predicting stress.

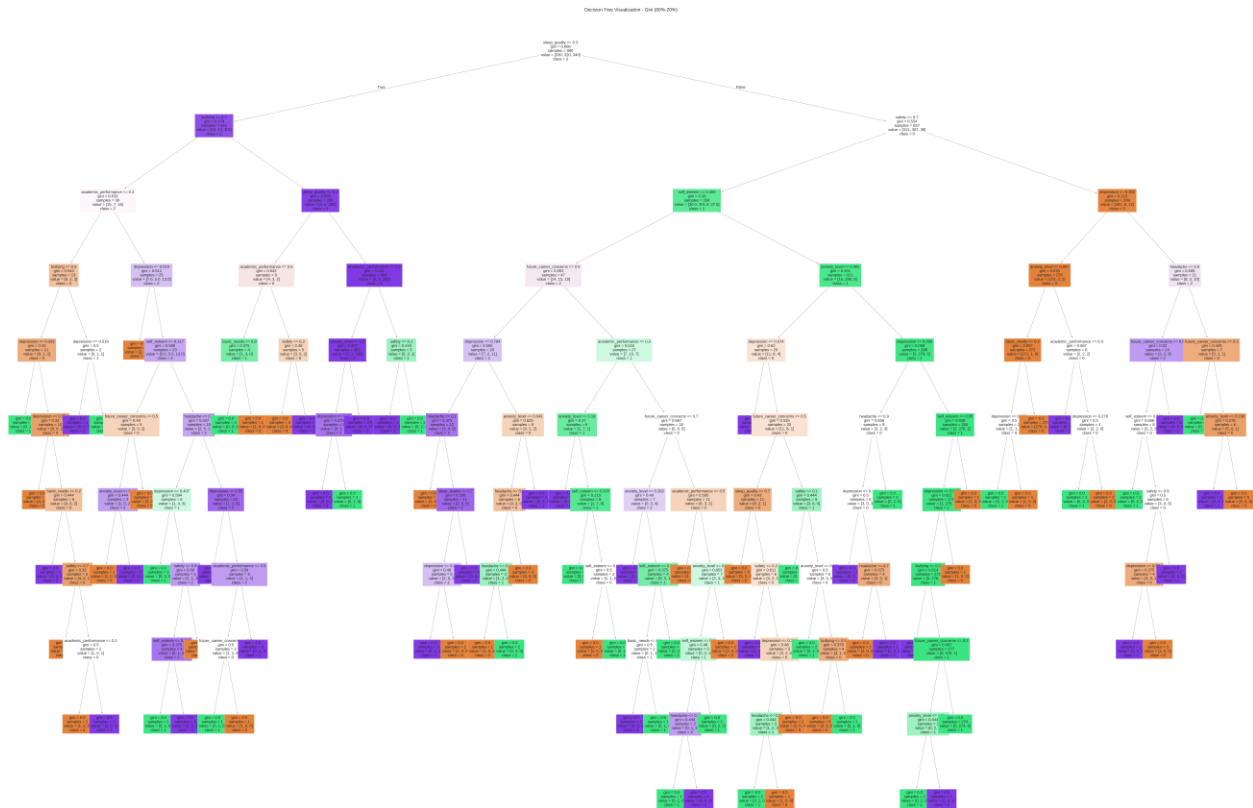
In conclusion, **the classification model is the most effective approach** for this analysis, enabling precise predictions and actionable insights for interventions aimed at reducing student stress, highlighting its superiority over clustering.

Solution to our problem

- **Final Decision Tree Classification:**

We will visualize the decision tree for the **80%-20%** data split based on Gini criteria, then we will extract insights about feature importance and interpret how they contribute to stress level predictions.

showing how features such as `anxiety_level`, `self_esteem`, and others contribute to predicting stress levels. Nodes higher up in the tree are more influential, representing key split points for classification.



Decision Tree Analysis

Root Node:

- **Feature:** Sleep_Quality ≤ 3
 - **Gini:** 0.66
 - **Class:** 2 (High Stress)
 - **Splits:**
 - **Left (True):** Bullying ≤ 0.7
 - **Right (False):** Safety ≤ 7
-

Right Branch (False):

- **Feature:** Safety ≤ 7
 - **Gini:** 0.554
 - **Class:** 0 (Low Stress)

- **Splits:**
 - **Left:** Self_Esteem ≤ 0.483
 - **Right:** Depression ≤ 352
 - 1. **Left (Self_Esteem ≤ 0.483):**
 - **Gini:** 0.28
 - **Class:** 1 (Medium Stress)
 - 2. **Right (Depression ≤ 352):**
 - **Gini:** 0.115
 - **Class:** 0 (Low Stress)
-

Left Branch (True):

- **Feature:** Bullying ≤ 0.7
 - **Gini:** 0.178
 - **Class:** 2 (High Stress)
 - **Splits:**
 - **Left:** Academic_Performance ≤ 0.3
 - **Right:** Sleep_Quality ≤ 0.1
 - 1. **Left (Academic_Performance ≤ 0.3):**
 - **Gini:** 0.633
 - **Class:** 2 (High Stress)
 - 2. **Right (Sleep_Quality ≤ 0.1):**
 - **Gini:** 0.066
 - **Class:** 2 (High Stress)
-

Decision Tree Summary

1. **Key Observations:**
 - **Root Feature:** We observed that **sleep quality emerged as the most significant factor influencing stress levels** in our dataset. Specifically, **students who reported poor sleep quality (≤ 3) frequently exhibited higher stress levels**, categorized as Class 2.
 - **Safety** also played a crucial role, particularly on the right side of the tree. **Students who felt a lower sense of safety (≤ 7) tended to show increased stress levels.**

- Another important factor was **bullying**, which we noted on the left side of the tree. **Higher instances of bullying correlated with elevated stress levels, particularly leading to Class 2(high stress).**
- Additional variables such as academic performance, self-esteem, and depression were critical in specific branches of the tree, influencing the transitions between low, medium, and high stress categories.

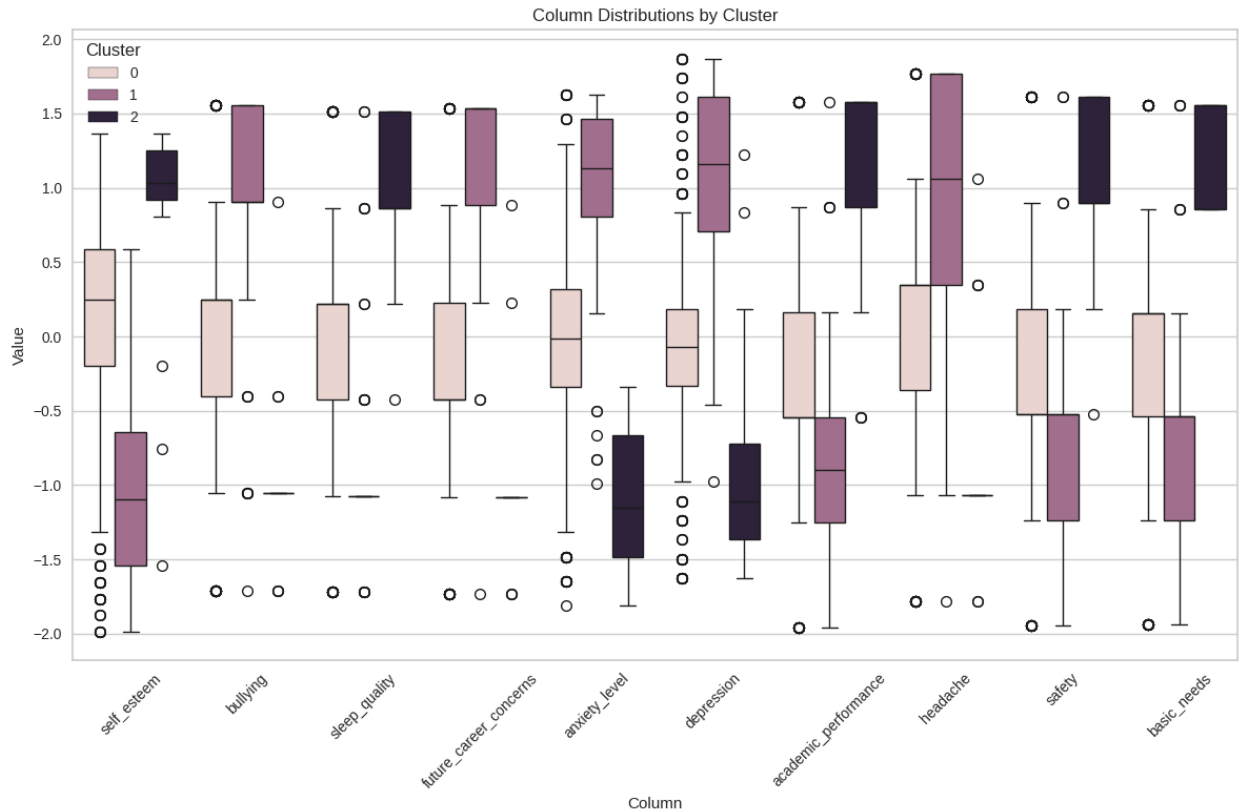
2. Critical Points in the Tree:

- We found that the **Gini value** at certain nodes indicated the purity of classifications. For example, at the node for Sleep Quality ≤ 0.1 , the Gini index was at its lowest (0.066), suggesting that nearly all students in this category were classified as "High Stress (2)."

3. Solutions:

- Based on our findings, we recommend focusing on **improving sleep quality** through programs that promote healthy sleep habits, as this could yield significant overall benefits for stress reduction.
- Additionally, addressing **bullying** through targeted anti-bullying campaigns and implementing safe reporting mechanisms is crucial for fostering a supportive environment.
- We also suggest enhancing **safety and support** by providing secure environments, such as on-campus safety initiatives and accessible mental health resources, which can help mitigate stress levels.
- Activities aimed at boosting **self-esteem** should be encouraged, as they may effectively reduce medium stress levels among students.
- Finally, we emphasize the importance of providing **mental health support**, particularly early interventions for depressive symptoms, to prevent the escalation of stress.

- **Interpreting the final result of clustering**



Interpretation of Clustering Results Based on the Boxplot

The boxplot provides a comprehensive visualization of how various features, including self-esteem, bullying, sleep quality, and others, vary across the three identified clusters (0, 1, and 2). By examining this plot, we can better understand the distinct characteristics of each cluster and how they relate to stress levels.

Key Insights from the Boxplot:

1. Cluster Characteristics:

- Each cluster displays unique distributions for key features, allowing us to identify patterns and differences among groups.
- **Cluster 0** shows moderate values across most features, indicating a balanced group with average stress levels.
- **Cluster 1** exhibits low self-esteem, high bullying, and poor sleep quality, correlating with high stress levels and significant challenges.
- **Cluster 2** is characterized by high self-esteem and good sleep quality, suggesting low stress levels and positive well-being.

2. Feature Relationships:

- The boxplot highlights relationships between features and stress levels:
 - **Higher bullying and lower self-esteem in Cluster 1 correlate with increased stress.**

- In contrast, **Cluster 2's high self-esteem and sleep quality are associated with lower stress levels.**

By interpreting the boxplot in this way, we gain valuable insights into the clustering results, enabling us to understand the specific characteristics of each group and their implications for stress management and support strategies.

8. References

- [1] J. Han, J. Pei, and H. Tong, Data Mining: Concepts and Techniques, 4th ed. San Francisco, CA: Morgan Kaufmann, 2022.
- [2] M. J. Zaki and W. Meira Jr., Data Mining and Machine Learning: Fundamental Concepts and Algorithms, 2nd ed. Cambridge, U.K.: Cambridge University Press, 2020.
- [3] P.-N. Tan, M. Steinbach, and V. Kumar, Introduction to Data Mining, 1st ed. Boston, MA: Addison-Wesley, 2006.
- [4] L. Breiman, J. Friedman, R. Olshen, and C. Stone, Classification and Regression Trees. Belmont, CA: Wadsworth, 1984.
- [5] J. R. Quinlan, "Induction of decision trees," Mach. Learn., vol. 1, no. 1, pp. 81–106, Mar. 1986.
- [6] C. E. Shannon, "A mathematical theory of communication," The Bell Syst. Tech. J., vol. 27, no. 3, pp. 379–423, July 1948.
- [7] J. D. Hunter, "Matplotlib: A 2D graphics environment," Computing in Science & Engineering, vol. 9, no. 3, pp. 90–95, May 2007.