

Py Project

Introduction

This project is about exploring factors affecting liver disease.

Dataset Description

The dataset contains information related to hepatitis, with 155 instances and 20 attributes, including the class attribute. The data was donated by G.Gong from Carnegie-Mellon University via Bojan Cestnik of Jozef Stefan Institute.

Variable Types and Statistical Measures

Categorical Variables: Type of variables that represent categories or groups. These variables can take on a limited, fixed number of distinct values or labels, and there is no inherent order or numerical significance among these categories. Class (DIE, LIVE), SEX (male, female), STEROID (no, yes), ANTIVirALS (no, yes), FATIGUE (no, yes), MALAISE (no, yes), ANOREXIA (no, yes), LIVER BIG (no, yes), LIVER FIRM (no, yes), SPLEEN PALPABLE (no, yes), SPIDERS (no, yes), ASCITES (no, yes), VARICES (no, yes), HISTOLOGY (no, yes)

Ordinal Variables: Type of categorical variable that, in addition to having distinct categories, also have a meaningful order or ranking among them. However, the intervals between the categories are not necessarily uniform or measurable. AGE (10, 20, 30, 40, 50, 60, 70, 80)

Continuous Variables: Quantitative variables that can take on an infinite number of values within a given range. These variables are typically measured on a continuous scale and can include decimal values. BILIRUBIN, ALK PHOSPHATE, SGOT, ALBUMIN, PROTOME

Graphical Representations and Statistical Measures:

Categorical Variables: Bar charts to show the distribution of each category. Class distribution can be visualized using a bar chart.

Ordinal Variables: Histogram to show the distribution of ages.

Continuous Variables: Box plots to identify outliers. Histograms for a visual representation of continuous variable distributions.

Dataset:

```
In [6]: import pandas as pd
df = pd.read_csv("hepatitis.txt")

print(df.info())
print(df.describe())

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 154 entries, 0 to 153
Data columns (total 20 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Class           154 non-null    int64  
 1   AGE             153 non-null    object  
 2   SEX             154 non-null    int64  
 3   STEROID         154 non-null    object  
 4   ANTIVirALS     154 non-null    int64  
 5   FATIGUE        153 non-null    object  
 6   MALAISE         154 non-null    object  
 7   ANOREXIA        154 non-null    object  
 8   LIVER BIG       154 non-null    object  
 9   LIVER FIRM      154 non-null    object  
 10  SPLEEN PALPABLE 154 non-null    object  
 11  SPIDERS         154 non-null    object  
 12  ASCITES         154 non-null    object  
 13  VARICES         154 non-null    float64 
 14  BILIRUBIN       152 non-null    object  
 15  ALK PHOSPHATE  153 non-null    object  
 16  SGOT            153 non-null    object  
 17  ALBUMIN         153 non-null    object  
 18  PROTOME         154 non-null    object  
 19  HISTOLOGY       153 non-null    float64 
dtypes: float64(1), int64(18), object(18)
memory usage: 24.2+ KB
None
```

Data Cleaning:

Data cleaning plays a pivotal role in the data preparation phase as it entails the identification and rectification of errors or inconsistencies within a dataset. This essential process guarantees the accuracy, comprehensiveness, and analysis readiness of the data. Typical data cleaning activities involve addressing missing values, eliminating duplicates, rectifying data types, and getting rid of irrelevant or redundant information. The primary objective is to improve the dataset's quality, alleviate potential biases or inaccuracies, and establish a more dependable basis for meaningful analysis and interpretation.

```
In [7]: import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

column_names = ["Class", "AGE", "SEX", "STEROID", "ANTIVIRALS", "FATIGUE", "MALAISE", "ANOREXIA", "LIVER BIG", "LIVER FIRM", "SPLEEN PALPABLE", "SPIDERS", "ASCITES", "VARICES", "BILIRUBIN", "ALK PHOSPHATE", "SGOT", "ALBUMIN", "PROTIME", "HISTOLOGY"]

missing_values = ["BILIRUBIN": "?", "ALK PHOSPHATE": "?", "SGOT": "?", "ALBUMIN": "?", "PROTIME": "?"]

df = pd.read_csv("hepatitis.txt", names=column_names, na_values=missing_values)
df.dropna(inplace=True)
df["BILIRUBIN"] = df["BILIRUBIN"].astype(float)

# Save the cleaned dataset
df.to_csv("hepatitis_cleaned.csv", index=False)
print(df.info())
print(df.describe())

<class 'pandas.core.frame.DataFrame'>
Int64Index: 77 entries, 3 to 154
Data columns (total 20 columns):
 #   Column          Non-Null Count  Dtype  
--- 
 0   Class           77 non-null    int64  
 1   AGE             77 non-null    object  
 2   SEX             77 non-null    int64  
 3   STEROID         77 non-null    object  
 4   ANTIVirALS     77 non-null    int64  
 5   FATIGUE        77 non-null    object  
 6   MALAISE         77 non-null    object  
 7   ANOREXIA        77 non-null    object  
 8   LIVER BIG       77 non-null    object  
 9   LIVER FIRM      77 non-null    object  
 10  SPLEEN PALPABLE 77 non-null    object  
 11  SPIDERS         77 non-null    object  
 12  ASCITES         77 non-null    object  
 13  VARICES         77 non-null    float64 
 14  BILIRUBIN       77 non-null    object  
 15  ALK PHOSPHATE  77 non-null    object  
 16  SGOT            77 non-null    float64 
 17  ALBUMIN         77 non-null    object  
 18  PROTIME         77 non-null    object  
 19  HISTOLOGY       77 non-null    float64 
dtypes: float64(3), int64(3), object(14)
memory usage: 12.6+ KB
None
```

Data Analysis:

Hypothesis Testing

```
In [8]: import pandas as pd
from scipy.stats import ttest_ind, chis2_contingency

df = pd.read_csv("hepatitis_cleaned.csv")
selected_factors = ['AGE', 'SEX', 'ANTIVIRALS', 'BILIRUBIN', 'SGOT', 'HISTOLOGY']

class_1_data = df[df['Class'] == 1]
class_2_data = df[df['Class'] == 2]

continuous_columns = ['BILIRUBIN', 'SGOT']

for column in continuous_columns:
    t_statistic, p_value = ttest_ind(class_1_data[column].dropna(), class_2_data[column].dropna())
    print(f"t-statistic: {t_statistic}, p-value: {p_value}")
    if p_value < 0.05:
        print("Reject the null hypothesis. There is a significant difference.\n")
    else:
        print("Fail to reject the null hypothesis. There is no significant difference.\n")

categorical_columns = ['SEX', 'STEROID', 'ANTIVIRALS', 'FATIGUE', 'MALAISE', 'LIVER BIG', 'LIVER FIRM', 'SPLEEN PALPABLE', 'SPIDERS', 'ASCITES', 'VARICES', 'HISTOLOGY']

# Perform chi-square tests for categorical variables
for column in categorical_columns:
    contingency_table = pd.crosstab(class_1_data[column], df[column])
    chi2_stat, p_value = chis2_contingency(contingency_table)
    print(f"Test for {column}:")
    print(f"\tChi-square statistic: ({chi2_stat})")
    print(f"\tP-value: {p_value}")
    if p_value < 0.05:
        print("Reject the null hypothesis. There is a significant difference.\n")
    else:
        print("Fail to reject the null hypothesis. There is no significant difference.\n")

Test for BILIRUBIN:
t-statistic: 1.0980899999999999
p-value: 0.29849459857684
Fail to reject the null hypothesis. There is no significant difference.

Test for SGOT:
t-statistic: 0.336928691254193
p-value: 0.7484720912320856
Fail to reject the null hypothesis. There is no significant difference.

Test for SEX:
Chi-square statistic: 1.17739217810863
p-value: 0.318210921844769
Fail to reject the null hypothesis. There is no significant difference.

Test for STEROID:
Chi-square statistic: 1.0565076164874556
p-value: 0.384013970929693575
Fail to reject the null hypothesis. There is no significant difference.

Test for ANTIVirALS:
Chi-square statistic: 2.7694394663480884
p-value: 0.09607980135369354
Fail to reject the null hypothesis. There is no significant difference.

Test for FATIGUE:
Chi-square statistic: 8.521843758430649
p-value: 0.01410628935254913
Fail to reject the null hypothesis. There is no significant difference.

Test for MALAISE:
Chi-square statistic: 8.2794354838799676
p-value: 0.0123208567684
Fail to reject the null hypothesis. There is no significant difference.

Test for ANOREXIA:
Chi-square statistic: 8.2794354838799676
p-value: 0.0123208567684
Fail to reject the null hypothesis. There is no significant difference.

Test for LIVER BIG:
Chi-square statistic: 10.6972024301075269
p-value: 0.004372093980572119
Reject the null hypothesis. There is a significant difference.

Test for SPLEEN PALPABLE:
Chi-square statistic: 7.459625390218524
p-value: 0.02399733023102146
Reject the null hypothesis. There is a significant difference.

Test for SPIDERS:
Chi-square statistic: 10.62284439034287
p-value: 0.013956313877031165
Reject the null hypothesis. There is a significant difference.

Test for VARICES:
Chi-square statistic: 13.467731147131726
p-value: 0.000955653779012229
Reject the null hypothesis. There is a significant difference.

Test for HISTOLOGY:
Chi-square statistic: 16.09584677419355
p-value: 3.92894508975144e-05
Reject the null hypothesis. There is a significant difference.
```

In summary, the tests indicate significant differences in several variables(BILIRUBIN, MALAISE, LIVER BIG, LIVER FIRM, SPLEEN PALPABLE, SPIDERS, ASCITES, VARICES, HISTOLOGY), suggesting meaningful associations or patterns in the dataset for those specific factors. Other factors (SEX, STEROID, ANTIVirALS, FATIGUE, ANOREXIA) do not have any significant difference.

Correlation Analysis

```
In [10]: import pandas as pd
import seaborn as sns
import matplotlib.pyplot as plt

df = pd.read_csv("hepatitis_cleaned.csv")
selected_columns = ['BILIRUBIN', 'MALAISE', 'HISTOLOGY']
selected_data = df[selected_columns]
correlation_matrix = selected_data.corr()
print(correlation_matrix)
```

```
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm', fmt=".2f", linewidths=.5)
plt.title('Correlation Matrix')
plt.show()
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.00    -0.03   0.27
MALAISE    -0.03    1.00    0.13
HISTOLOGY  0.27    0.13   1.00
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
Correlation Matrix
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.000000
```

```
BILIRUBIN  MALAISE  HISTOLOGY
BILIRUBIN  1.000000 -0.034376  0.268529
MALAISE    -0.034376  1.000000  0.132362
HISTOLOGY  0.268529  0.132362  1.
```