# Relationship of Different Biochemical and Clinical Markers with Liver Cirrhosis Stages

Rahul Somabhai Chaudhary
Data Science
University of North Texas
Denton, TX, USA
RahulSomabhaiChaudhary@my.unt.edu

**Introduction**

Liver Cirrhosis is a chronic, progressive disease spread across the liver, making it failure to work properly over prolonged period of time. This condition, often resulting from alcohol abuse, hepatitis infection or fatty liver disease significantly impairing the liver function.

This study provides the aims to understands the effect of different components affecting the Liver resulting in Liver Cirrhosis such as Age, Cholesterol and many more by comparing data of 25000 individuals. In this study, the advance tools such as Tableau and languages like Python has been used.

Through this study, I am to identify the important trends and relations that offer insightful knowledge into complex interplay of factors affecting liver cirrhosis progression.

## I. Dataset

### A. Collecting Dataset

Dataset has been collected from Kaggle. It contains the information of 25000 individuals, with their data such as Age, Sex, Ascites, Hepatomegaly etc., this data provides a key information on understanding the level of liver cirrhosis of that individual.

Kaggel Dataset link: Liver Cirrhosis Stage Classification

### B. Dataset Information

Dataset contains 25000 rows and 19 columns, each row represents different person while each column represents different characteristics affecting liver. 19 columns contains as follows:

1. Sex: Gender of the patients (Male/Female)

2. Age: Age of the patient in years

3. Ascites: It is condition of too much fluid in the abdomen. Here it represents the presence of Ascites or not (Y/N)

4. Hepatomegaly: Condition of enlarged liver, here it is represented by yes or no(Y/N)

5. Spiders: Presence of spider angiomas(Y/N)

6. Edema: It occurs when the fluid accumulates in the leg[1]. presence of edema N (no edema and no diuretic therapy for edema), S (edema present without diuretics, or edema resolved by diuretics), or Y (edema despite diuretic therapy)[2]

7. Stages: histologic stage of disease ( 1, 2, or 3 )[2]

8. Status: status of the patient C (censored), CL (censored due to liver tx), or D (death)[2]

9. Drug: Type of drug used for treatment(D-penicillamine or placebo)

10. N_Days: Number of days between registration and the earlier of death, transplantation, or study analysis time in 1986[2]

11. Bilirubin: Serum bilirubin in[mg/dl]

12. Cholesterol: Blood cholesterol level of patient in [mg/dl]

13. Albumin: Albumin level in [gm/dl]

14. Copper: Urine copper level in [ug/day]

15. Alk_Phos: Alkaline phosphate level in [U/liter]

16. SGOT: Aspartate aminotransferase [U/ml]

17. Triglicerides: Triglycerides in [mg/dl]

18. Platelets: Platelets count per cubic [ml/1000]

19. Prothrombin: Prothrombin time in seconds

### C. Tools used for the dataset

- Python: For data cleaning, Exploratory Data Analysis and visualization
- Tableau: For visualization

## II. Data Cleaning And exploratory data analysis

To continue with the dataset and start our hypothesis, we need to do some cleaning such as removing null or missing values from our dataset and perform EDA to get knowledge of our dataset.

### A. Data Cleaning

In this portion, we will perform the cleaning of data using python, we will start with

1. Importing necessary modules and the dataset



2. Checking for any null values

```
#Checkin the null value
df.isnull().sum()

N_Days          0
Status          0
Drug            0
Sex             0
Ascites         0
Hepatomegaly    0
Spiders         0
Edema           0
Bilirubin       0
Cholesterol     0
Albumin         0
Copper          0
Alk_Phos        0
SGOT            0
Tryglicerides   0
Platelets       0
Prothrombin     0
Stage           0
Age             0
dtype: int64
```

As the data doesn't contains any null values we can move forward and perform EDA

## B. Performing Exploratory Data Analysis

Now to understand the dataset better and to find the relationship among the attributes we will perform EDA.



Here by function head and tail we can print some of the first and last rows of dataset.



The function shape provides the size of dataset while the description provides the statistical calculation of dataset.



Method info provides information such as type and non-null values of columns.



Unique provides the number of unique values in columns.

## C. Finding the relations among attributes

We will implement the relationship among the elements using python,

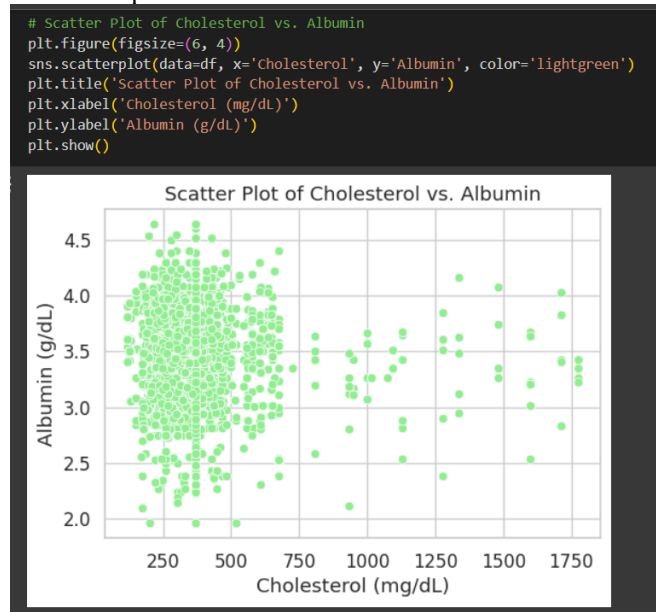1. What is the ratio of male to female in dataset?



The ratio of male to female is 11.5:88.5 respectively.

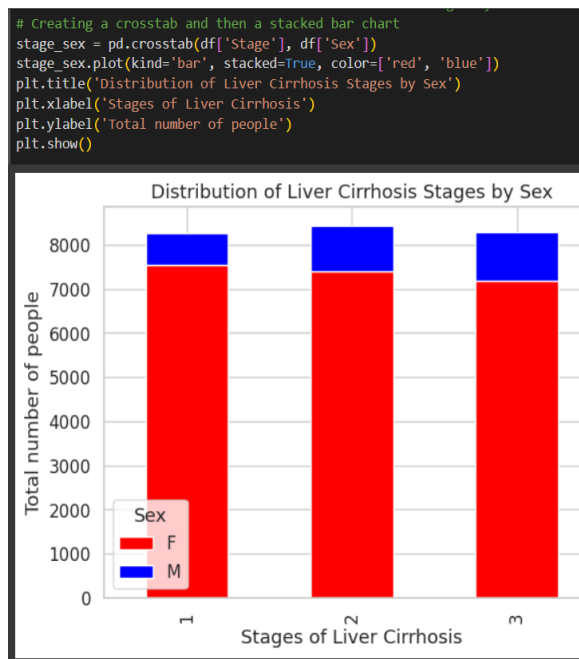2. What is the distribution of drugs?



From the figure, we can say that the 15,827 have given placebo while the rest of the 9,173 have been treated with D-penicillamine.
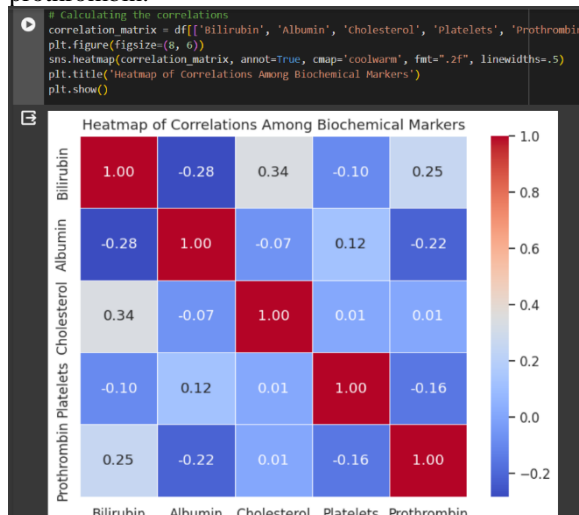
3. Relationship between cholesterol and albumin.



4. Relationship between male and female individual regarding their liver cirrhosis stages.

```
# Creating a crosstab and then a stacked bar chart
stage_sex = pd.crosstab(df['Stage'], df['Sex'])
stage_sex.plot(kind='bar', stacked=True, color=['red', 'blue'])
plt.title('Distribution of Liver Cirrhosis Stages by Sex')
plt.xlabel('Stages of Liver Cirrhosis')
plt.ylabel('Total number of people')
plt.show()
```



5. Correlation matrix showing relationship between bilirubin, albumin, cholesterol, platelets, prothrombin.



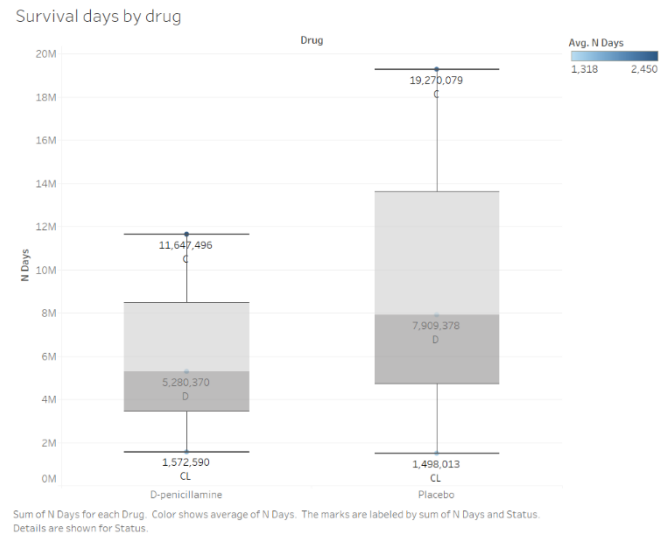### III. HYPOTHESIS TESTING USING TABLEAU

From the above visualizations, we can conclude that the attributes are related to each other, so with the use of tableau, we will now focusing on proving the hypothesis. Below are the mentioned hypothesis that we are going to focus upon.

1. Patients receiving different drugs have different survival rates.

2. The presence of ascites, hepatomegaly or edema is more frequent at higher stages of liver cirrhosis.

3. Cholesterol levels are not significantly different between patients with and without edema.

4. There is no significant relationship between age and prothrombin time.

A. *Patients receiving different drugs have different survival rates.*
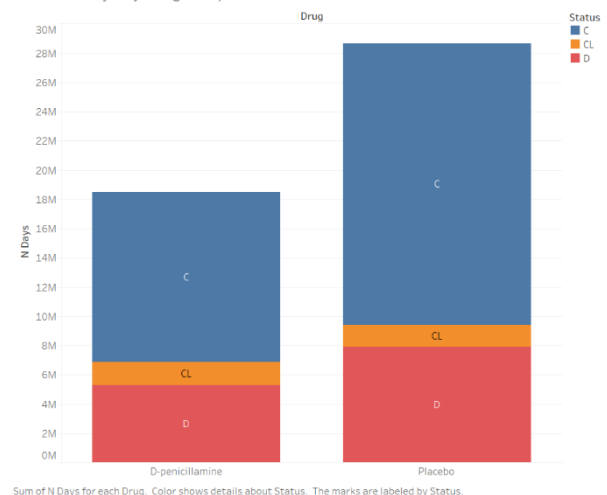
Visualization 1:
Here, I have used the box-and-whisker plot for visualization. Box plot is the way of visualizing data based on five-number summary, first quartile (Q1), median (Q2), third quartile (Q3), minimum and maximum.



Result: From this plot we can easily estimate that the difference in days of recovery from placebo is way higher than those from D-penicillamine, hence our hypothesis can be valid from this visualization.
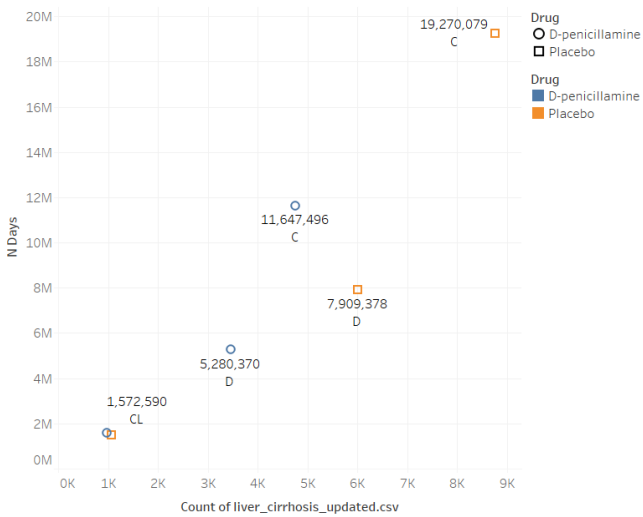
Visualization 2:



Here, the visualization used is a bar chart.
Bar chart uses the length or height of bars to represent the quantity of data [3].

Result: From the bars we can conclude that, for all the status whether it is C, CL, or D the placebo has been used more as a recovery drug than those of D-penicillamine, proving the validation of our hypothesis.

Visualization 3:

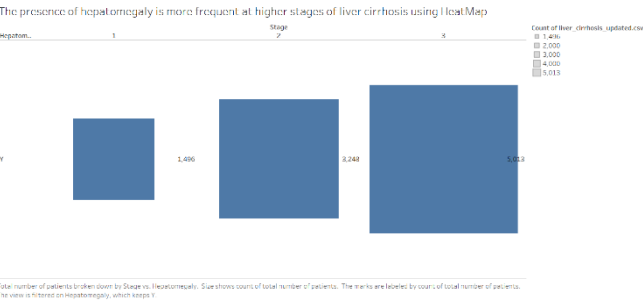## Distribution of Survival Days by Drug Type and Patient Status



Total number of patients vs. sum of N Days. Color shows details about Status. Shape shows details about Drug. The marks are labeled by sum of N Days and Status.

The visualization used is called scatter plot.

Scatter plot uses dots to represent the data. The dots on the cartesian plain representing data are placed on x and y axis.[4]

Result: From the visualization, we can conclude that the placebo has higher number of survival day respect to the patient count compared to those of D-Penicillamine. Hence the hypothesis is valid.

*B. The presence of ascites, hepatomegaly, or edema is more frequent at higher stages of liver cirrhosis.*

Visualization 1:



The presence of hepatomegaly is more frequent at higher stages of liver cirrhosis using Heat Map

Total number of patients broken down by Stage vs. Hepatomegaly. Size shows count of total number of patients. The marks are labeled by count of total number of patients. The view is filtered on Hepatomegaly, which keeps Y.
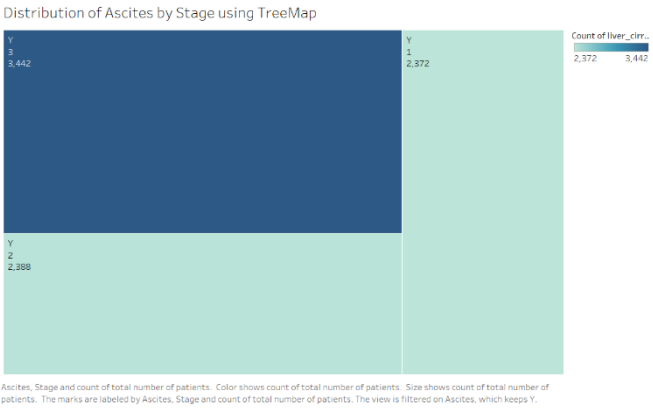
Visualization used to prove this hypothesis is called heat map. A graphical representation of data in the form of color coding to represents the different values of that dataset.[5]

Result: From the visualization, after filtering the data by the patient having hepatomegaly, we can conclude that as the liver cirrhosis stages increases, the number of patient with hepatomegaly increases as well, making our hypothesis valid.
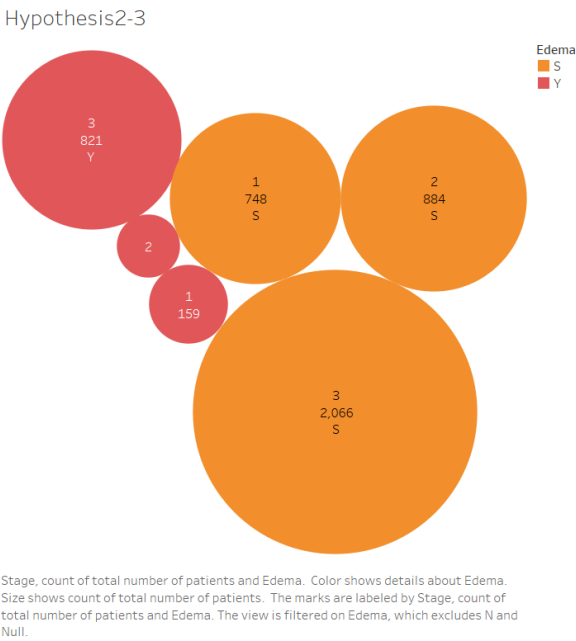
Visualization 2:

The visualization used is tree map for validating the hypothesis.

Tree map represents the data into hierarchical and rectangular form, each rectangle representing the data and size of rectangle is based on size of data.



Distribution of Ascites by Stage using TreeMap

Ascites, Stage and count of total number of patients. Color shows count of total number of patients. Size shows count of total number of patients. The marks are labeled by Ascites, Stage and count of total number of patients. The view is filtered on Ascites, which keeps Y.

Result: From the visualization, after sorting the patients with the ascites as present, we can see that the number of patients are higher in higher stages of liver cirrhosis and the number decreases as the stages of liver cirrhosis decreases, making our hypothesis valid.

Visualization 3:



Hypothesis2-3

Stage, count of total number of patients and Edema. Color shows details about Edema. Size shows count of total number of patients. The marks are labeled by Stage, count of total number of patients and Edema. The view is filtered on Edema, which excludes N and Null.

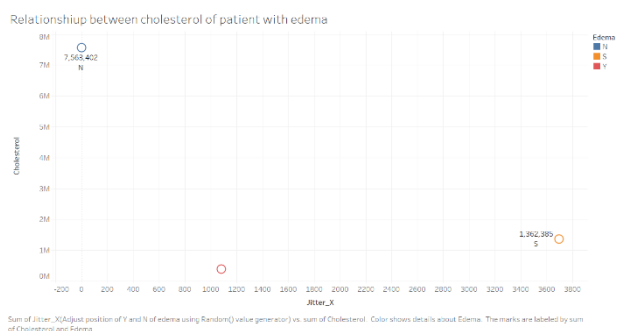Here, packed bubble is used for visualization.

In packed bubble chart, data are displayed in the form of clusters of circles, with bigger the data represented by the larger circle and vice versa.[6]

Result: After sorting the edema level to Y and S, we can conclude that the rise in number of edema patient is higher in patients with higher stages of liver cirrhosis, concluding that the hypothesis is valid.

*C. Serum cholesterol levels are not significantly different between patients with and without edema.*
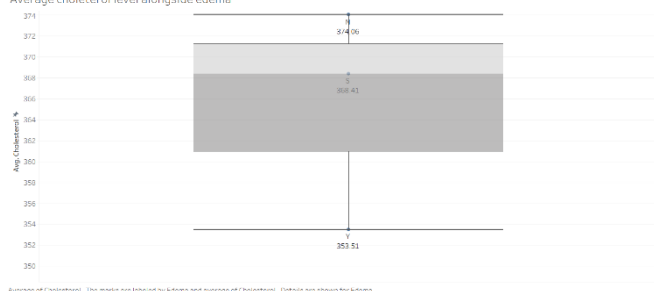
Visualization 1:

The visualization used is scatter plot.

## Visualization 1:


Relationshiup between cholesterol of patient with edema

Sum of Jitter_X(Adjust position of Y and N of edema using Random() value generator) vs. sum of Cholesterol. Color shows details about Edema. The marks are labeled by sum of Cholesterol and Edema.

Result: From the visualization, we can clearly see there is a large difference in level of cholesterol with respect to patient with all stages of edema, hence we can conclude that the hypothesis is invalid.
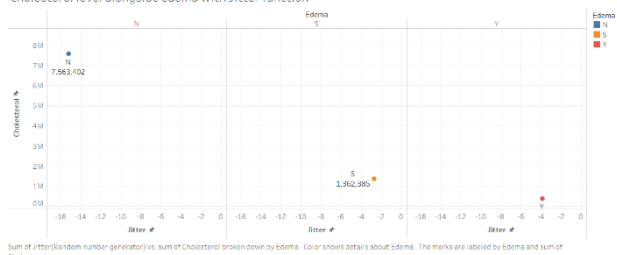
## Visualization 2:


Average choleterol level alongside edema

Average of Cholesterol. The marks are labeled by Edema and average of Cholesterol. Details are shown for Edema.

For this visualization, I have used box-and-whisker plot.
Result: For the average cholesterol level of patient with edema is 353.51 while the patient without edema it is 374.06, implying to the significant difference, from the visualization we can nullify our hypothesis.

## Visualization 3:


Cholesterol level alongside edema with Jitter function

Sum of Jitter(Random number generator) vs. sum of Cholesterol broken down by Edema. Color shows details about Edema. The marks are labeled by Edema and sum of Cholesterol.

Here, visualization used is called Scatter plot with Jitter. Jitter is a function that will adjust a position of 'Y' and 'N' of edema using Random () function. Here we will use an amplifier with a value of 0.1.
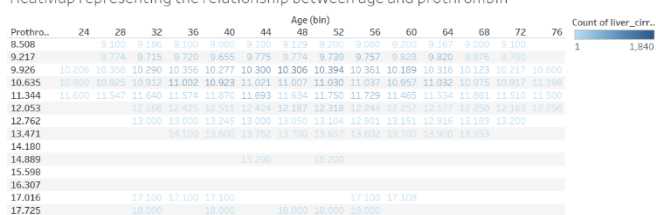
Result: From the visualization, it can clearly conclude that there is a significant difference between cholesterol levels of patients with and without edema. Hence our hypothesis is invalid.

*D. There is no significant relationship between age and prothrombin time.*

Visualization 1:
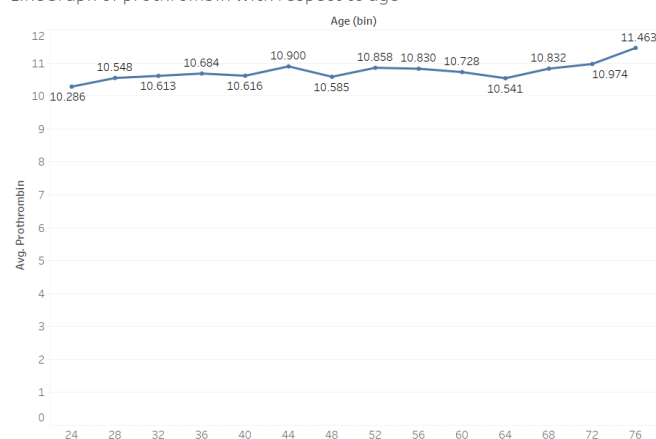Visualization used to check the validity of hypothesis is heat map.

## HeatMap representing the relationship between age and prothrombin


Average of Prothrombin broken down by Age (bin) vs. Prothrombin (bin). Color shows count of total patients.

Result: From the visualization, the average prothrombin time increases with increase in age group, this proves that our hypothesis is invalid.

## Visualization 2:


LineGraph of prothrombin with respect to age

The trend of average of Prothrombin for Age (bin). The marks are labeled by average of Prothrombin.
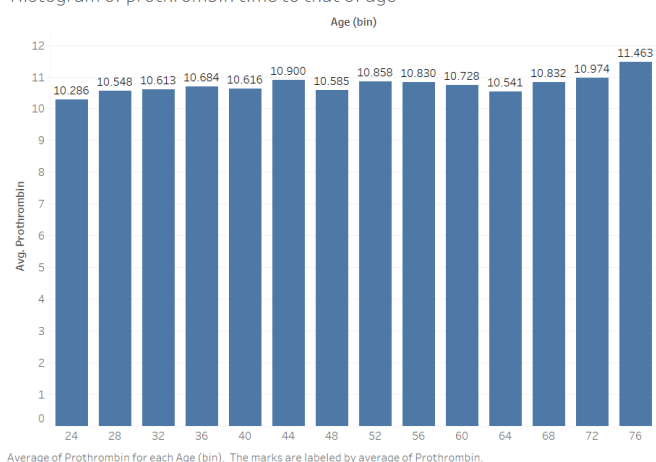
Visualization used is called line chart.
Line chart is representation of data with points and joining the points to observe the trends using lines.[7]

Result: From the visualization, the lowest prothrombin time is 10.286 by the age group of 20-24 (lowest age group), while the highest prothrombin levels are 11.463 of the age group 73-76 (highest age group), implicating the rise of prothrombin level with increase in age, hence our hypothesis is invalid.

## Visualization 3:


Histogram of prothrombin time to that of age

Average of Prothrombin for each Age (bin). The marks are labeled by average of Prothrombin.

Visualization used is histogram.
Histogram is the distribution of data by making bins in the form of bars.[8]

Result:

From the visualization, it can be clearly observed that the relationship between age and prothrombin exists, making our hypothesis invalid.

## CONCLUSION

From the above results we can conclude that,

- Patient receiving different drugs have different survival rate.
- The ascites, hepatomegaly, or edema levels are directly proportional to liver cirrhosis stages.
- Cholesterol levels are significantly different between the patient with edema and without edema.
- Age and prothrombin time are related, with increase in age the prothrombin time increases.

## REFERENCES

[1]   Mayo Clinic
[2]   Liver Cirrhosis Stage Classification Dataset by Aadarsh Velu
[3]   Bar Chart wikipedia-https://en.wikipedia.org/wiki/Bar_chart
[4]   A complete guide to scatter plot- Atlassian
[5]   Heat map- Optimizely.
[6]   Build a packed bubble chart by tableau.
[7]   A complete guide to line chart- Atlassian
[8]   Histogram Unveiled- Atlassian