# About Data ⚕

The data contains the information collected from the Mayo Clinic trial in primary biliary cirrhosis (PBC) of the liver conducted between 1974 and 1984. A description of the clinical background for the trial and the covariates recorded here is in Chapter 0, especially Section 0.2 of Fleming and Harrington, Counting Processes and Survival Analysis, Wiley, 1991. A more extended discussion can be found in Dickson, et al., Hepatology 10:1-7 (1989) and in Markus, et al., N Eng J of Med 320:1709-13 (1989).

A total of 424 PBC patients, referred to Mayo Clinic during that ten-year interval, met eligibility criteria for the randomized placebo-controlled trial of the drug D-penicillamine. The first 312 cases in the dataset participated in the randomized trial and contain largely complete data. The additional 112 cases did not participate in the clinical trial but consented to have basic measurements recorded and to be followed for survival. Six of those cases were lost to follow-up shortly after diagnosis, so the data here are on an additional 106 cases as well as the 312 randomized participants.

The dataset consists of following columns :

1. ID: unique identifier
2. N_Days: number of days between registration and the earlier of death, transplantation, or study analysis time in July 1986
3. Status: status of the patient C (censored), CL (censored due to liver tx), or D (death)
4. Drug: type of drug D-penicillamine or placebo
5. Age: age in [days]
6. Sex: M (male) or F (female)
7. Ascites: presence of ascites N (No) or Y (Yes)
8. Hepatomegaly: presence of hepatomegaly N (No) or Y (Yes)
9. Spiders: presence of spiders N (No) or Y (Yes)
10. Edema: 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)
11. Bilirubin: serum bilirubin in [mg/dl]
12. Cholesterol: serum cholesterol in [mg/dl]
13. Albumin: albumin in [gm/dl]
14. Copper: urine copper in [ug/day]
15. Alk_Phos: alkaline phosphatase in [U/liter]
16. SGOT: SGOT in [U/ml]
17. Triglycerides: triglicerides in [mg/dl]
18. Platelets: platelets per cubic [ml/1000]
19. Prothrombin: prothrombin time in seconds [s]
20. Stage: histologic stage of disease (1, 2, 3, or 4)

# Imports

In [1]:

```python
# %load_ext google.colab.data_table
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

In [2]:

```python
#@title Default title text
df = pd.read_csv(r'C:\Users\91701\Desktop\livr prediction/cirrhosis.csv',index_col='ID')
df.head()
```

Out[2]:

| ID | N_Days | Status | Drug | Age | Sex | Ascites | Hepatomegaly | Spiders | Edema | Bilirubi |
|----|--------|--------|------|-----|-----|---------|--------------|---------|-------|----------|
| 1 | 400 | D | D-penicillamine | 21464 | F | Y | Y | Y | Y | 14. |
| 2 | 4500 | C | D-penicillamine | 20617 | F | N | Y | Y | N | 1. |
| 3 | 1012 | D | D-penicillamine | 25594 | M | N | N | N | S | 1. |
| 4 | 1925 | D | D-penicillamine | 19994 | F | N | Y | Y | S | 1. |
| 5 | 1504 | CL | Placebo | 13918 | F | N | Y | Y | N | 3. |

In [3]:

```python
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 418 entries, 1 to 418
Data columns (total 19 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   N_Days         418 non-null    int64
 1   Status         418 non-null    object
 2   Drug           312 non-null    object
 3   Age            418 non-null    int64
 4   Sex            418 non-null    object
 5   Ascites        312 non-null    object
 6   Hepatomegaly   312 non-null    object
 7   Spiders        312 non-null    object
 8   Edema          418 non-null    object
 9   Bilirubin      418 non-null    float64
 10  Cholesterol    284 non-null    float64
 11  Albumin        418 non-null    float64
 12  Copper         310 non-null    float64
 13  Alk_Phos       312 non-null    float64
 14  SGOT           312 non-null    float64
 15  Tryglicerides  282 non-null    float64
 16  Platelets      407 non-null    float64
 17  Prothrombin    416 non-null    float64
 18  Stage          412 non-null    float64
dtypes: float64(10), int64(2), object(7)
memory usage: 65.3+ KB
```

This data set has about 19 features. These features are related to the patient's details like age, sex, etc. and patient's blood tests like prothrombin, triglycerides, platelets levels, etc. All these factors help in understanding a patient's chances of liver cirrhosis.

we have some NA values in our data, lets look at some statistical summary of numerical columns in out dataset.

In [4]:

```python
df.describe()
```

Out[4]:

|        | N_Days      | Age          | Bilirubin  | Cholesterol | Albumin    | Copper     | Alk_Ph     |
|--------|-------------|--------------|------------|-------------|------------|------------|------------|
| count  | 418.000000  | 418.000000   | 418.000000 | 284.000000  | 418.000000 | 310.000000 | 312.0000   |
| mean   | 1917.782297 | 18533.351675 | 3.220813   | 369.510563  | 3.497440   | 97.648387  | 1982.6557  |
| std    | 1104.672992 | 3815.845055  | 4.407506   | 231.944545  | 0.424972   | 85.613920  | 2140.3888  |
| min    | 41.000000   | 9598.000000  | 0.300000   | 120.000000  | 1.960000   | 4.000000   | 289.0000   |
| 25%    | 1092.750000 | 15644.500000 | 0.800000   | 249.500000  | 3.242500   | 41.250000  | 871.5000   |
| 50%    | 1730.000000 | 18628.000000 | 1.400000   | 309.500000  | 3.530000   | 73.000000  | 1259.0000  |
| 75%    | 2613.500000 | 21272.500000 | 3.400000   | 400.000000  | 3.770000   | 123.000000 | 1980.0000  |
| max    | 4795.000000 | 28650.000000 | 28.000000  | 1775.000000 | 4.640000   | 588.000000 | 13862.4000 |

## We have some missing values in our data, lets see how many and in which columns.

In [5]:

```
df.isna().sum()
```

Out[5]:

```
N_Days             0
Status             0
Drug             106
Age                0
Sex                0
Ascites          106
Hepatomegaly     106
Spiders          106
Edema              0
Bilirubin          0
Cholesterol      134
Albumin            0
Copper           108
Alk_Phos         106
SGOT             106
Tryglicerides    136
Platelets         11
Prothrombin        2
Stage              6
dtype: int64
```

# Handling Missing Values

This is a problem, we could just get rid of all examples with NA values, but in this case our case of small dataset we cannot afford that.

We will impute the missing entries with some statistical calculations.

# We have two different types of data

1. Numerical data ( Age, Cholesterol, Platelets.. etc)
2. Categorical Data ( Drug, Sex, Spiders..etc)

We will have to use different imputation for each type

1. For the numerical type we can use mean or median. In this case we will go with median to avoid skewing in the presence of outliers
2. For Categorical type we will impute the most frequent class.

In [6]:

```python
# For Numerical Type
df.select_dtypes(include=(['int64', 'float64'])).isna().sum()
```

Out[6]:

```
N_Days           0
Age              0
Bilirubin        0
Cholesterol    134
Albumin          0
Copper         108
Alk_Phos       106
SGOT           106
Tryglicerides  136
Platelets       11
Prothrombin      2
Stage            6
dtype: int64
```

In [7]:

```python
df.select_dtypes(include=(['int64', 'float64'])).isna().sum()
df_num_col = df.select_dtypes(include=(['int64', 'float64'])).columns
for c in df_num_col:
    df[c].fillna(df[c].median(), inplace=True)

df.select_dtypes(include=(['int64', 'float64'])).isna().sum()
```

Out[7]:

```
N_Days           0
Age              0
Bilirubin        0
Cholesterol      0
Albumin          0
Copper           0
Alk_Phos         0
SGOT             0
Tryglicerides    0
Platelets        0
Prothrombin      0
Stage            0
dtype: int64
```

In [8]:

```python
# For Categorical type
df.select_dtypes(include=('object')).isna().sum()
```

Out[8]:

```
Status          0
Drug          106
Sex             0
Ascites       106
Hepatomegaly  106
Spiders       106
Edema           0
dtype: int64
```

In [9]:

```python
df_cat_col = df.select_dtypes(include=('object')).columns
for c in df_cat_col:
    df[c].fillna(df[c].mode().values[0], inplace=True)

df.select_dtypes(include=('object')).isna().sum()
```

Out[9]:

```
Status          0
Drug            0
Sex             0
Ascites         0
Hepatomegaly    0
Spiders         0
Edema           0
dtype: int64
```

# Exploratory Data Analysys

## Lets dive into the data and visualize it, this often revels interesting patterns.

First lets take a look at how many examples per calss do we have in our dataset.

In [10]:

```python
df['Stage'].value_counts()
```
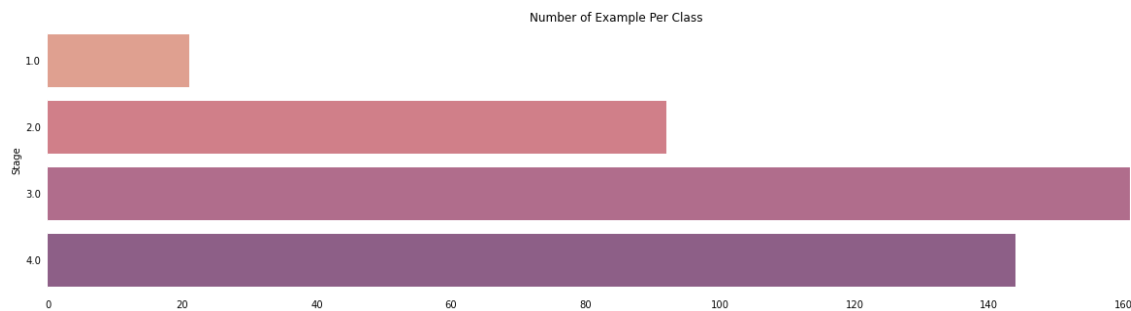
Out[10]:

```
3.0    161
4.0    144
2.0     92
1.0     21
Name: Stage, dtype: int64
```

In [11]:

```
plt.figure(figsize=(21,5))
sns.countplot(y=df['Stage'], palette="flare", alpha=0.8, )
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Number of Example Per Class')
```
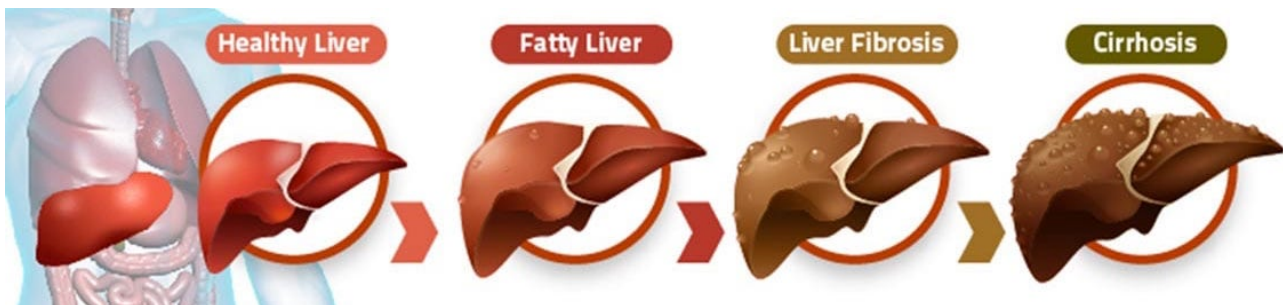
Out[11]:

```
Text(0.5, 1.0, 'Number of Example Per Class')
```



As we can observe we have class imbalances in our dataset i.e some classes have more examples than other. This could make it difficult for our model to train and achieve desired score. No worries, we can fix that later.

# Setting up Target and Features

**For this demonstration we will keep things simple by predicting one of the two classes i.e (Cirrhosis and No Cirrhois). I have another project that predicts the stage of the disease, converting this problem into a multicalss classification.**



In [12]:

```
# Converting Target categories into intigers 1 for Cirrhosis, 0 otherwise
df['Stage'] = np.where(df['Stage'] == 4,1,0)
```

## Lets observe some Features with their relation with the disease

In [13]:

```python
plt.figure(figsize=(21.2,10))

plt.subplot(2,3,1)
sns.countplot(x=df['Stage'], hue=df['Sex'], palette='Blues', alpha=0.9)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Disease Stage Across Gender')

plt.subplot(2,3,2)
sns.countplot(x=df['Stage'], hue=df['Ascites'], palette='Purples', alpha=0.9)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Ascites proportion across Stages')

plt.subplot(2,3,3)
sns.countplot(x=df['Stage'], hue=df['Drug'], palette='Blues', alpha=0.9)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Medications prescribed across Stages');

plt.subplot(2,3,4)
sns.countplot(x=df['Stage'], hue=df['Hepatomegaly'], palette='Purples', alpha=0.9)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Hepatomegaly');

plt.subplot(2,3,5)
sns.countplot(x=df['Stage'], hue=df['Spiders'], palette='Blues', alpha=0.9)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Presence of Spiders across stages');

plt.subplot(2,3,6)
sns.countplot(x=df['Stage'], hue=df['Edema'], palette='Purples', alpha=0.9)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Edema');
```
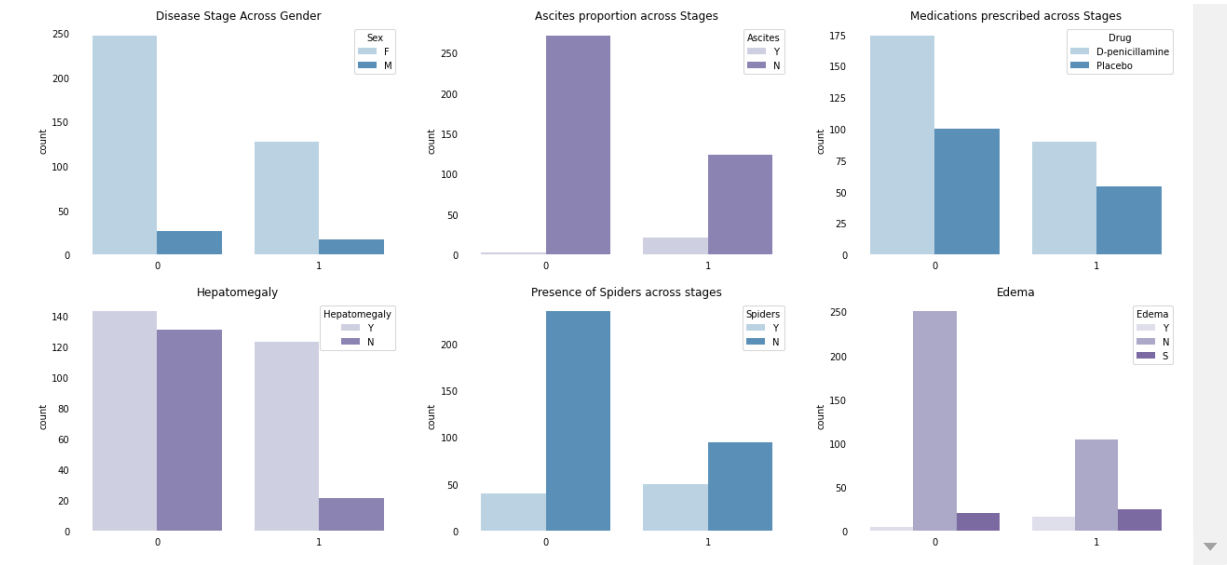
There are some interesting insights if we observe closely. Take the case at Ascites, we observe that the rist of disease is higher with increase in Ascites. also presence of spiders has a positive relation with disease risk.

In [14]:

```python
#@title Distribution Polts
plt.figure(figsize=(20.6,15))

plt.subplot(3,3,1)
sns.kdeplot(df['Cholesterol'], hue=df['Stage'], fill=True, palette='Purples')
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Cholesterol Distribution in stages');

plt.subplot(3,3,2)
sns.kdeplot(df['Bilirubin'], hue=df['Stage'], fill=True, palette='Blues', common_norm=True)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Bilirubin');

plt.subplot(3,3,3)
sns.kdeplot(df['Tryglicerides'], hue=df['Stage'], fill=True, palette='Purples', common_norm
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Tryglicerides');

plt.subplot(3,3,4)
sns.kdeplot(df['Age'], hue=df['Stage'], fill=True, palette='Blues', common_norm=True)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Age Distribution in stages');

plt.subplot(3,3,5)
sns.kdeplot(df['Prothrombin'], hue=df['Stage'], fill=True, palette='Purples', common_norm=T
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Prothrombin');

plt.subplot(3,3,6)
sns.kdeplot(df['Copper'], hue=df['Stage'], fill=True, palette='Blues', common_norm=True)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Copper');

plt.subplot(3,3,7)
sns.kdeplot(df['Platelets'], hue=df['Stage'], fill=True, palette='Purples')
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Platelets in stages');

plt.subplot(3,3,8)
sns.kdeplot(df['Albumin'], hue=df['Stage'], fill=True, palette='Blues', common_norm=True)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('Albumin');
```
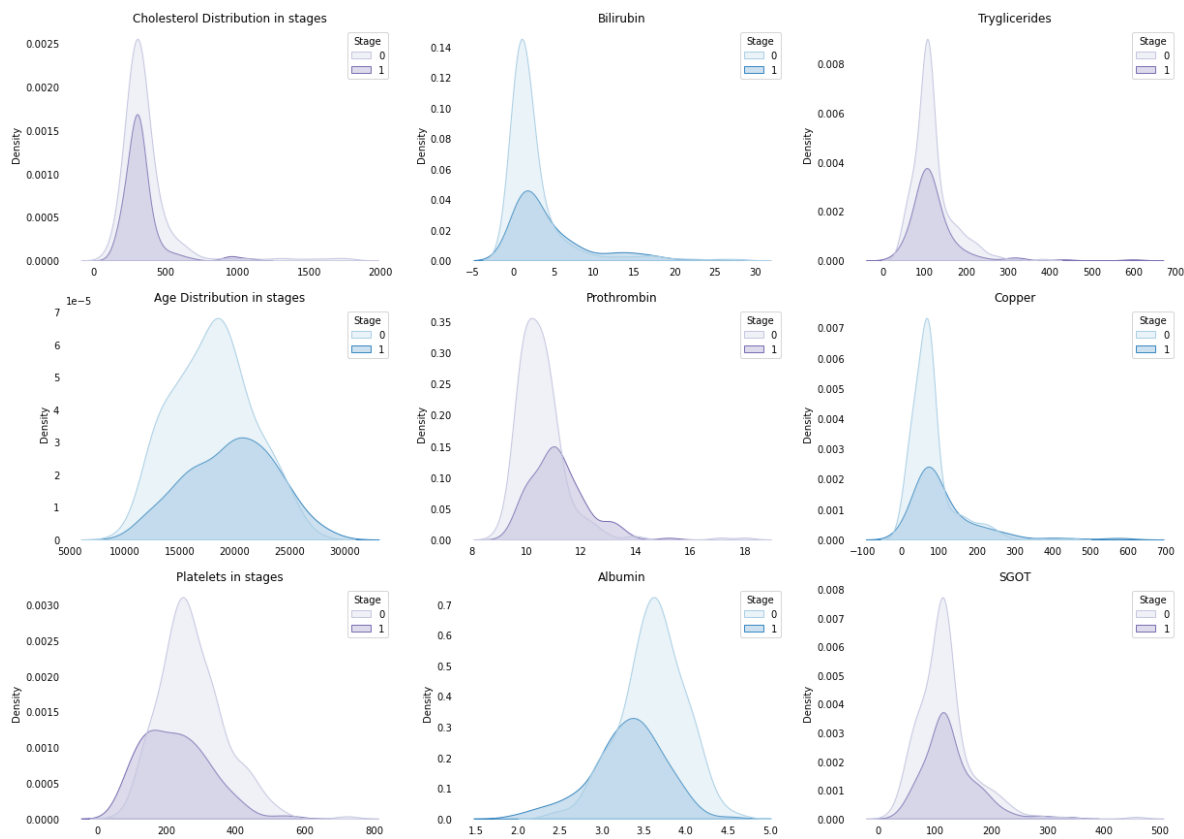
```
plt.subplot(3,3,9)
sns.kdeplot(df['SGOT'], hue=df['Stage'], fill=True, palette='Purples', common_norm=True)
sns.despine(top=True, right=True, bottom=True, left=True)
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('')
plt.title('SGOT');
```



**Looking at the feature distribution we can observe that in features such as Age, Prothrombin, Copper the risk of the disease increase with increase in feature value, thus having a positive co-relation on with the disease probability. Lets fit a regression line to check.**
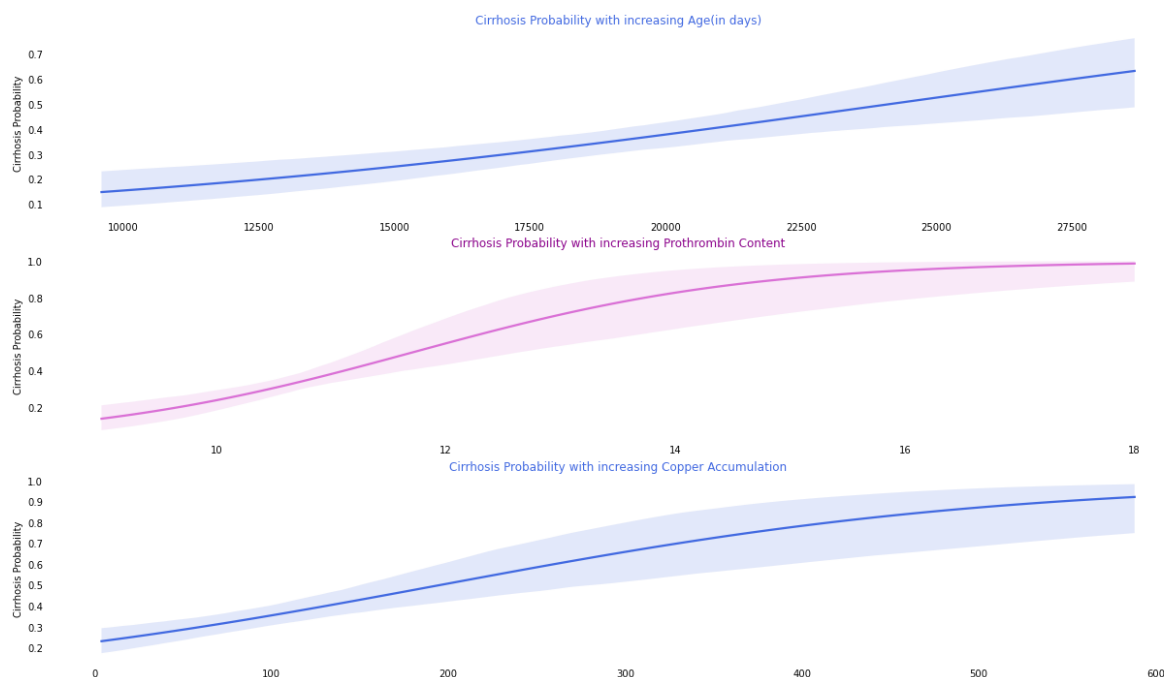
In [15]:

```python
#@title Regression Plots of Positive Correlated Features.
plt.figure(figsize=(21,12))

plt.subplot(3,1,1)
sns.regplot(x=df['Age'], y=df['Stage'], scatter=False, logistic=True, color='royalblue')
sns.despine(fig=None, ax=None, top=True, right=True, left=True, bottom=True, offset=None, t
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('');
plt.ylabel('Cirrhosis Probability');
plt.setp(plt.title('Cirrhosis Probability with increasing Age(in days)'), color='royalblue'

plt.subplot(3,1,2)
sns.regplot(x=df['Prothrombin'], y=df['Stage'], scatter=False, logistic=True, color='orchid
sns.despine(fig=None, ax=None, top=True, right=True, left=True, bottom=True, offset=None, t
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('');
plt.ylabel('Cirrhosis Probability');
plt.setp(plt.title('Cirrhosis Probability with increasing Prothrombin Content'), color='dar

plt.subplot(3,1,3)
sns.regplot(x=df['Copper'], y=df['Stage'], scatter=False, logistic=True, color='royalblue')
sns.despine(fig=None, ax=None, top=True, right=True, left=True, bottom=True, offset=None, t
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('');
plt.ylabel('Cirrhosis Probability');
plt.setp(plt.title('Cirrhosis Probability with increasing Copper Accumulation'), color='roy
```



Looks like the data checks with our intuition. These parameters indeed increase the risk of the disease.

## We can also observe some features such as Platelets, Albumin, Cholesterol where the probability of disease decrease with increase in feature value. Lets tally that with some more regression plots.

In [16]:

```python
#@title Regression Plots of negatively correlated Features.
plt.figure(figsize=(21,12))

plt.subplot(3,1,1)
sns.regplot(x=df['Platelets'], y=df['Stage'], scatter=False, logistic=True, color='orchid')
sns.despine(fig=None, ax=None, top=True, right=True, left=True, bottom=True, offset=None, t
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('');
plt.ylabel('Cirrhosis Probability');
plt.setp(plt.title('Cirrhosis Probability with Platelets'), color='darkmagenta');

plt.subplot(3,1,2)
sns.regplot(x=df['Albumin'], y=df['Stage'], scatter=False, logistic=True, color='royalblue'
sns.despine(fig=None, ax=None, top=True, right=True, left=True, bottom=True, offset=None, t
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('');
plt.ylabel('Cirrhosis Probability');
plt.setp(plt.title('Cirrhosis Probability with Albumin Content'), color='royalblue');

plt.subplot(3,1,3)
sns.regplot(x=df['Cholesterol'], y=df['Stage'], scatter=False, logistic=True, color='orchid
sns.despine(fig=None, ax=None, top=True, right=True, left=True, bottom=True, offset=None, t
plt.tick_params(axis='both', which='both', bottom=False, top=False, left=False)
plt.xlabel('');
plt.ylabel('Cirrhosis Probability');
plt.setp(plt.title('Cirrhosis Probability Cholesterol'), color='darkmagenta') ;
```



Platelets, Albumin checks with our logic the findings about Cholesterol seems interesting! Looks like people with high Cholesterol have lower risk of Cirrhosis, this might not sound correct but our data certainly shows so.

This should help our model predict the target. We will be looking at what features contribute the most in later part of the project.

# Preprocessing data

In [17]:

```python
# replacing catagorical data with intigers.
df['Sex'] = df['Sex'].replace({'M':0, 'F':1})                          # Male : 0 , F
df['Ascites'] = df['Ascites'].replace({'N':0, 'Y':1})                  # N : 0, Y : 1
df['Drug'] = df['Drug'].replace({'D-penicillamine':0, 'Placebo':1})    # D-penicillam
df['Hepatomegaly'] = df['Hepatomegaly'].replace({'N':0, 'Y':1})        # N : 0, Y : 1
df['Spiders'] = df['Spiders'].replace({'N':0, 'Y':1})                  # N : 0, Y : 1
df['Edema'] = df['Edema'].replace({'N':0, 'Y':1, 'S':-1})              # N : 0, Y : 1
df['Status'] = df['Status'].replace({'C':0, 'CL':1, 'D':-1})           # 'C':0, 'CL':
```

## We will not be using 'Status' and 'N_days' as our features since this will cause data Leakage.

In [18]:

```python
# Setting up Features and Target
X = df.drop(['Status', 'N_Days', 'Stage'], axis=1)
y = df.pop('Stage')
```

## Earlier while examining the distribution of target in the data, we found that the data is unevenly distributed, that is there are more examples of a certain class than other.

###To tackel this imbalance we will use **Stratified k-Fold** Cross Validation.

# Model Selection.

lets first try out a quick Logestic regression classifier and see how it performs.

In [20]:

```python
from sklearn.metrics import classification_report
log_model_predict = log_model.predict(test)
log_model_predict_proba = log_model.predict_proba(test)

print(classification_report(y.iloc[test_index], log_model_predict))
```

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0 | 0.70 | 0.78 | 0.74 | 27 |
| 1 | 0.45 | 0.36 | 0.40 | 14 |
| accuracy |  |  | 0.63 | 41 |
| macro avg | 0.58 | 0.57 | 0.57 | 41 |
| weighted avg | 0.62 | 0.63 | 0.62 | 41 |

In [21]:

```python
from sklearn.metrics import roc_auc_score
from sklearn.metrics import roc_curve, auc

fpr, tpr, threshold = roc_curve(y.iloc[test_index], log_model_predict_proba[:,1])
roc_auc = auc(fpr, tpr)

print('AUC : ', roc_auc_score(y.iloc[test_index], log_model_predict_proba[:,1]))
```

AUC :   0.6507936507936508

In [73]:

```python
def trainer(X_train, y_train, X_test, y_test):

    models= [[' SVM ',SVC()],
             [' Decision Tree ', DecisionTreeClassifier()],
             [' Random Forest ', RandomForestClassifier()],
             [' Logistic Regression ', LogisticRegression(max_iter=200)],
             [' AdaBoost ', AdaBoostClassifier()],
             [' KNN ', KNeighborsClassifier()]]

    scores = []

    for model_name, model in models:

        model = model
        model.fit(X_train, y_train)
        pred = model.predict(X_test)
        cm_model = confusion_matrix(y_test, pred)
        scores.append(accuracy_score(y_test, model.predict(X_test)))


        print(Back.YELLOW + Fore.BLACK + Style.BRIGHT + model_name)
        print(Back.RESET)
        print(cm_model)
        print('\n' + Fore.BLUE + 'Training Acc.  : ' + Fore.GREEN + str(round(accuracy_scor
        print(Fore.BLUE + 'Validation Acc.: ' + Fore.GREEN + str(round(accuracy_score(y_tes
        print(Fore.CYAN + classification_report(y_test, pred))

        visualizer = ROCAUC(model)
        visualizer.fit(X_train, y_train)
        visualizer.score(X_test, y_test)
        visualizer.show()

        print('\n' + Fore.BLACK + Back.WHITE + '※※※※※※※※※※※※※※※※※※※※※※※※※※※※※※※


    return scores
```

In [63]:

```python
import warnings
import itertools

import pandas as pd
import numpy as np

import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler, LabelEncoder
from sklearn.metrics import confusion_matrix, accuracy_score, classification_report, roc_cu
from sklearn.pipeline import Pipeline
from sklearn.tree import DecisionTreeClassifier
from sklearn.linear_model import LogisticRegression
from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
from sklearn.neighbors import KNeighborsClassifier
from sklearn.svm import SVC



from xgboost import XGBClassifier
from catboost import CatBoostClassifier



from colorama import Fore, Back, Style
#from yellowbrick.classifier import ROCAUC
```

In [62]:

```python
conda install -c districtdatalabs yellowbrick
```

Note: you may need to restart the kernel to use updated packages.

usage: conda-script.py [-h] [-V] command ...
conda-script.py: error: unrecognized arguments: yellowbrick

In [43]:

```python
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.2, random_state = 2
```

In [96]:

```python
def score_vis(score):

    names = ['Random Forest', 'Logistic Regression']

    plt.rcParams['figure.figsize']=8,10
    ax = sns.barplot(x=names, y=score, palette = "plasma", saturation =2.0)
    plt.xlabel('Model', fontsize = 20 )
    plt.ylabel('Accuracy(%)', fontsize = 20)
    plt.title('Model Comparison - Test set(random forest vs logistic regression)', fontsize
    plt.xticks(fontsize = 12, horizontalalignment = 'center', rotation = 8)
    plt.yticks(fontsize = 12)
    for i in ax.patches:
        width, height = i.get_width(), i.get_height()
        x, y = i.get_xy()
        ax.annotate(f'{round(height,2)}%', (x + width/2, y + height*1.02), ha='center', fon
    plt.show()
```

In [91]:

```python
def trainer(X_train, y_train, X_test, y_test):

    models= [[' Random Forest ', RandomForestClassifier()],
             [' Logistic Regression ', LogisticRegression(max_iter=200)]]
    scores = []

    for model_name, model in models:

        model = model
        model.fit(X_train, y_train)
        pred = model.predict(X_test)
        cm_model = confusion_matrix(y_test, pred)
        scores.append(accuracy_score(y_test, model.predict(X_test)))
        print(Back.YELLOW + Fore.BLACK + Style.BRIGHT + model_name)
        print(Back.RESET)
        print(cm_model)
        print('\n' + Fore.BLUE + 'Training Acc.  : ' + Fore.GREEN + str(round(accuracy_scor
        print(Fore.BLUE + 'Validation Acc.: ' + Fore.GREEN + str(round(accuracy_score(y_tes
        print(Fore.CYAN + classification_report(y_test, pred))
        print('\n' + Fore.BLACK + Back.WHITE + '☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀☀


    return scores
```

In [92]:

```
scores = trainer(X_train, y_train, X_test, y_test)
```

 Random Forest

[[49  2]
 [18 15]]

Training Acc.  : 100.0%
Validation Acc.: 76.19%

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0          | 0.73      | 0.96   | 0.83     | 51      |
| 1          | 0.88      | 0.45   | 0.60     | 33      |
|            |           |        |          |         |
| accuracy   |           |        | 0.76     | 84      |
| macro avg  | 0.81      | 0.71   | 0.72     | 84      |
| weighted avg | 0.79    | 0.76   | 0.74     | 84      |


⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂

 Logistic Regression

[[49  2]
 [23 10]]

Training Acc.  : 71.86%
Validation Acc.: 70.24%

|            | precision | recall | f1-score | support |
|------------|-----------|--------|----------|---------|
| 0          | 0.68      | 0.96   | 0.80     | 51      |
| 1          | 0.83      | 0.30   | 0.44     | 33      |
|            |           |        |          |         |
| accuracy   |           |        | 0.70     | 84      |
| macro avg  | 0.76      | 0.63   | 0.62     | 84      |
| weighted avg | 0.74    | 0.70   | 0.66     | 84      |


⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂⁂


C:\ProgramData\Anaconda3\New folder\a3\lib\site-packages\sklearn\linear_mode
l\_logistic.py:763: ConvergenceWarning: lbfgs failed to converge (status=1):
STOP: TOTAL NO. of ITERATIONS REACHED LIMIT.

Increase the number of iterations (max_iter) or scale the data as shown in:
    https://scikit-learn.org/stable/modules/preprocessing.html (https://scik
it-learn.org/stable/modules/preprocessing.html)
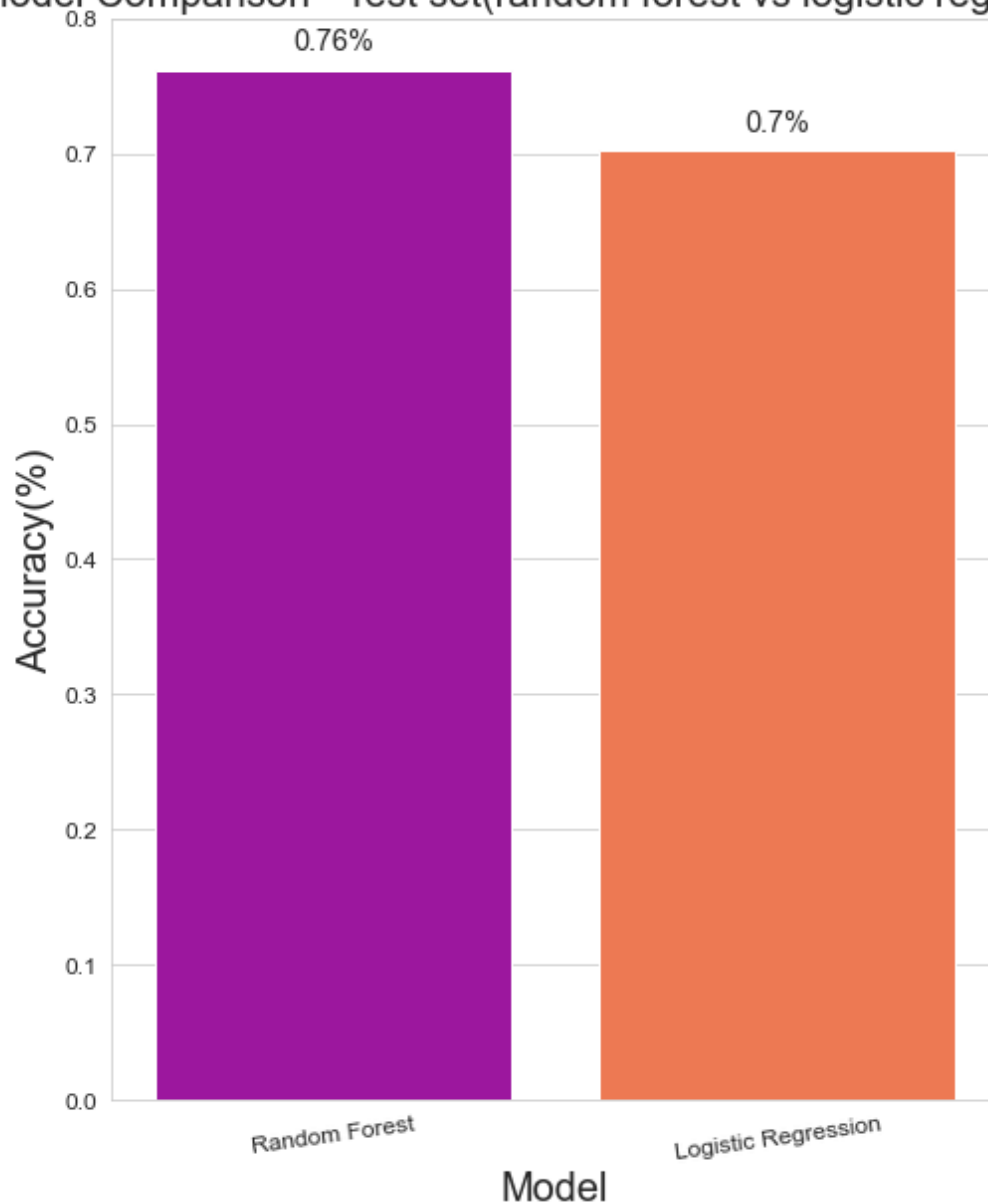Please also refer to the documentation for alternative solver options:
    https://scikit-learn.org/stable/modules/linear_model.html#logistic-regre
ssion (https://scikit-learn.org/stable/modules/linear_model.html#logistic-re
gression)
  n_iter_i = _check_optimize_result(

In [97]:

```
score_vis(scores)
```



Model Comparison - Test set(random forest vs logistic regression)

In [ ]:

In [ ]: