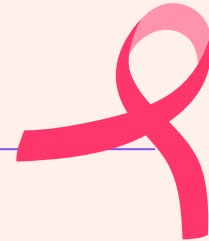


---

# METASTATIC CANCER & EQUITY IN HEALTHCARE

By. Ajibola-Elias Maryam





# CONTENTS



**01**

**Data  
Introduction**



**02**

**Data  
Understanding**



**03**

**Data  
Preprocessing**



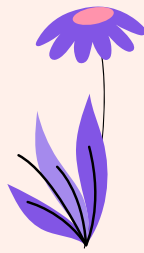
**04**

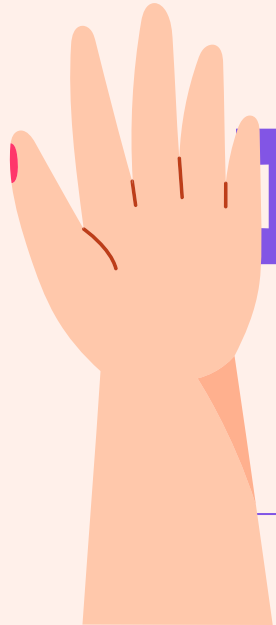
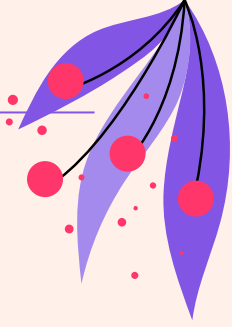
**Model Building**



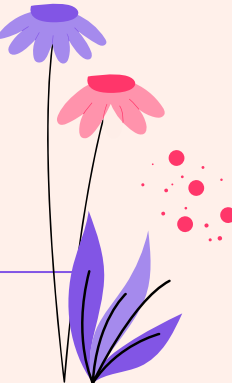
**05**

**Model Evaluation**



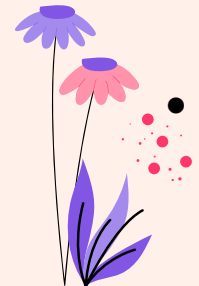


# Data Introduction



# Introduction

- Metastatic TNBC is considered the most aggressive TNBC & requires most urgent and timely treatment. Unnecessary delays in diagnosis & subsequent treatment can have devastating effects in these difficult cancers
- The primary goal is to detect relationships between demographics of the patients and the likelihood of getting timely treatment
- The secondary goal is to see if environmental hazard impact proper diagnosis and treatment.
- However the goal of our model is to just predict whether the patients received metastatic cancer diagnosis with 90 days or not.



# Introduction

## Columns

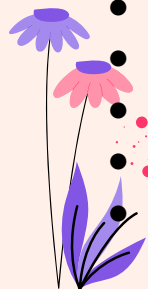
- **patient\_id**-unique identification number of patient
- **patient\_race**-Asian, African, American, Hispanic/Latino, White, Other Race
- **payer\_type**-payer type at Medicaid, Commercial, Medicare on the metastatic date
- **patient\_state**- Patient State (e.g. AL, AK, AZ, AR, CA, CO etc...) on the metastatic date
- **patient\_zip3**- Patient Zip3 (e.g. 190) on the metastatic date
- **patient\_age** - Derived from Patient Year of Birth (index year minus year of birth)
- **patient\_gender** - F, M on the metastatic date
- **bmi** - If Available, will show available BMI information (Earliest BMI recording post metastatic date)
- **breast\_cancer\_diagnosis\_code** - ICD10 or ICD9 diagnoses code
- **breast\_cancer\_diagnosis\_desc** - ICD10 or ICD9 code description. This column is raw text and may require NLP/ processing and cleaning
- **metastatic\_cancer\_diagnosis\_code** - ICD10 diagnoses code
- **metastatic\_first\_novel\_treatment** - Generic drug name of the first novel treatment (e.g. "Cisplatin") after metastatic diagnosis



# Introduction

## Columns

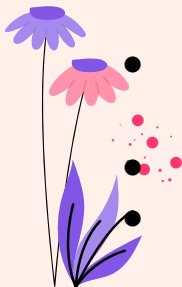
- `metastatic_first_novel_treatment_type` - Description of Treatment (e.g. Antineoplastic) of first novel treatment after metastatic diagnosis
- `region` - Region of patient location
- `division` - Division of patient location
- `population` - An estimate of the zip code's population.
- `density` - The estimated population per square kilometer.
- `age_median` - The median age of residents in the zip code.
- `male` - The percentage of residents who report being male (e.g. 55.1).
- `female` - The percentage of residents who report being female (e.g. 44.9).
- `married` - The percentage of residents who report being married (e.g. 44.9).
- `family_size` - The average size of resident families (e.g. 3.22).
- `income_household_median` - Median household income in USD.
- `income_household_six_figure` - Percentage of households that earn at least \$100,000 (e.g. 25.3)
- `home_ownership` - Percentage of households that own (rather than rent) their residence.



# Introduction

## Columns

- **housing\_units** - The number of housing units (or households) in the zip code.
- **home\_value** - The median value of homes that are owned by residents.
- **rent\_median** - The median rent paid by renters.
- **education\_college\_or\_above** - The percentage of residents with at least a 4-year degree.
- **labor\_force\_participation** - The percentage of residents 16 and older in the labor force.
- **unemployment\_rate** - The percentage of residents unemployed.
- **race\_white** - The percentage of residents who report their race White.
- **race\_black** - The percentage of residents who report their race as Black or African American.
- **race\_asian** - The percentage of residents who report their race as Asian.
- **race\_native** - The percentage of residents who report their race as American Indian and Alaska Native.
- **race\_pacific** - The percentage of residents who report their race as Native Hawaiian and Other Pacific Islander.
- **race\_other** - The percentage of residents who report their race as Some other race.
- **race\_multiple** - The percentage of residents who report their race as Two or more races.



# Introduction

## Columns

- **hispanic** - The percentage of residents who report being Hispanic. Note: Hispanic is considered to be an ethnicity and not a race.
- **age\_under\_10** - The percentage of residents aged 0-9.
- **age\_10\_to\_19** - The percentage of residents aged 10-19.
- **age\_20s** - The percentage of residents aged 20-29.
- **age\_30s** - The percentage of residents aged 30-39.
- **age\_40s** - The percentage of residents aged 40-49.
- **age\_50s** - The percentage of residents aged 50-59.
- **age\_60s** - The percentage of residents aged 60-69.
- **age\_70s** - The percentage of residents aged 70-79.
- **age\_over\_80** - The percentage of residents aged over 80.
- **divorced** - The percentage of residents divorced.
- **never\_married** - The percentage of residents never married.
- **widowed** - The percentage of residents never widowed.
- **family\_dual\_income** - The percentage of families with dual income earners.

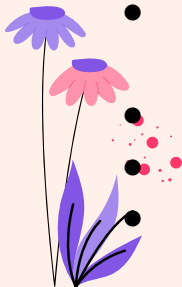




# Introduction

## Columns

- `income_household_under_5` - The percentage of households with income under \$5,000.
- `income_household_5_to_10` - The percentage of households with income from \$5,000-\$10,000.
- `income_household_10_to_15` - The percentage of households with income from \$10,000-\$15,000.
- `income_household_15_to_20` - The percentage of households with income from \$15,000-\$20,000.
- `income_household_20_to_25` - The percentage of households with income from \$20,000-\$25,000.
- `income_household_25_to_35` - The percentage of households with income from \$25,000-\$35,000.
- `income_household_35_to_50` - The percentage of households with income from \$35,000-\$50,000.
- `income_household_50_to_75` - The percentage of households with income from \$50,000-\$75,000.
- `income_household_75_to_100` - The percentage of households with income from \$75,000-\$100,000.
- `income_household_100_to_150` - The percentage of households with income from \$100,000-\$150,000.
- `income_household_150_over` - The percentage of households with income over \$150,000.
- `income_individual_median` - The median income of individuals in the zip code.



# Introduction

## Columns

- **poverty** - The median value of owner occupied homes.
- **rent\_burden** - The median rent as a percentage of the median renter's household income.
- **education\_less\_highschool** - The percentage of residents with less than a high school education.
- **education\_highschool** - The percentage of residents with a high school diploma but no more.
- **education\_some\_college** - The percentage of residents with some college but no more.
- **education\_bachelors** - The percentage of residents with a bachelor's degree (or equivalent) but no more.
- **education\_graduate** - The percentage of residents with a graduate degree.
- **education\_stem\_degree** - The percentage of college graduates with a Bachelor's degree or higher in a Science and Engineering (or related) field.
- **self-employed** - The percentage of households reporting self-employment income on their 2016 IRS tax return.
- **farmer** - The percentage of households reporting farm income on their 2016 IRS tax return.
- **disabled** - The percentage of residents who report a disability.
- **limited\_english** - The percentage of residents who only speak limited English.
- **commute\_time** - The median commute time of resident workers in minutes.

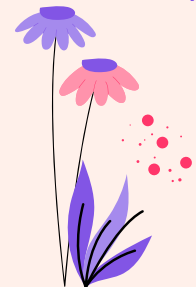
# Introduction

## Columns

- **health\_uninsured** - The percentage of residents who report not having health insurance.
- **veteran** - The percentage of residents who are veterans.
- **ozone** - Annual Ozone (O3) concentration data at Zip3 level. This data shows how air quality data may impact health.
- **PM25** - Annual Fine Particulate Matter (PM2.5) concentration data at Zip3 level. This data shows how air quality data may impact health.
- **NO2** - Annual Nitrogen Dioxide (NO2) concentration data at Zip3 level. This data shows how air quality data may impact health.

## Target

- **DiagPeriodL90D** - Diagnosis Period Less Than 90 Days. This is an indication of whether the cancer was diagnosed within 90 Days.





**Data**

**Understanding**



# Understanding the data

- Target: 1 - diagnosed within 90 day ;  
0-didn't get diagnosed within 90 days
- Binary classification
- Three csv files: Train, Test, and sample\_submission
- Size: 16.36MB
- 83 columns, 12906 rows
- 8327 missing values
- No duplicates
- Numerical & categorical values
- No Invalid Entries

# Understanding the data

## Outliers Detected

- The mean being bigger than the max values, suggested that that we might have outliers in the dataset
- Will be handles with log transformation..

```
df.describe()
```

	patient_id	patient_zip3	patient_age	bmi	population	density
count	12906.000000	12906.000000	12906.000000	3941.000000	12905.000000	12905.000000
mean	547381.196033	573.754300	59.183326	28.984539	20744.441237	1581.950419
std	260404.959974	275.447534	13.335216	5.696906	13886.903756	2966.305306
min	100063.000000	101.000000	18.000000	14.000000	635.545455	0.916667
25%	321517.000000	331.000000	50.000000	24.660000	9463.896552	171.857143
50%	543522.000000	554.000000	59.000000	28.190000	19154.190480	700.337500
75%	772671.750000	846.000000	67.000000	32.920000	30021.278690	1666.515385
max	999896.000000	999.000000	91.000000	85.000000	71374.131580	21172.000000

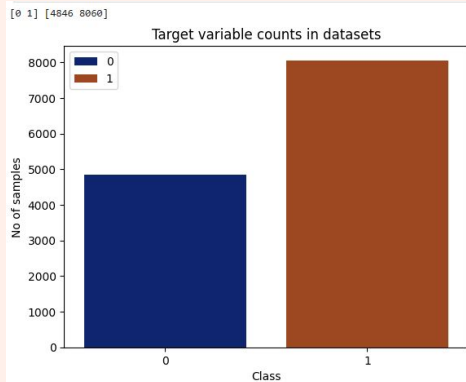
# Understanding the data

## Imbalanced data

- ❖ Not Diagnosed within 90 days[ 0 ]:  
4846
- ❖ Diagnosed within 90 day[ 1 ]: 8060
- ❖ Will be handled by resampling

```
target_column_name = 'DiagPeriodL90D'  
target_variable = df_encoded[target_column_name]  
  
# Adding the target variable to the selected_df  
selected_df[target_column_name] = target_variable
```

```
### Checking and plotting for the class in the target variable  
(unique, counts) = np.unique(selected_df['DiagPeriodL90D'], return_counts=True)  
print(unique, counts)  
sns.barplot(x=unique, y=counts, hue=unique, palette='dark', legend=True)  
plt.xlabel("Class")  
plt.ylabel("No of samples")  
plt.xticks()  
plt.title("Target variable counts in datasets")  
plt.show()  
plt.close()
```

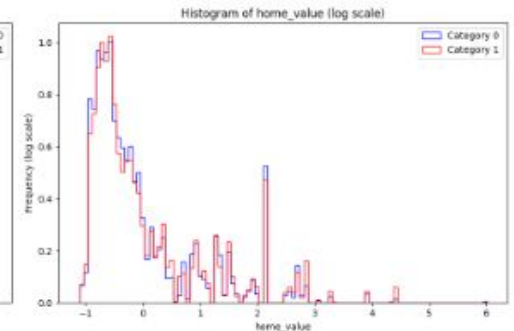
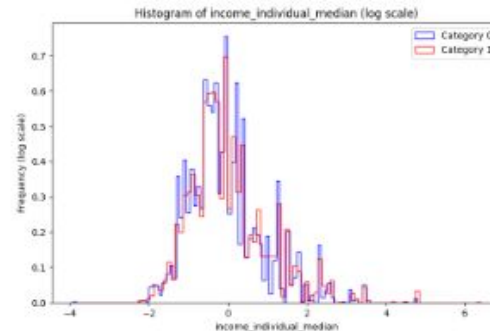
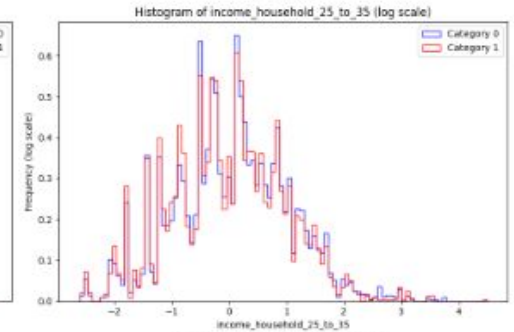
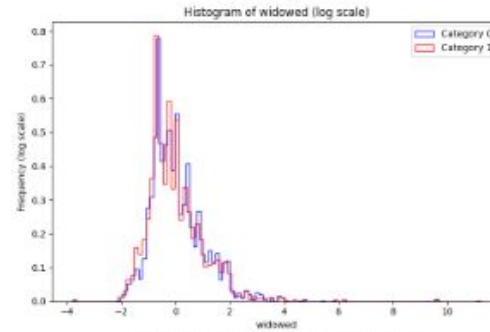
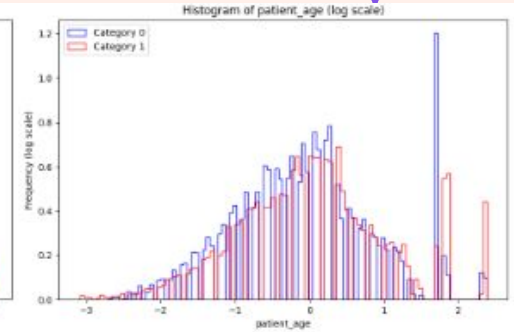
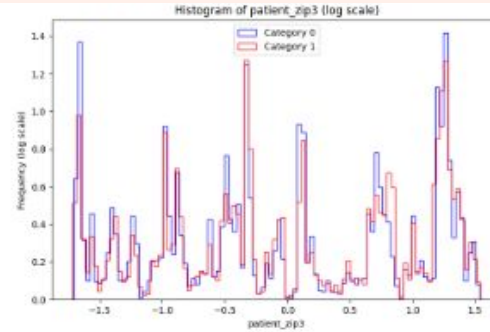


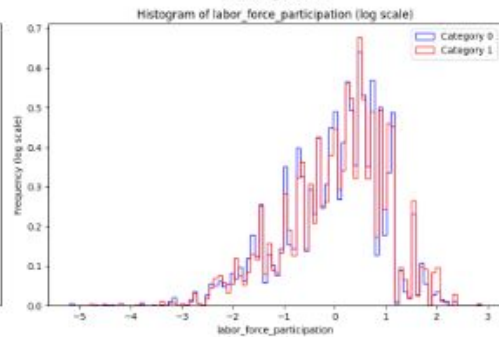
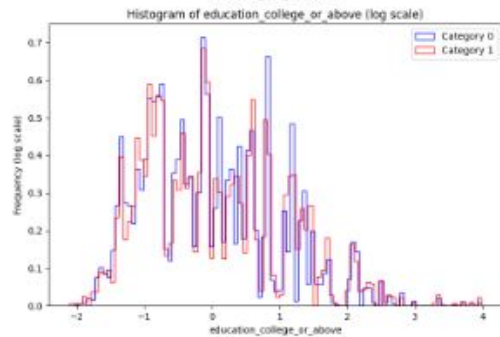
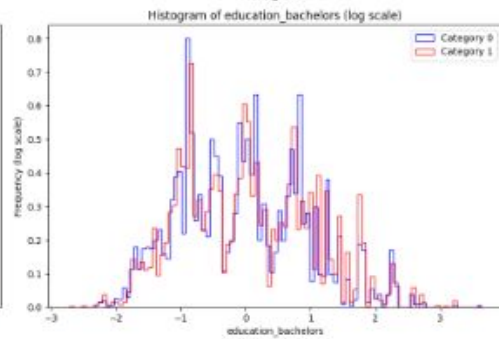
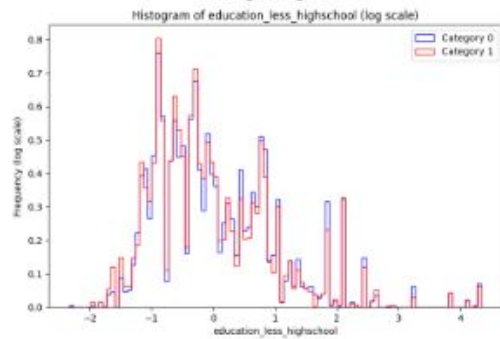
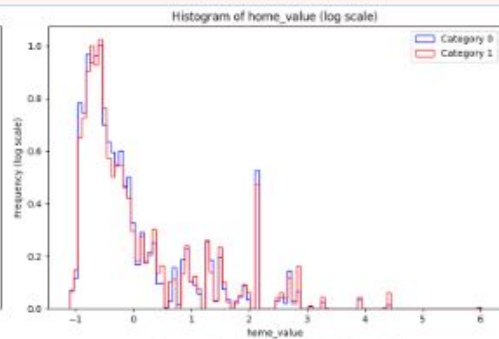
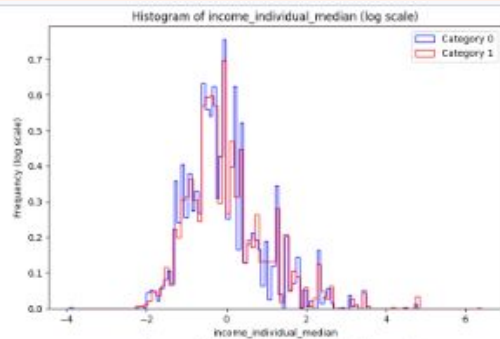


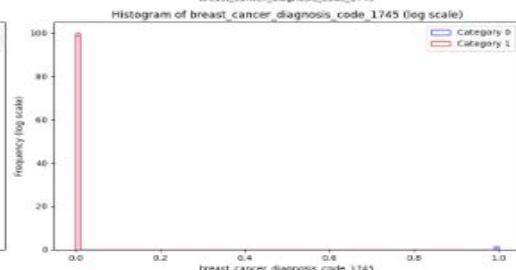
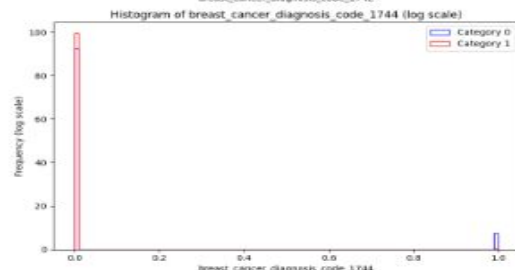
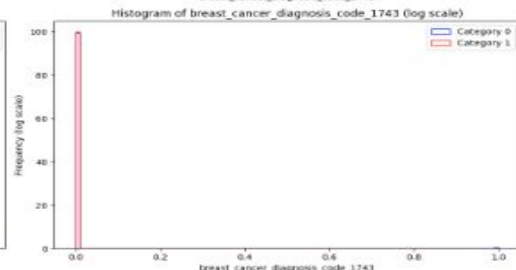
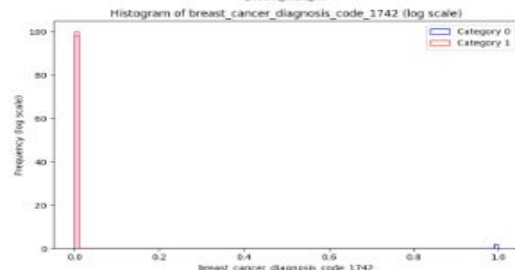
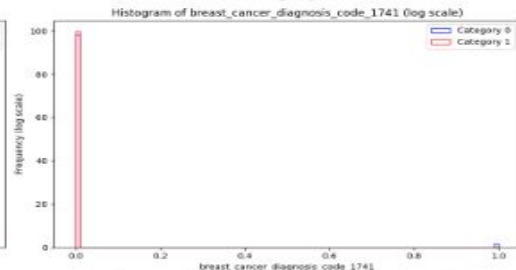
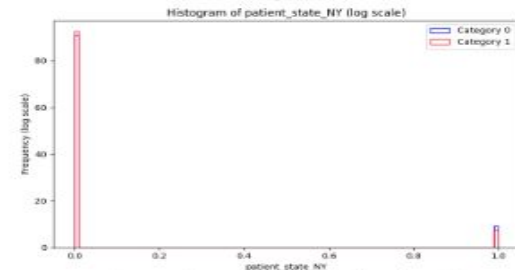
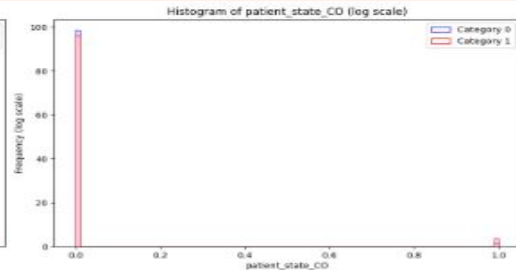
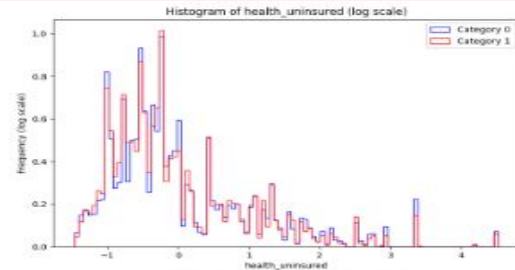
# DATA VISUALIZATION

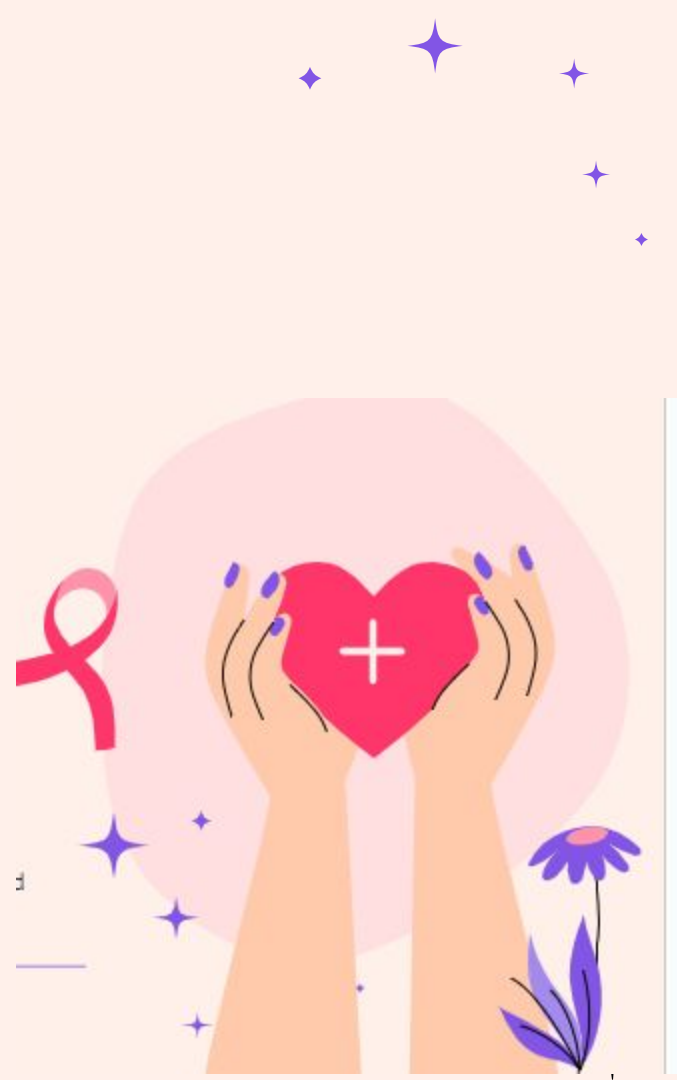
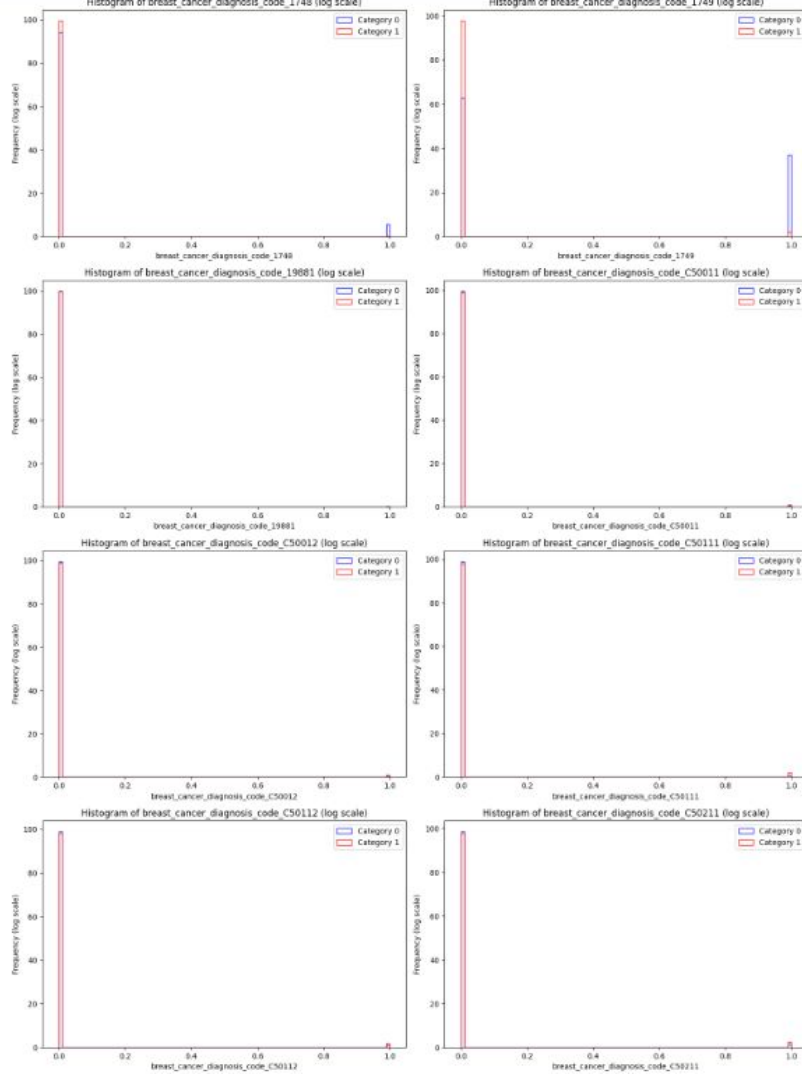


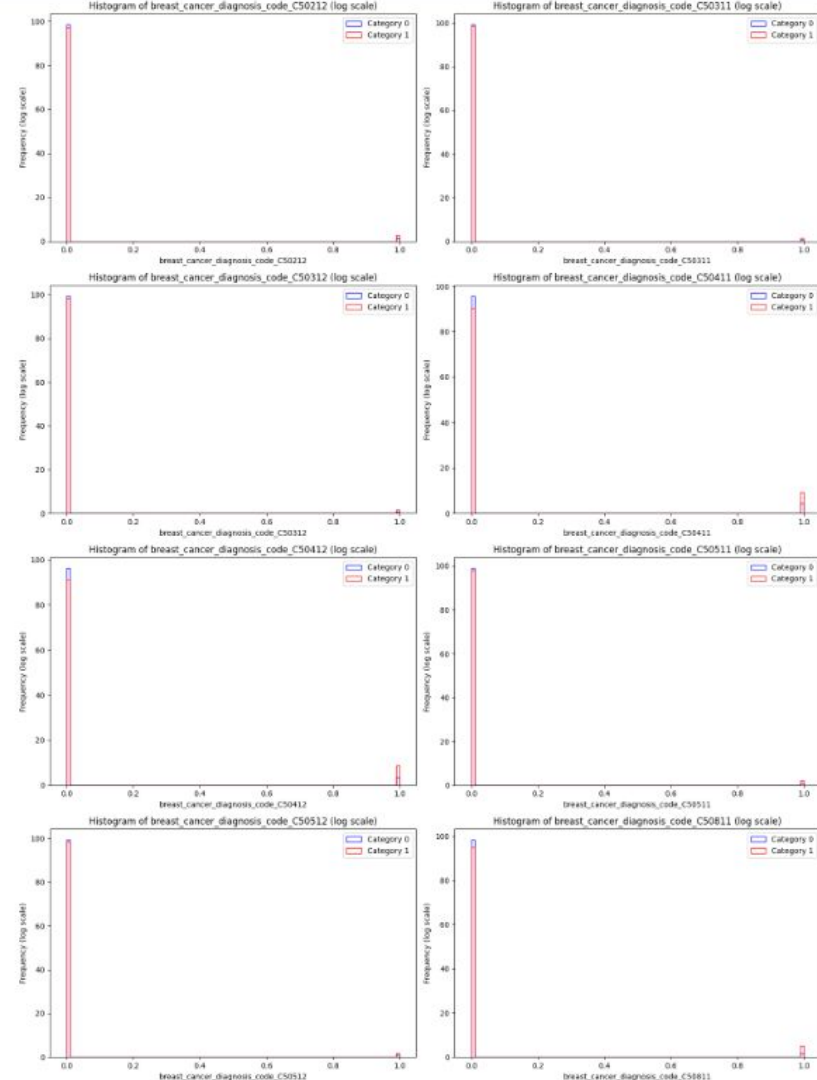
- ❖ Numerical and categorical datasets
- ❖ These Histograms are showcasing the distribution of the target variable in each feature

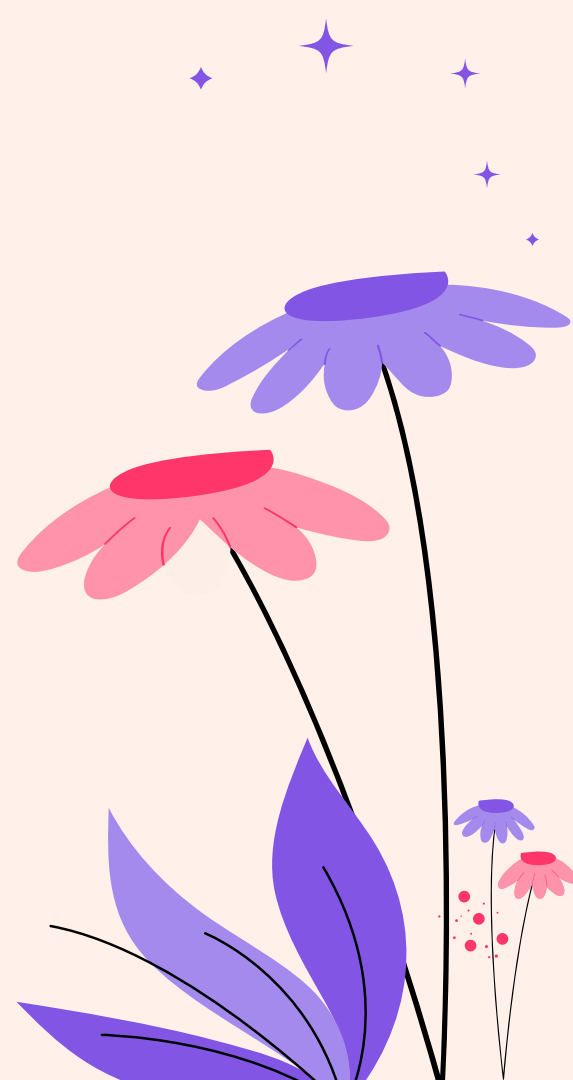
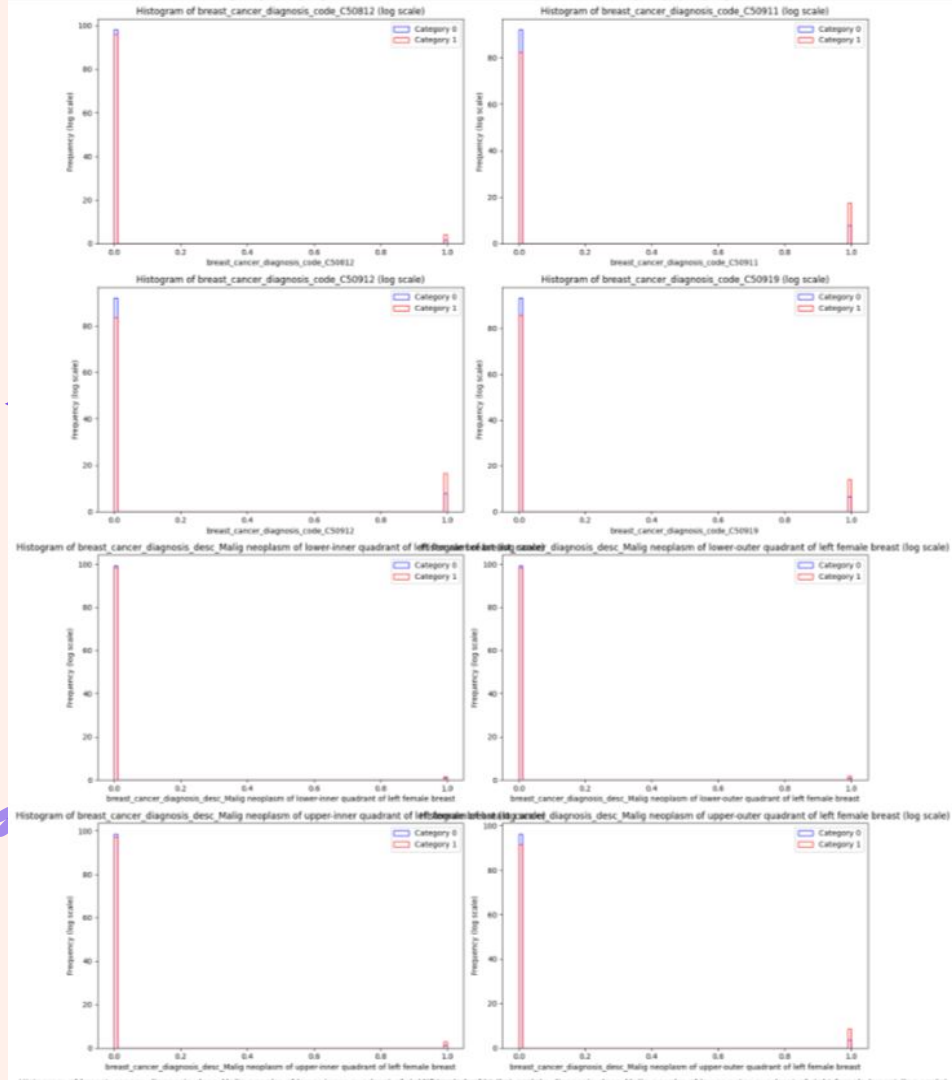


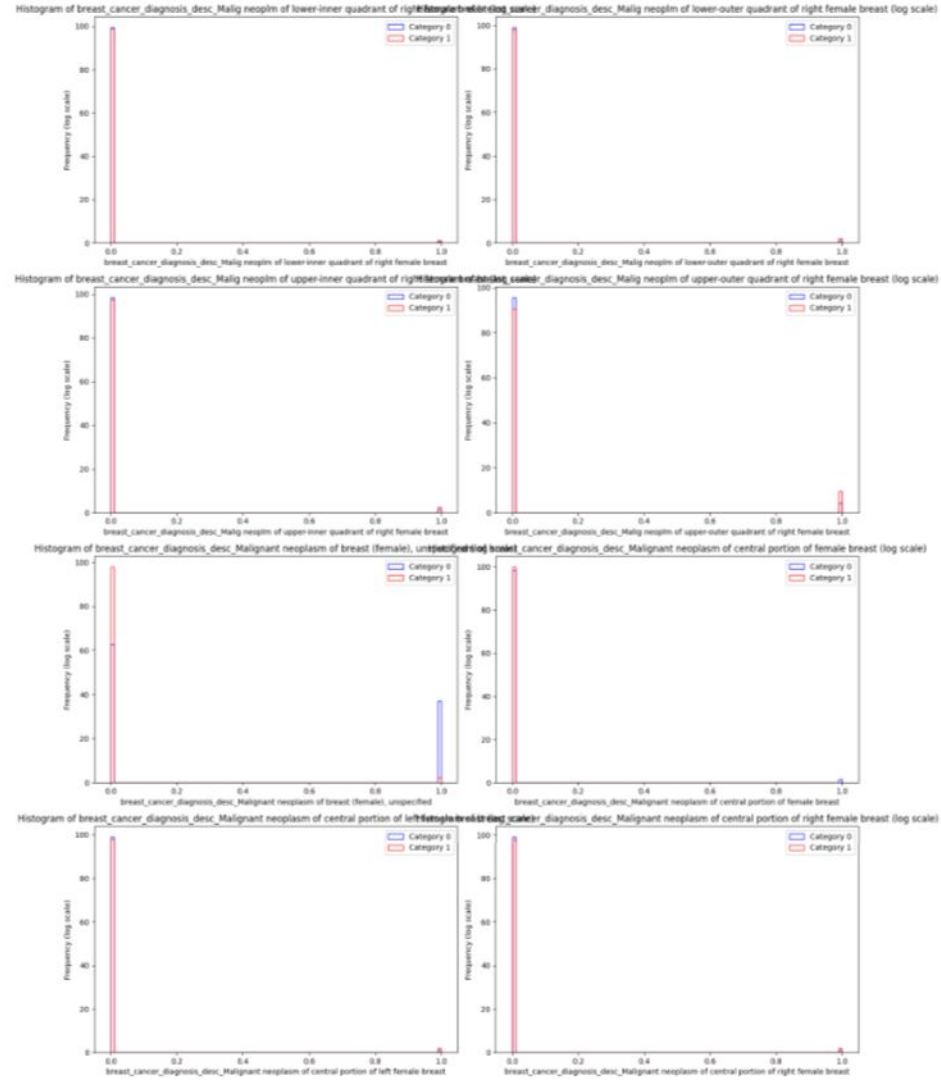


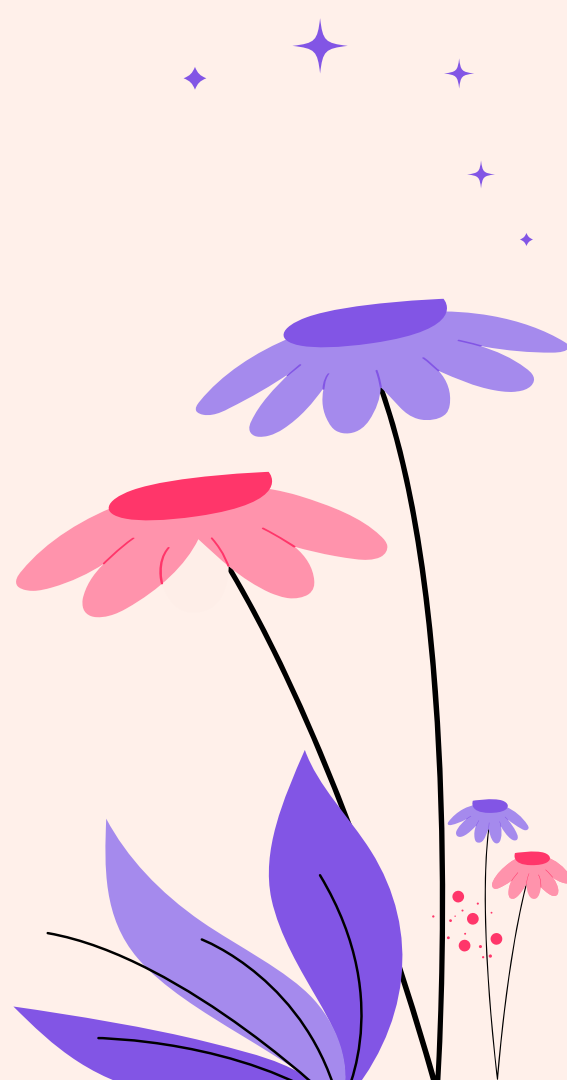
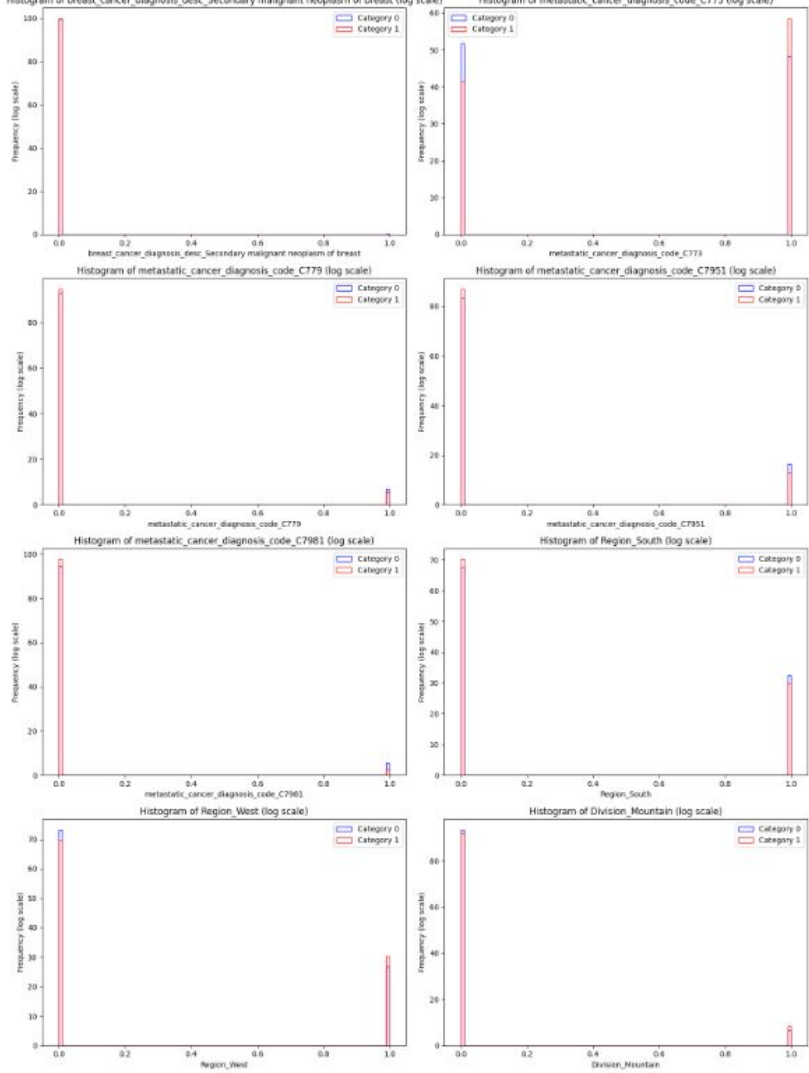




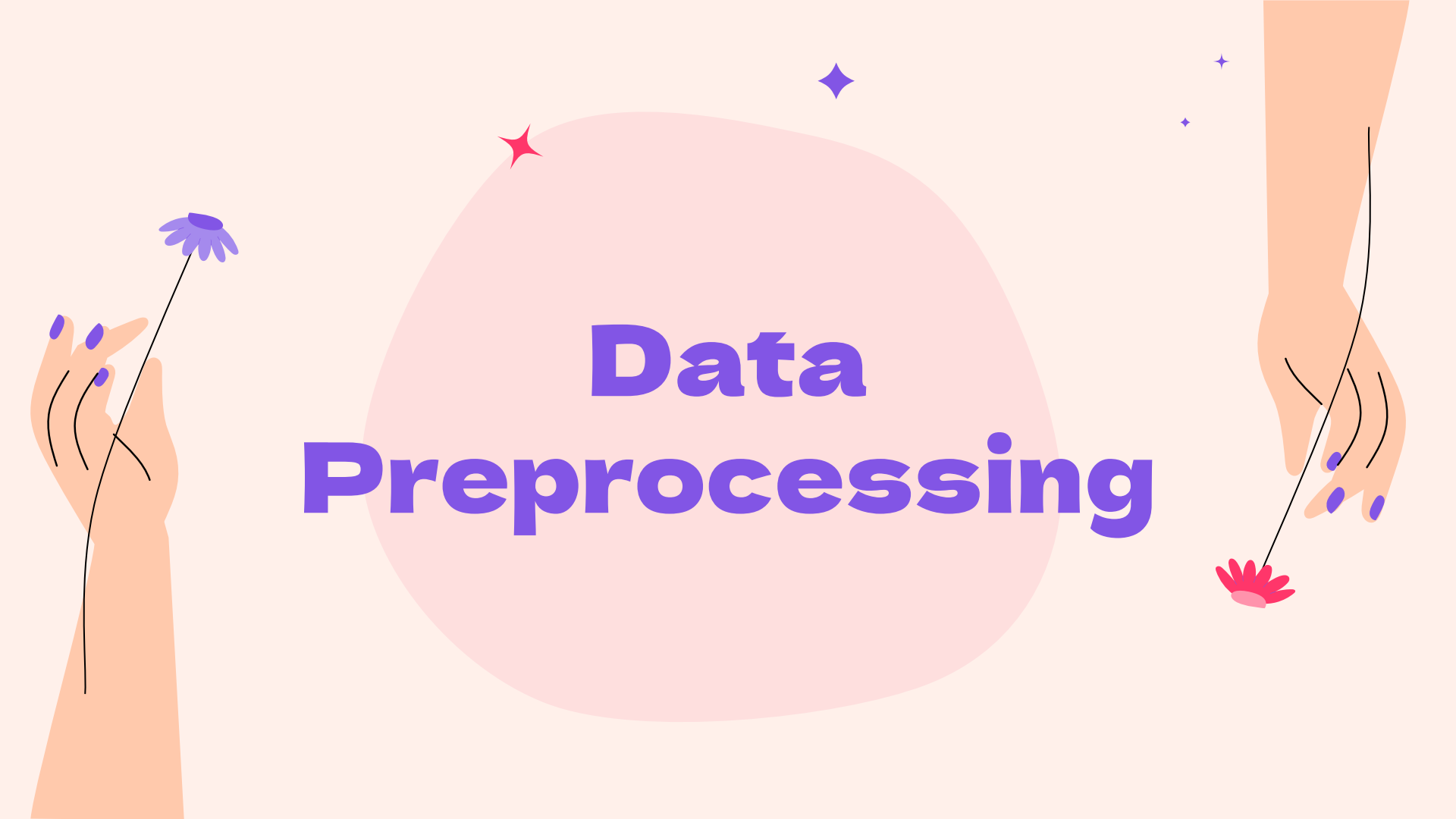






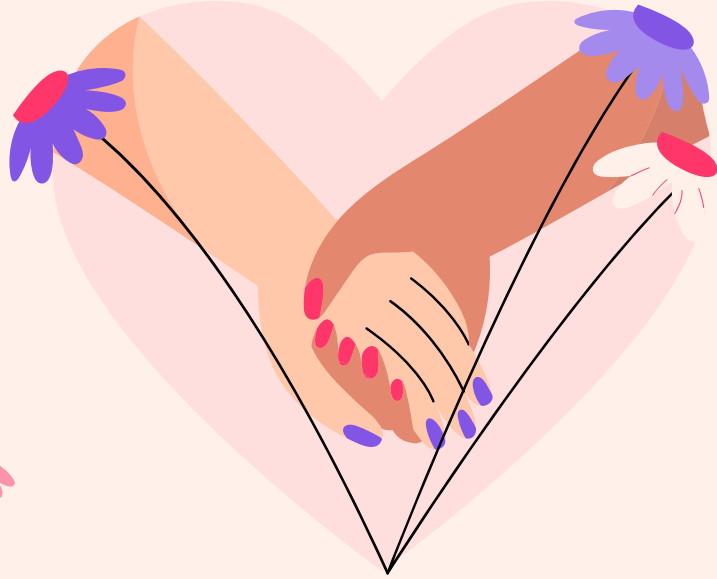




An illustration featuring two hands, one on the left and one on the right, both with orange skin and purple-painted fingernails. The left hand holds a purple flower on a thin black stem, while the right hand holds a pink flower on a similar stem. In the center, a large, light pink, rounded shape contains the text 'Data Preprocessing' in a bold, purple, sans-serif font. Several small, four-pointed starburst shapes are scattered around the central bubble: a red one at the top left, a purple one at the top center, and two small purple ones at the top right.

# Data Preprocessing

# DATA PREPROCESSING



**Missing Data Handled**



**Outliers handled with log transformation**

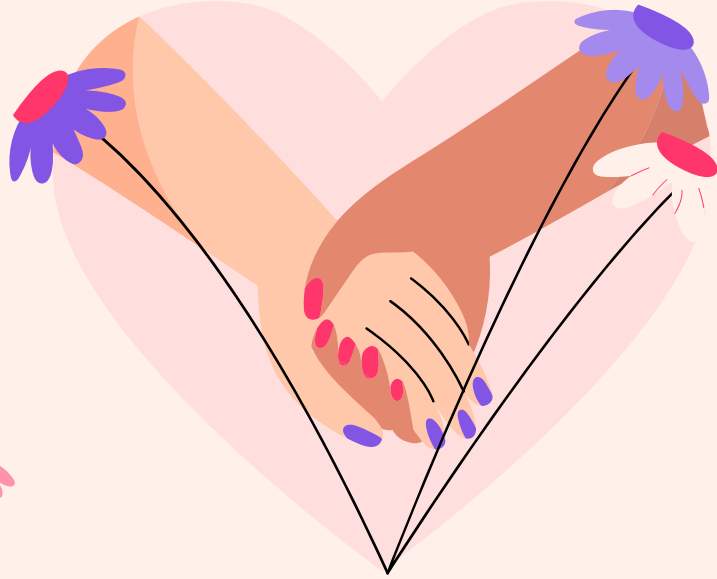


**One hot encoded the categorical variables**



**Dropped unnecessary columns like "ID"**

# DATA PREPROCESSING



**Using variance  
Thresholding for  
feature selections**



**Repeated the same  
steps for the testing  
data**

An illustration featuring two hands holding flowers. The hand on the left is holding a purple flower, and the hand on the right is holding a pink flower. Both hands have purple-painted fingernails. The background is a light pink color with a large, irregular pink shape in the center. There are also several small, four-pointed stars in purple and red scattered around the central shape.

# Model Building

# Logistics Regression Model



## Strengths

- Works well with both numerical and categorical data.
- Can handle large dataset with low computational cost
- Less susceptible to overfitting
- Provides easily interpretable results.



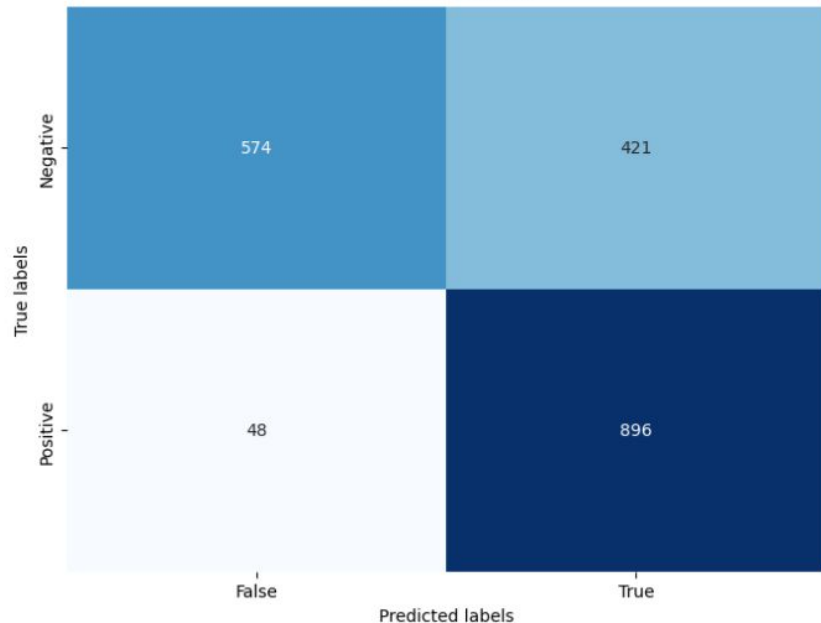
## Weaknesses

- May not capture complex relationship within the data.
- Sensitive to multicollinearity
- Assumes a linear relationship

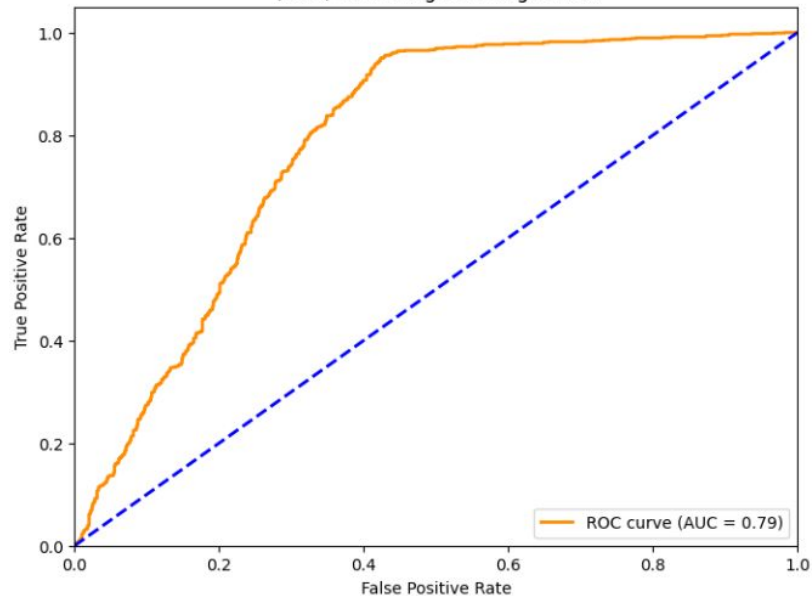
# Confusion Matrix & ROC Curve for Logistic Regression

Accuracy: 0.7581227436823105  
Precision: 0.680334092634776  
Recall: 0.9491525423728814  
F1-score: 0.7925696594427245

Confusion Matrix



(ROC) Curve-Logistics Regression



# Gradient Boosting Model



## Strengths

- Achieve high accuracy
- Handles complex relationships
- Helps Identify feature that contribute the most to the model's predictions
- Robust to Outliers



## Weaknesses

- Computationally expensive.
- Prone to overfitting
- Requires careful overfitting of Hyperparameters

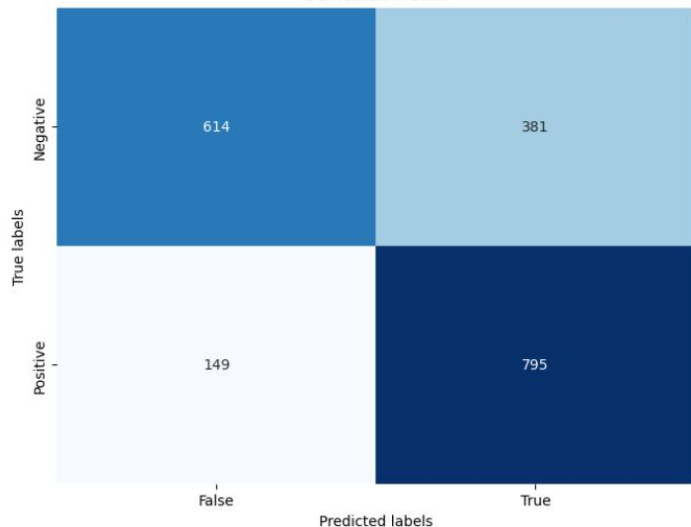
# Confusion Matrix & ROC Curve for XGBoost

Accuracy: 0.7266632284682826

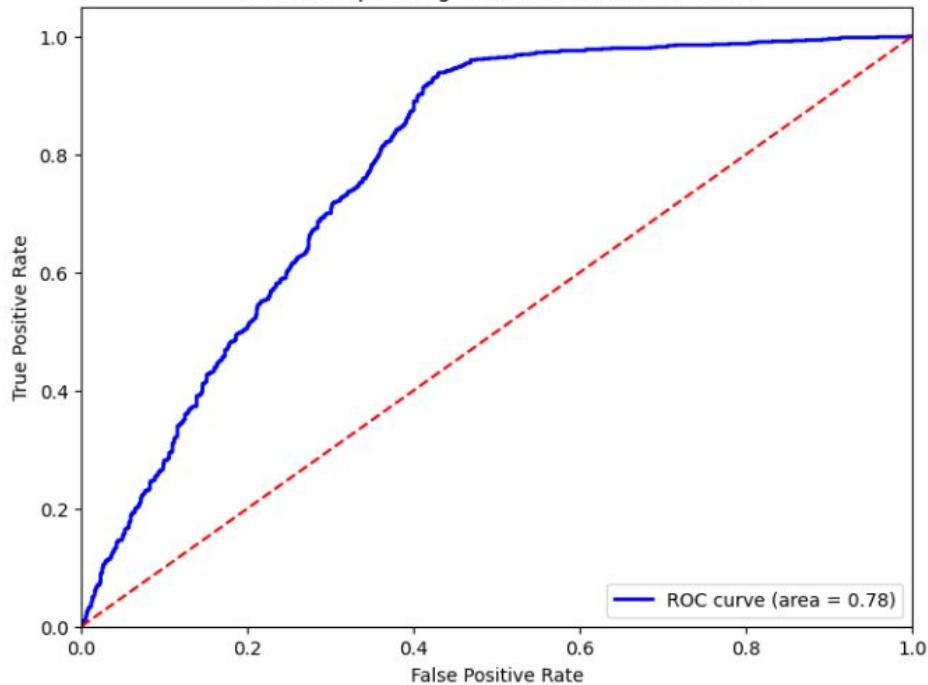
Classification Report:

	precision	recall	f1-score	support
0	0.80	0.62	0.70	995
1	0.68	0.84	0.75	944
accuracy			0.73	1939
macro avg	0.74	0.73	0.72	1939
weighted avg	0.74	0.73	0.72	1939

Confusion Matrix



Receiver Operating Characteristic (ROC) Curve







# Deep Learning

# FNN Model



## Strengths

- Can deal with non-linearity in a dataset
- Can handle numerical and categorical dataset
- Can handle large datasets



## Weaknesses

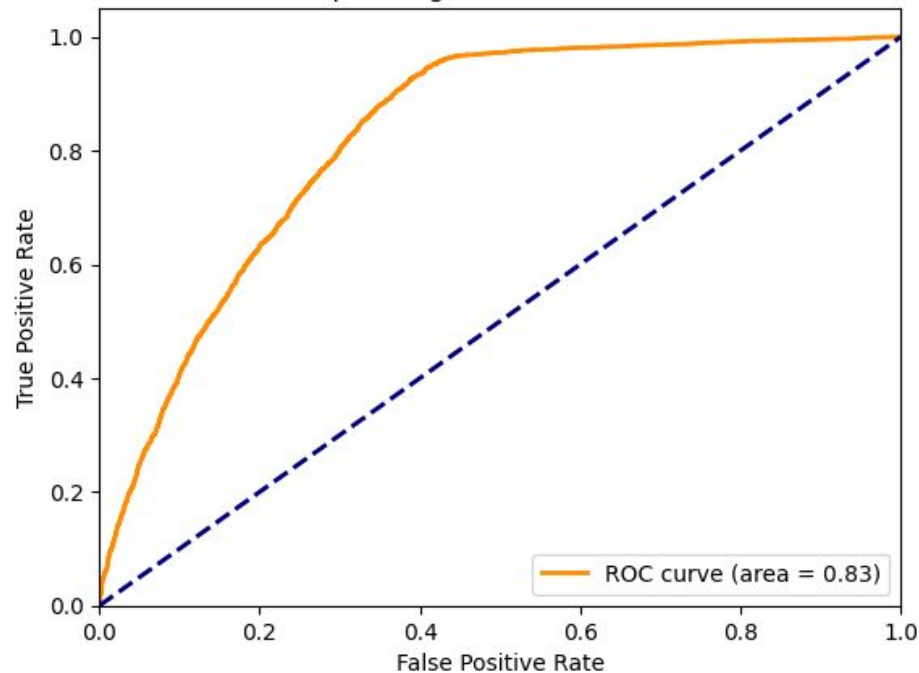
- Computationally expensive.
- Prone to overfitting

# ROC Curve of the FNN model

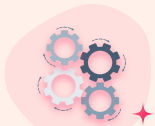
- FNN Model
- 10 Epochs

303/303 [=====] - 1s 4ms/step

Receiver Operating Characteristic (ROC) Curve



# MODEL EVALUATION



## Logistics Regression

Accuracy: 0.7581227436823105  
Precision: 0.680334092634776  
Recall: 0.9491525423728814  
F1-score: 0.7925696594427245

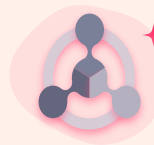


## XGBoost

Accuracy: 0.7266632284682826

Classification Report:

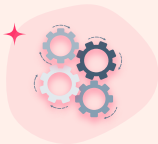
	precision	recall	f1-score
0	0.80	0.62	0.70
1	0.68	0.84	0.75



## FNN Model

Accuracy: 0.7407140135765076  
Precision: 0.7372381687164307  
Recall: 0.7480396032333374

## Kaggle Score of Each Model on the Test data



submission\_lr.csv

Complete (after deadline) · 14...

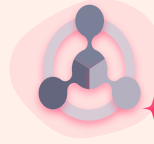
0.718



submission\_xgb1.csv

Complete (after deadline) · 15...

0.733



submission\_fnn.csv

Complete (after deadline) · 14...

0.787

# MODEL EVALUATION



- All of my model performed well, all their accuracy ranging in 70-80%, though the best model out of all of them was my deep learning model with a score of 0.74 on the training data & 0.78 on the test data.
- Though my model performed well, It definitely can be improved. I might be able to achieve that by scaling my datasets, and selecting more features for my models.





**THE END**

