**Women in Data Science Worldwide** | Afghanistan

**WOMEN in tech** ®
Afghanistan

# Upskill Workshop I

June 24, 2024

# Table of Content

# Technical Part

**Environment:**
- Google Colab

**Language:**
- Python

**Libraries:**
- numpy
- pandas
- scipy.stats
- matplotlib.pyplot
- seaborn
- math
- re
- nltk
- wordcloud

# WiDS Dataset 2024 Overview

**Dataset Source:**

- Provided by Gilead Sciences through Health Verity, featuring patients diagnosed with metastatic triple negative breast cancer in the US.

**Dataset Enrichment:**

- Includes demographics (age, race, BMI), diagnosis details (cancer codes, treatment), and insurance information.
- Geo-demographic data enriched from US Zip Codes Database for socio-economic insights.
- Zip code level climate data added to explore climate's impact on healthcare access.
- Target Column: The dataset includes a target variable metastatic_diagnosis_period, which measures the period (in days) between the initial breast cancer diagnosis and the subsequent diagnosis of metastatic cancer.

# Dataset Availability

The dataset is provided in two sets:

- **Training Dataset:** Named train.csv, this dataset contains observed values of the metastatic diagnosis period for model training.

- **Testing Dataset:** Named test.csv, this dataset withholds the observed values and is used for evaluating model predictions.

Datasets are available at :
https://www.kaggle.com/competitions/widsdatathon2024-challenge2/data

# Datasets Size

**Train Data Shape:**
- Number of Rows: 13173
- Number of Columns: 152
- Numerical Columns: 141
- Categorical Columns: 11

**Test Data Shape:**
- Number of Rows: 5646
- Number of Columns: 151
- Numerical Columns: 140
- Categorical Columns: 11

```python
# Get the shape of the DataFrame
shape_train = train.shape

# Print the shape
print("Shape of train data:", shape_train)
```

```
Shape of train data: (13173, 152)
```

```python
# Get the shape of the DataFrame
shape_test = test.shape

# Print the shape
print("Shape of test data:", shape_test)
```
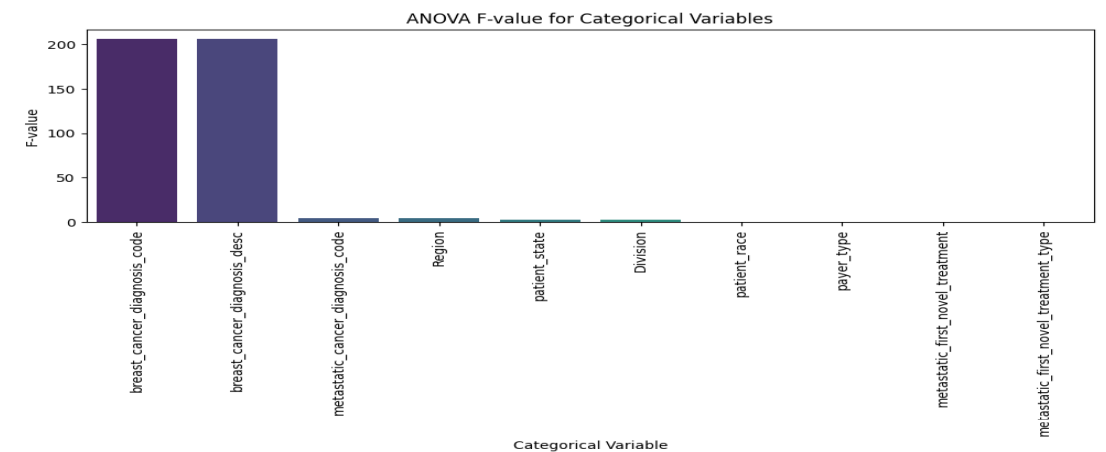
```
Shape of test data: (5646, 151)
```
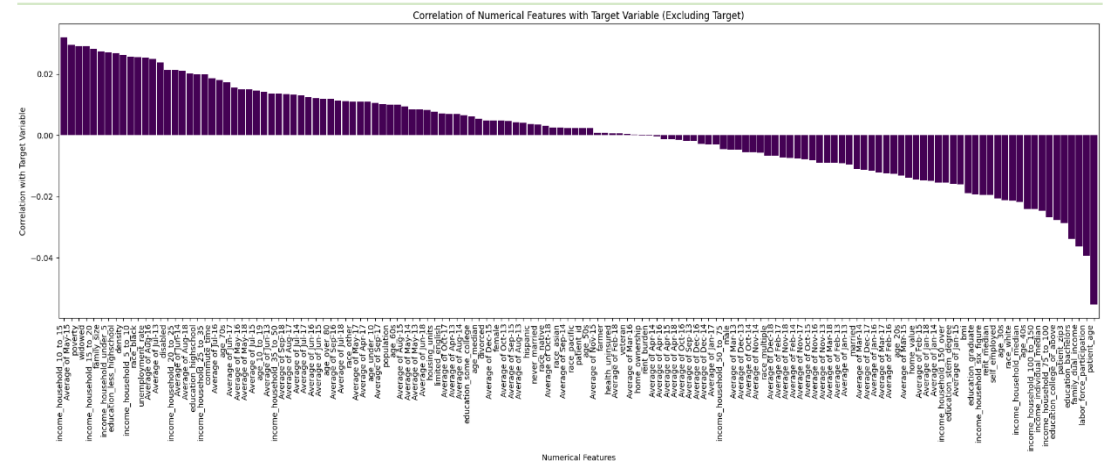
# Analyzing Relationships with the Target Variable

**Numerical Columns**

- Pearson correlation coefficient used to measure linear correlation with the target variable.

**Categorical Columns:**

- One-way ANOVA employed to assess statistical significance of categorical variables concerning the target variable.



Correlation of Numerical Features with Target Variable (Excluding Target)



ANOVA F-value for Categorical Variables

# Identifying the Most Related Columns to the Target Variable

## Subset Selection Method:

- Utilized Pearson correlation for numerical columns to identify the top 10 most correlated features with the target variable.

- Employed one-way ANOVA for categorical columns to select the 7 most relevant features in relation to the target variable.

```python
# Exclude the target variable and sort by absolute correlation values
corr_with_target = corr_matrix[target_variable].drop(target_variable).abs().sort_values(ascending=False).head(10)

# Print top 10 absolute correlation values (excluding the target variable)
print(f"Top 10 absolute correlations with '{target_variable}' (excluding itself):")
for col, corr_val in corr_with_target.items():
    print(f"{col}: {corr_val}")
```

```
Top 10 absolute correlations with 'metastatic_diagnosis_period' (excluding itself):
patient_age: 0.055231359481934965
labor_force_participation: 0.039368608644549864
family_dual_income: 0.03640348751137413
education_bachelors: 0.033842293355421896
income_household_10_to_15: 0.03197485109301944
Average of May-15: 0.029483726364884503
poverty: 0.029135428404253472
widowed: 0.029086316867978728
patient_zip3: 0.028686057480996174
income_household_15_to_20: 0.028111285426823878
```

```python
# Sort by F-value in descending order and get the top 7
top_anova_df = anova_df.sort_values(by='f_value', ascending=False).head(7)

# Print top 7 categorical variables by ANOVA F-value
print("Top 7 categorical variables by ANOVA F-value:")
print(top_anova_df)
```

```
Top 7 categorical variables by ANOVA F-value:
                categorical_variable      f_value
5           breast_cancer_diagnosis_code  206.133897
6           breast_cancer_diagnosis_desc  206.133897
7       metastatic_cancer_diagnosis_code    4.731484
3                                 Region    4.696268
2                          patient_state    2.920630
4                               Division    2.823486
0                           patient_race         NaN
```

# Descriptions of Columns in the Provided Subset

- patient_age: Age of the patient derived from their year of birth.
- labor_force_participation: Percentage of residents aged 16 and older participating in the labor force.
- family_dual_income: Percentage of families with dual income earners.
- education_bachelors: Percentage of residents with a bachelor's degree (or equivalent) but no more.
- Average of May-15: Average temperature for the patient's zip code in May 2015.
- poverty: Median value of owner-occupied homes.
- patient_zip3: Zip code of the patient's residence (first three digits).
- Region: Region of the patient's location.
- patient_state: State abbreviation of the patient's residence.
- Division: Division of the patient's location.
- patient_race: Race of the patient.
- breast_cancer_diagnosis_code: ICD10 or ICD9 diagnoses code for breast cancer.
- breast_cancer_diagnosis_desc: Description of the breast cancer diagnosis code.
- metastatic_cancer_diagnosis_code: ICD10 diagnoses code for metastatic cancer.
- metastatic_diagnosis_period: Period (in days) between breast cancer diagnosis and metastatic cancer diagnosis.

# Before We Start…

Checklist for Dataset Analysis

**Key Questions:**

What specific **questions** should this dataset answer?

**Tracking Parameters:**

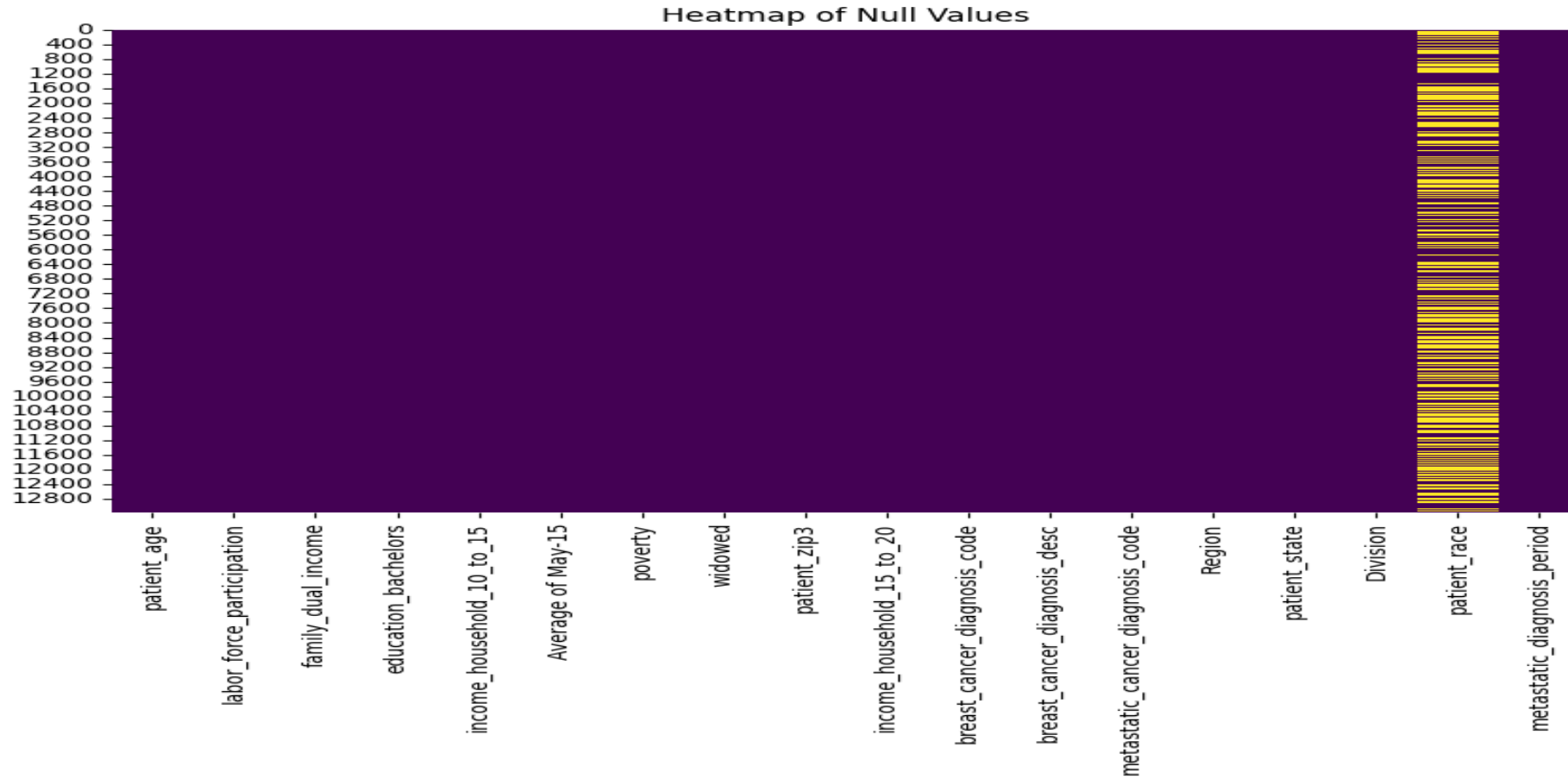Which **parameters** do we need to monitor to answer these questions?
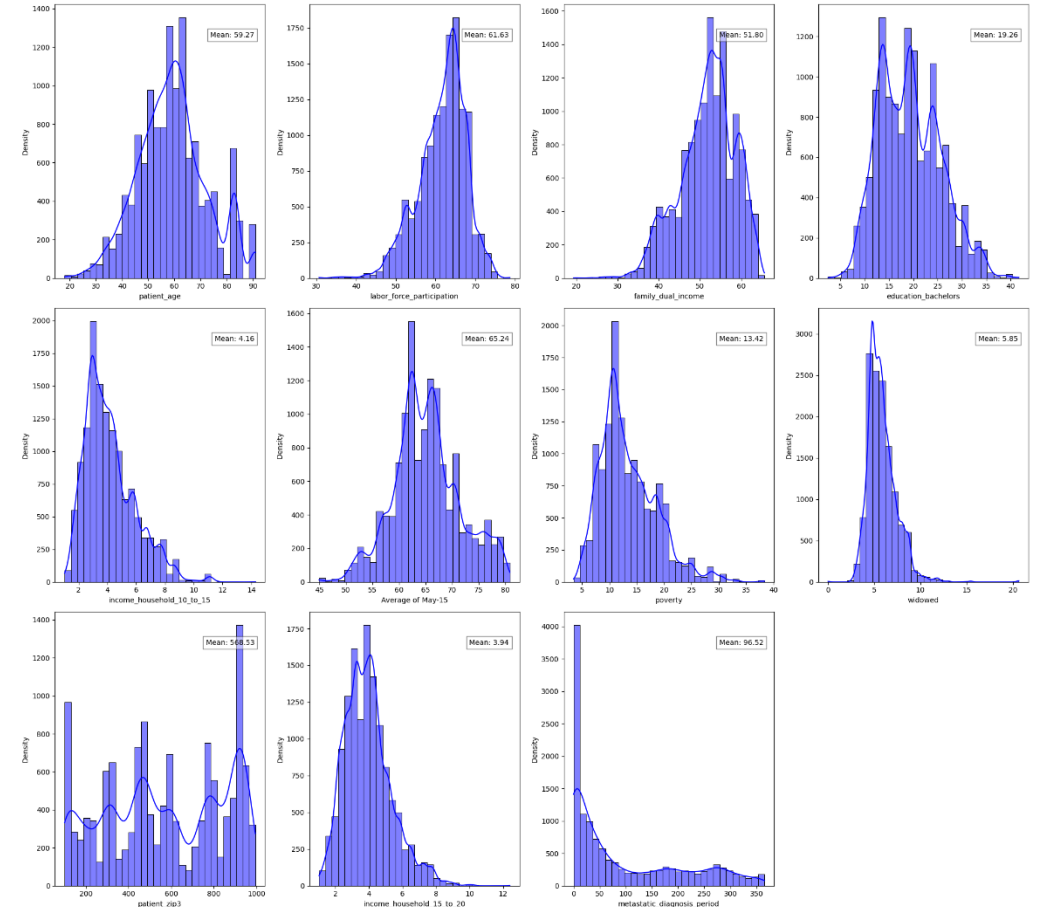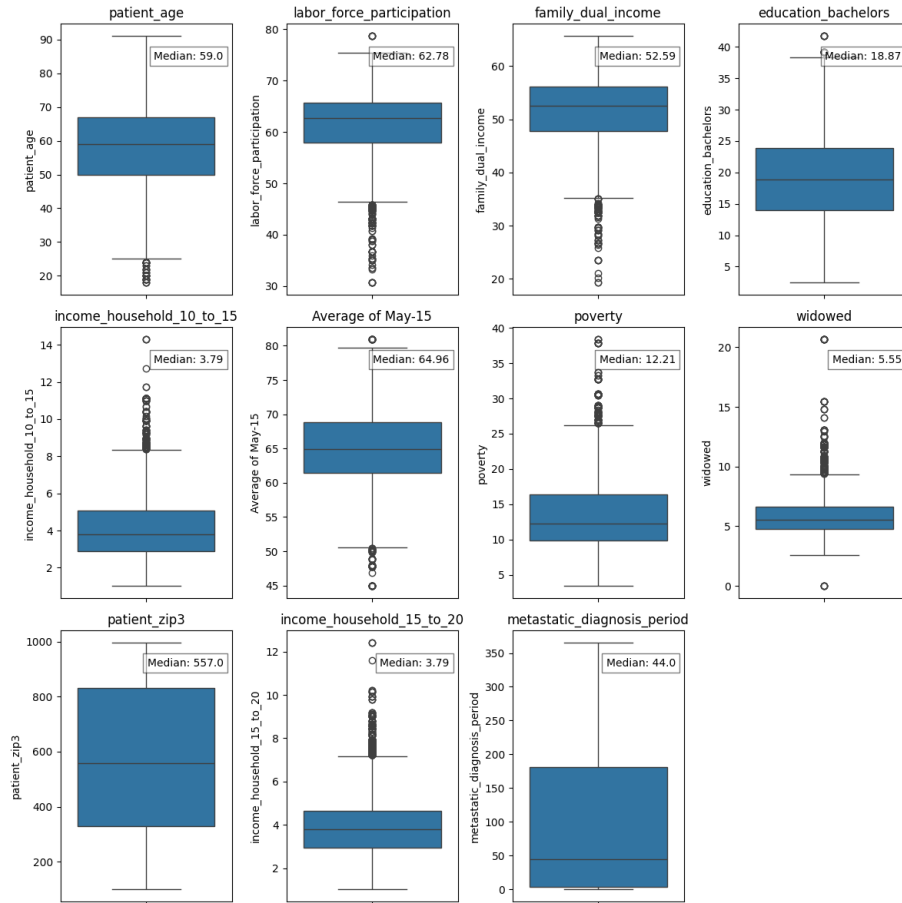
# Summary Statics

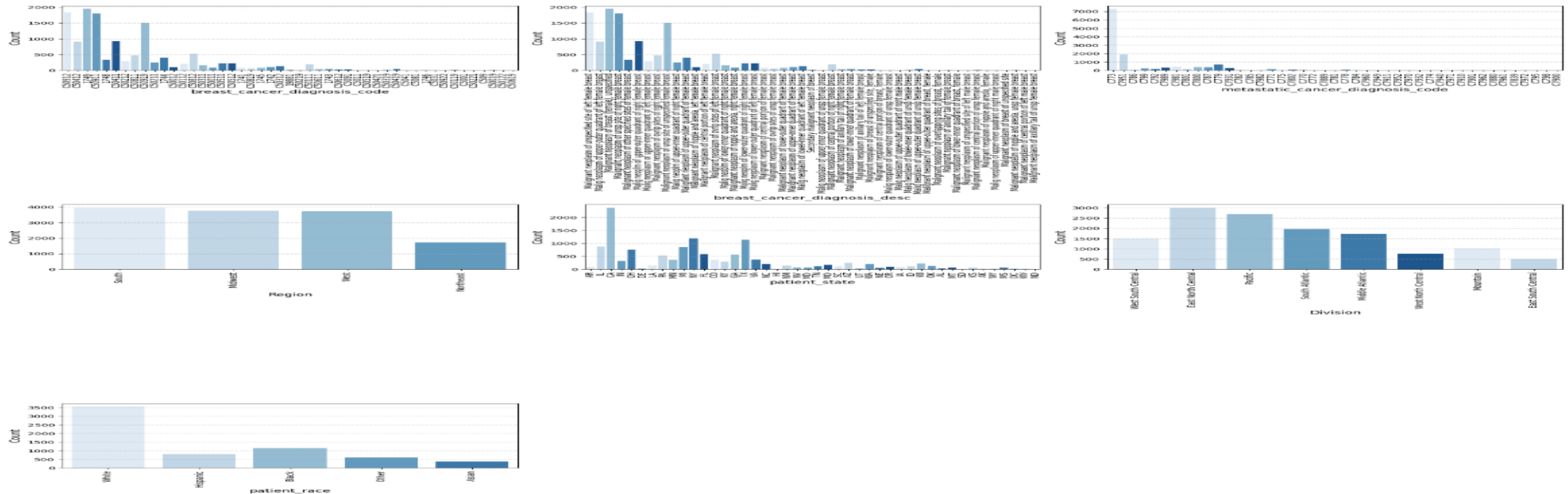| | type | count | nunique | unique% | null | null% | mean | std | min | 25% | 50% | 75% | max |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| patient_age | int64 | 13173 | 67 | 0.508616 | 0 | 0.000000 | 59.271313 | 13.218883 | 18.000000 | 50.000000 | 59.000000 | 67.000000 | 91.000000 |
| patient_zip3 | int64 | 13173 | 751 | 5.701055 | 0 | 0.000000 | 568.530859 | 275.758485 | 100.000000 | 330.000000 | 557.000000 | 832.000000 | 995.000000 |
| metastatic_diagnosis_period | int64 | 13173 | 366 | 2.778410 | 0 | 0.000000 | 96.515221 | 108.969873 | 0.000000 | 3.000000 | 44.000000 | 181.000000 | 365.000000 |
| labor_force_participation | float64 | 13173 | 656 | 4.979883 | 0 | 0.000000 | 61.633658 | 5.977344 | 30.700000 | 57.960000 | 62.780000 | 65.680000 | 78.670000 |
| education_bachelors | float64 | 13173 | 630 | 4.782510 | 0 | 0.000000 | 19.263585 | 6.255266 | 2.470000 | 13.980000 | 18.870000 | 23.890000 | 41.700000 |
| Average of May-15 | float64 | 13173 | 615 | 4.668640 | 0 | 0.000000 | 65.244507 | 6.306477 | 44.950000 | 61.460000 | 64.960000 | 68.800000 | 80.900000 |
| widowed | float64 | 13173 | 439 | 3.332574 | 0 | 0.000000 | 5.846155 | 1.556496 | 0.000000 | 4.770000 | 5.550000 | 6.610000 | 20.650000 |
| family_dual_income | float64 | 13168 | 666 | 5.055796 | 5 | 0.037956 | 51.800184 | 6.696196 | 19.310000 | 47.732500 | 52.590000 | 56.160000 | 65.640000 |
| income_household_10_to_15 | float64 | 13168 | 451 | 3.423670 | 5 | 0.037956 | 4.159681 | 1.751091 | 1.020000 | 2.900000 | 3.790000 | 5.090000 | 14.280000 |
| poverty | float64 | 13168 | 615 | 4.668640 | 5 | 0.037956 | 13.417748 | 5.105035 | 3.430000 | 9.870000 | 12.210000 | 16.410000 | 38.350000 |
| income_household_15_to_20 | float64 | 13168 | 430 | 3.264253 | 5 | 0.037956 | 3.943212 | 1.402426 | 1.030000 | 2.940000 | 3.790000 | 4.640000 | 12.400000 |
| breast_cancer_diagnosis_code | object | 13173 | 47 | 0.356790 | 0 | 0.000000 | nan | nan | nan | nan | nan | nan | nan |
| breast_cancer_diagnosis_desc | object | 13173 | 47 | 0.356790 | 0 | 0.000000 | nan | nan | nan | nan | nan | nan | nan |
| metastatic_cancer_diagnosis_code | object | 13173 | 43 | 0.326425 | 0 | 0.000000 | nan | nan | nan | nan | nan | nan | nan |
| Region | object | 13173 | 4 | 0.030365 | 0 | 0.000000 | nan | nan | nan | nan | nan | nan | nan |
| patient_state | object | 13173 | 44 | 0.334017 | 0 | 0.000000 | nan | nan | nan | nan | nan | nan | nan |
| Division | object | 13173 | 8 | 0.060730 | 0 | 0.000000 | nan | nan | nan | nan | nan | nan | nan |
| patient_race | object | 6516 | 5 | 0.037956 | 6657 | 50.535186 | nan | nan | nan | nan | nan | nan | nan |

# Heatmap of Missing Values

# Distribution of Numerical Data



**How are the numerical variables distributed across the dataset, and what insights can be derived from their distributions?**
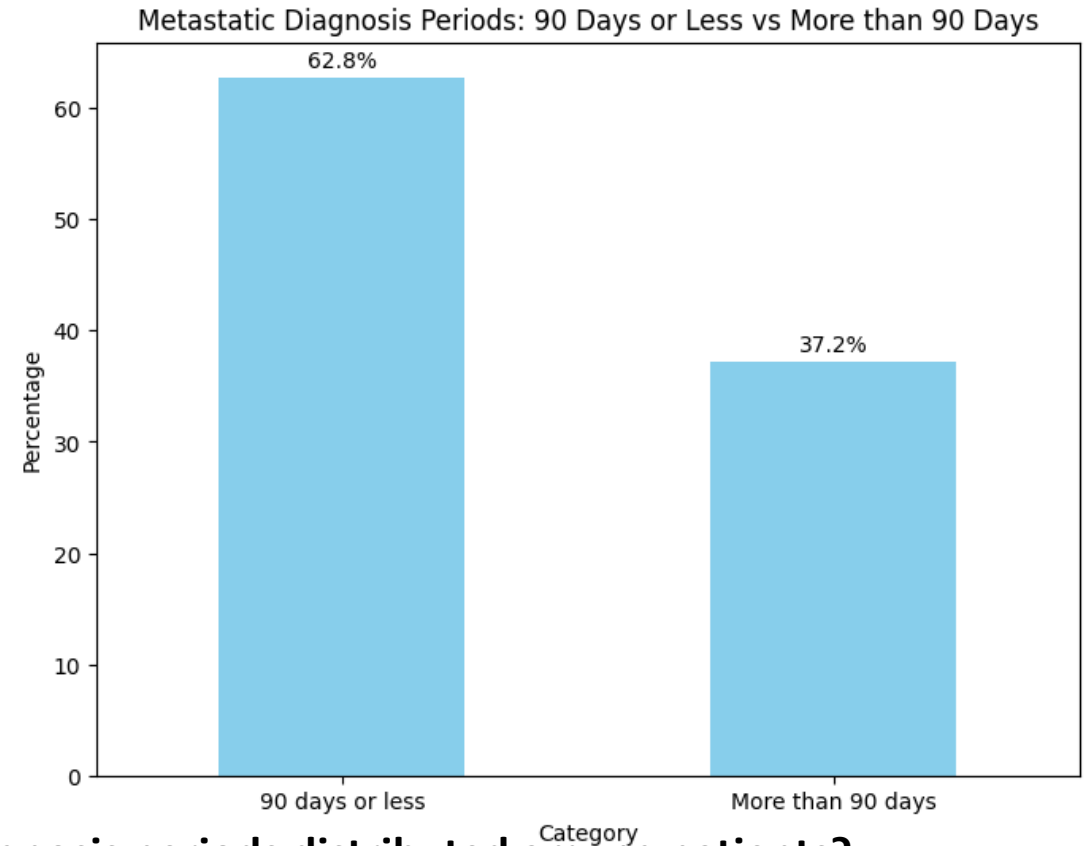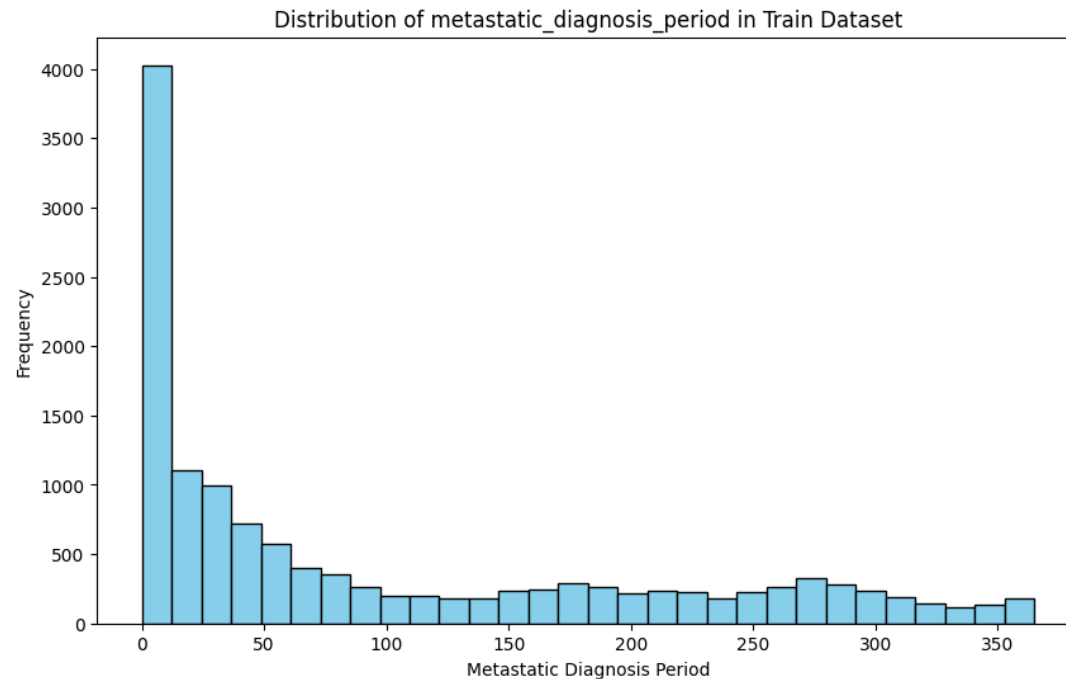
# Bar Chart of Categorical Data



**How do different categories within each categorical variable compare in terms of their frequency or distribution?**
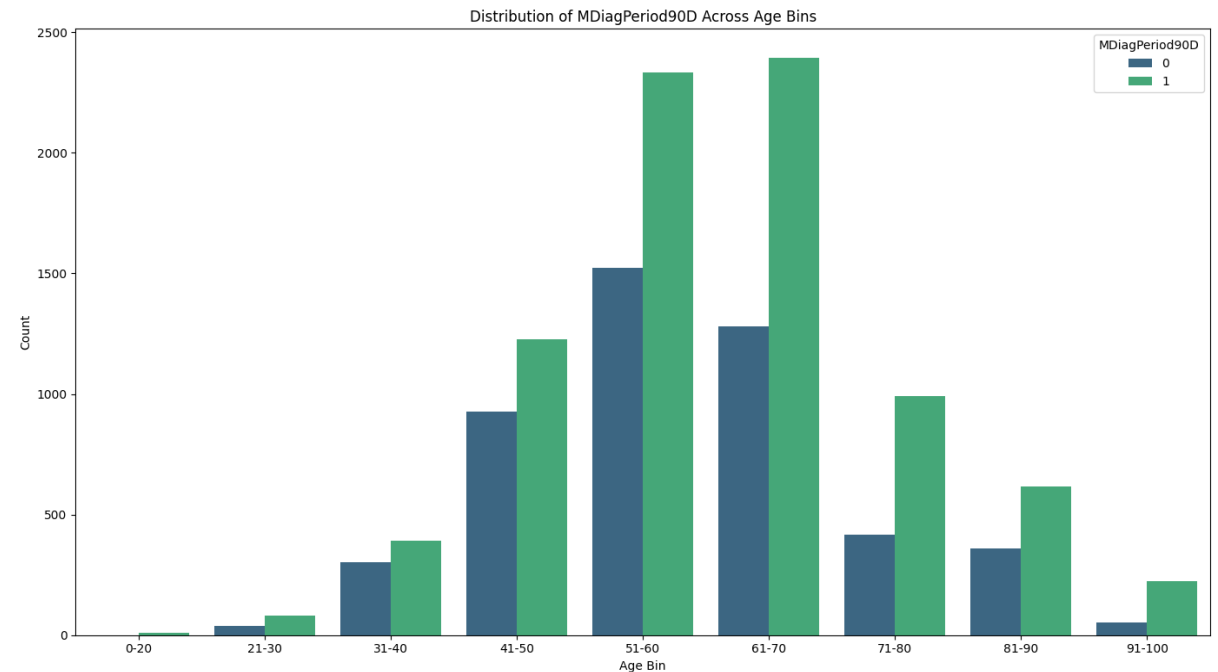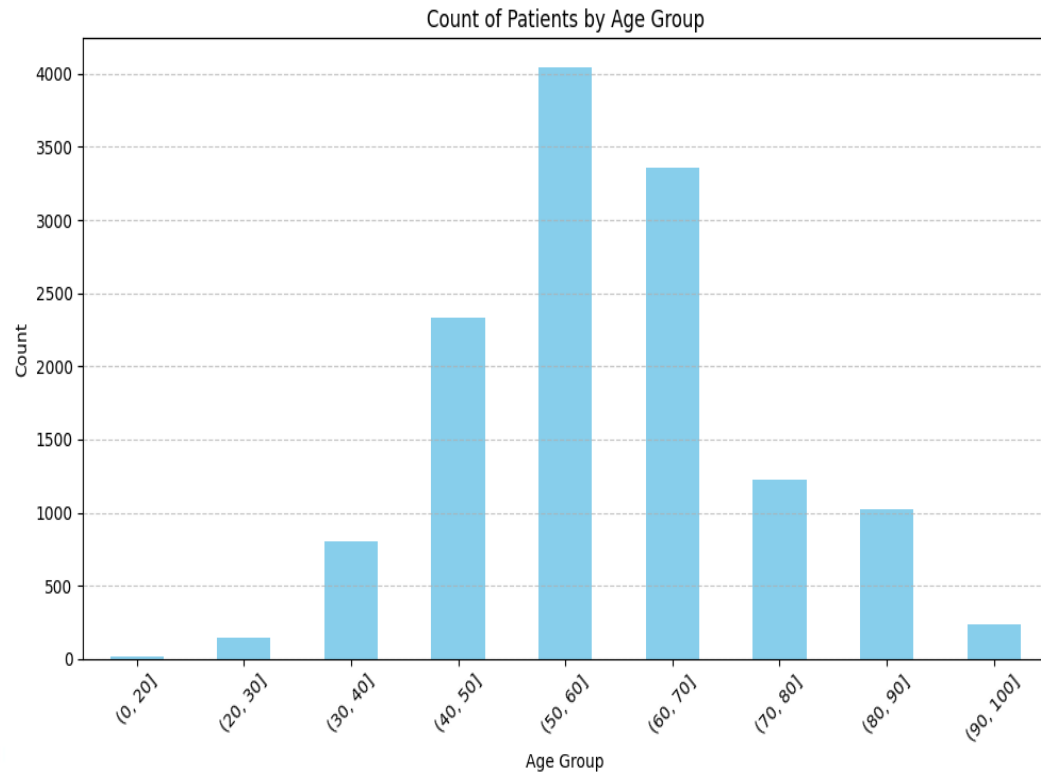
# Analysis of Target Column



- **How is the frequency distribution of metastatic diagnosis periods distributed among patients?**
- **How is the distribution of metastatic diagnosis periods categorized as 90 days or less versus more than 90 days among patients?**

# Analysis of Age

- **What is the distribution of patients across different age groups in the dataset?**
- **How does the distribution of `MDiagPeriod90D` (Metastatic Diagnosis Period within 90 days) vary across different age groups?**
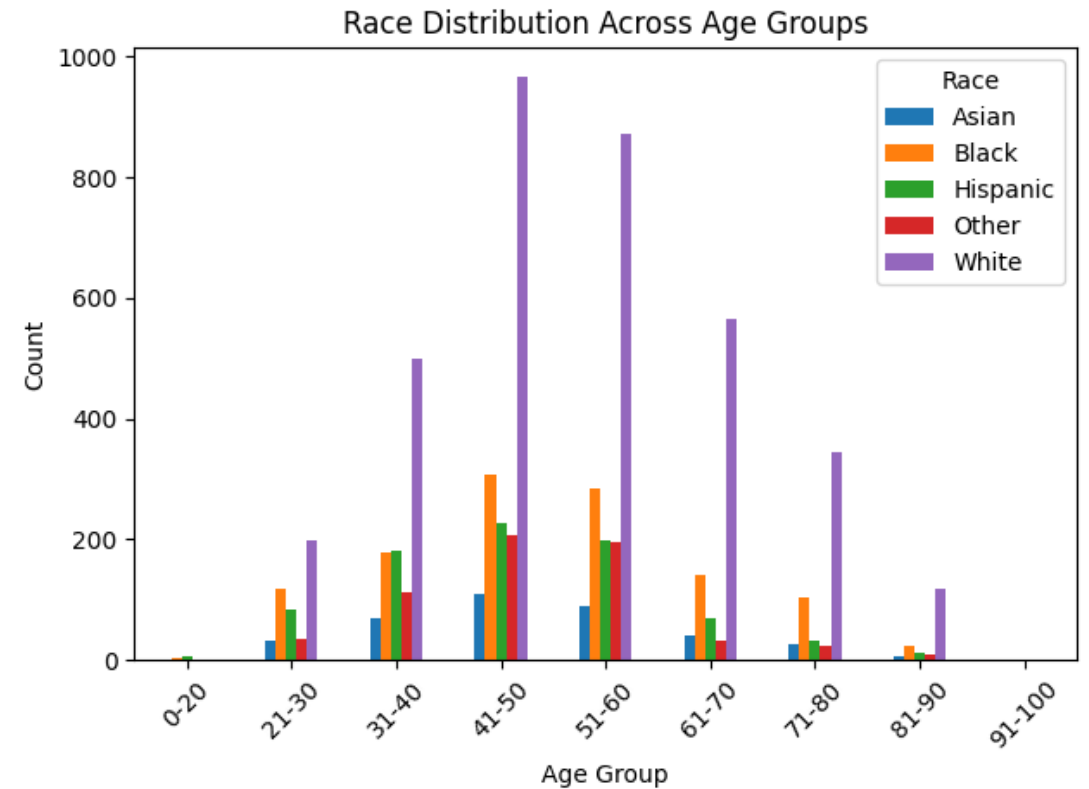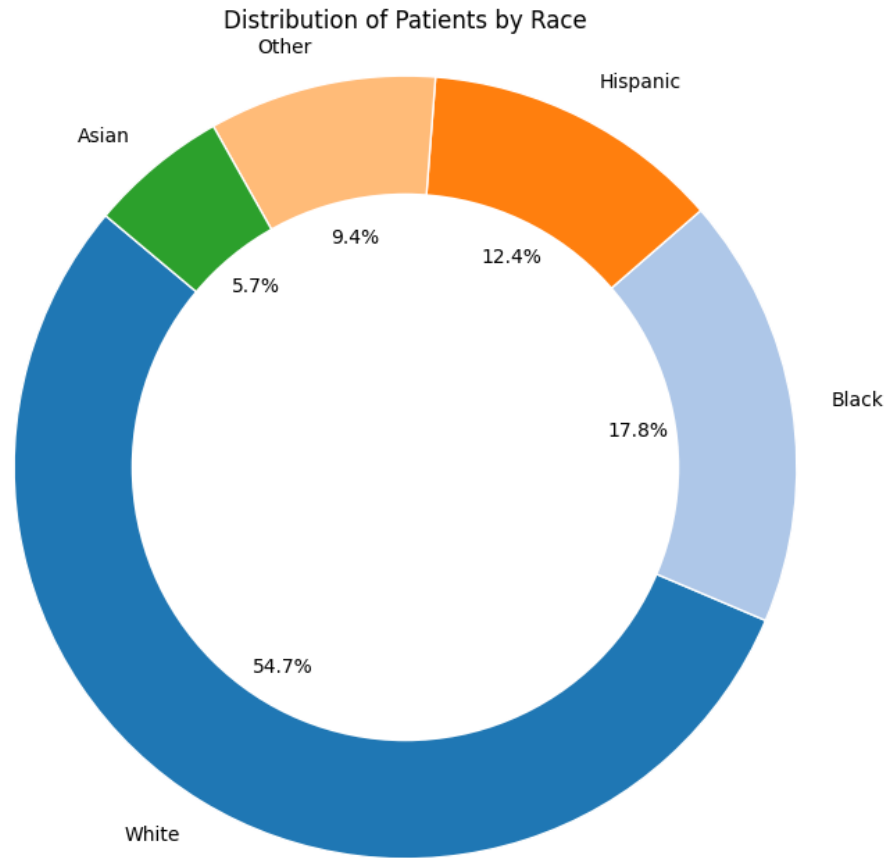
# Analysis of Region



Distribution of Patients by Region



Density of Patients Across Regions

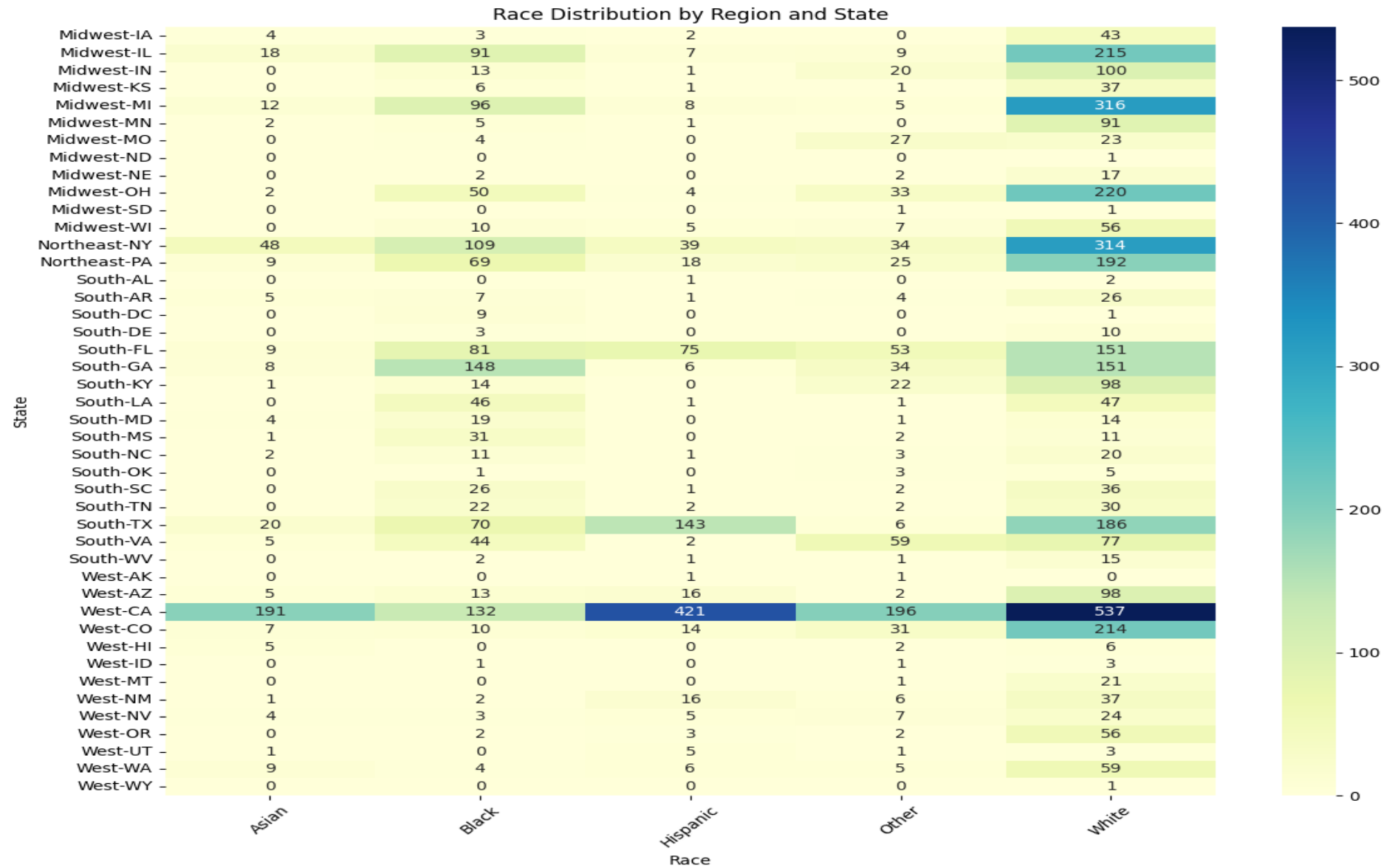**How are patients distributed across various regions?**

# Analysis of Patient Race



- **What is the distribution of patients across different racial groups in the dataset?**
- **How does the distribution of patient races vary across different age groups?**
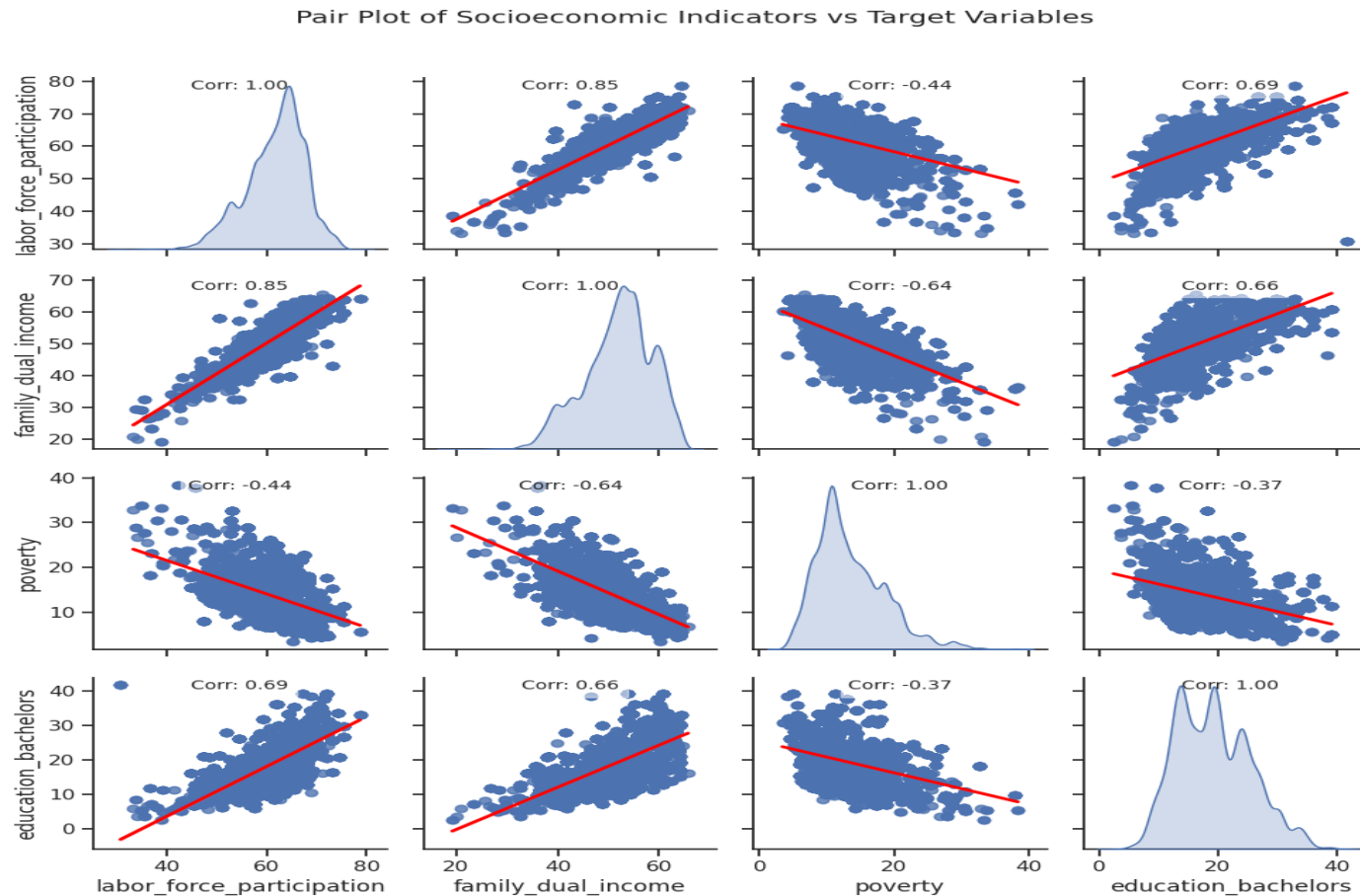
# Cont.



Race Distribution by Region and State

**How does the distribution of patient races vary across different regions and states?**
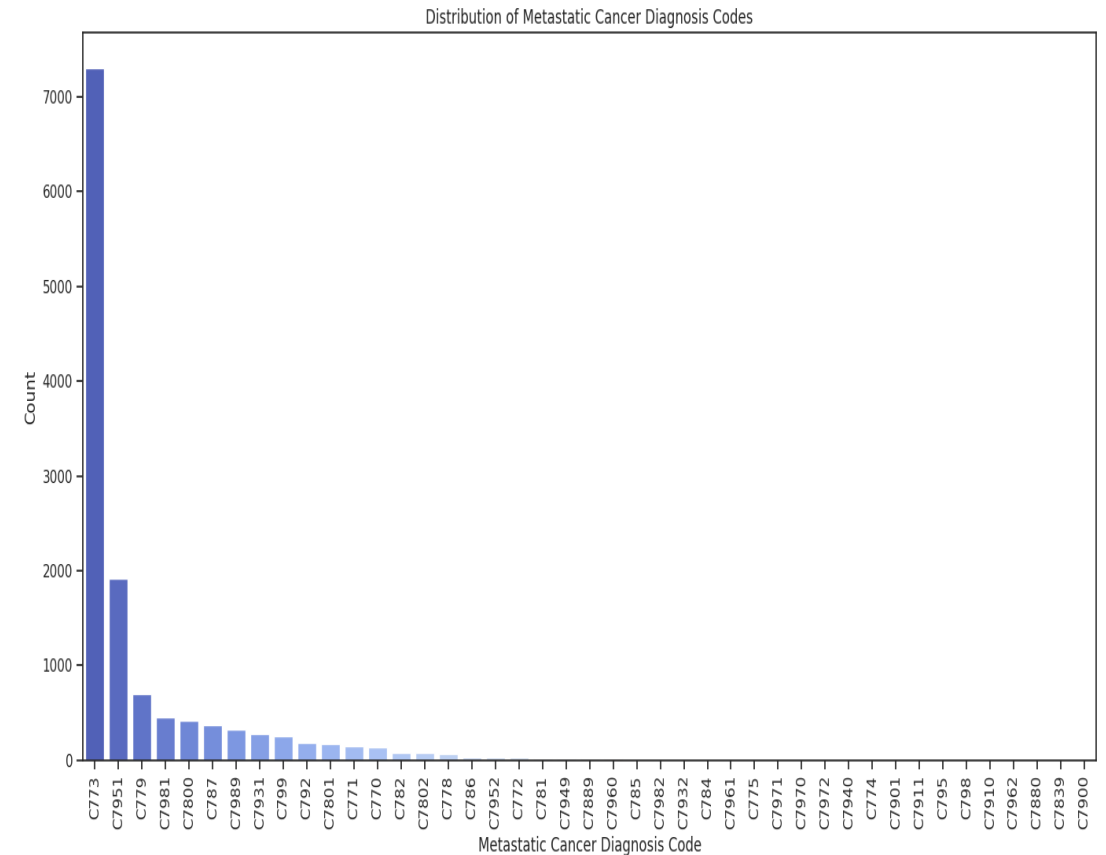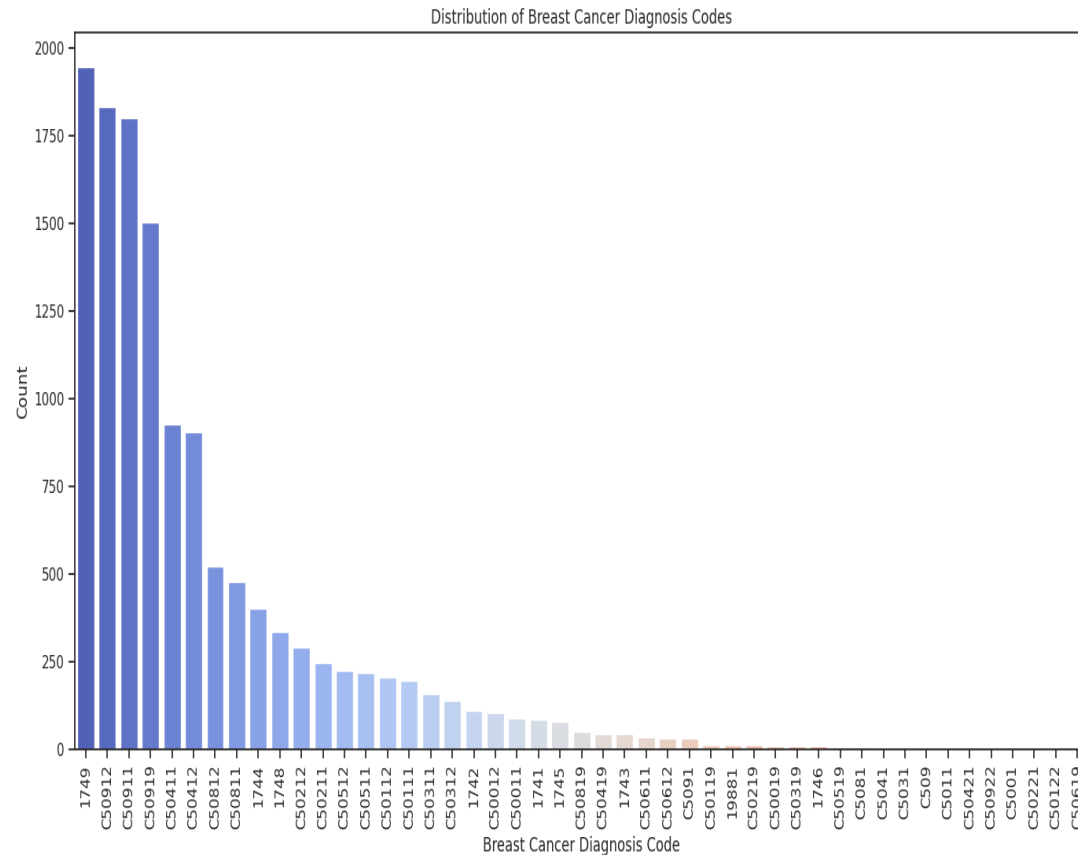
# Socioeconomic Indicators



Pair Plot of Socioeconomic Indicators vs Target Variables

How can we explore the relationships among socioeconomic indicators?

# Analysis of Cancer Codes Data



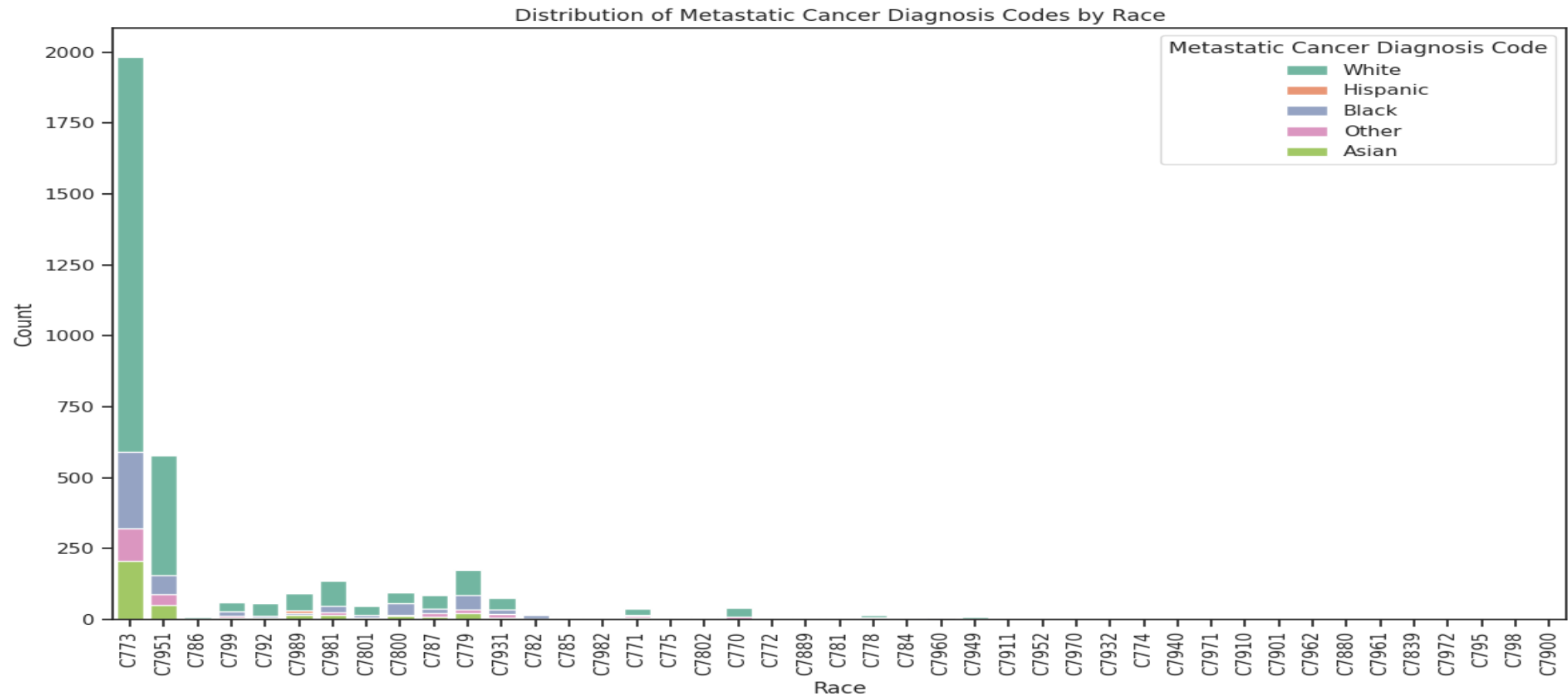Distribution of Breast Cancer Diagnosis Codes

Distribution of Metastatic Cancer Diagnosis Codes

**What are the distributions of breast cancer diagnosis codes and metastatic cancer diagnosis codes in the dataset?**

# Cont.



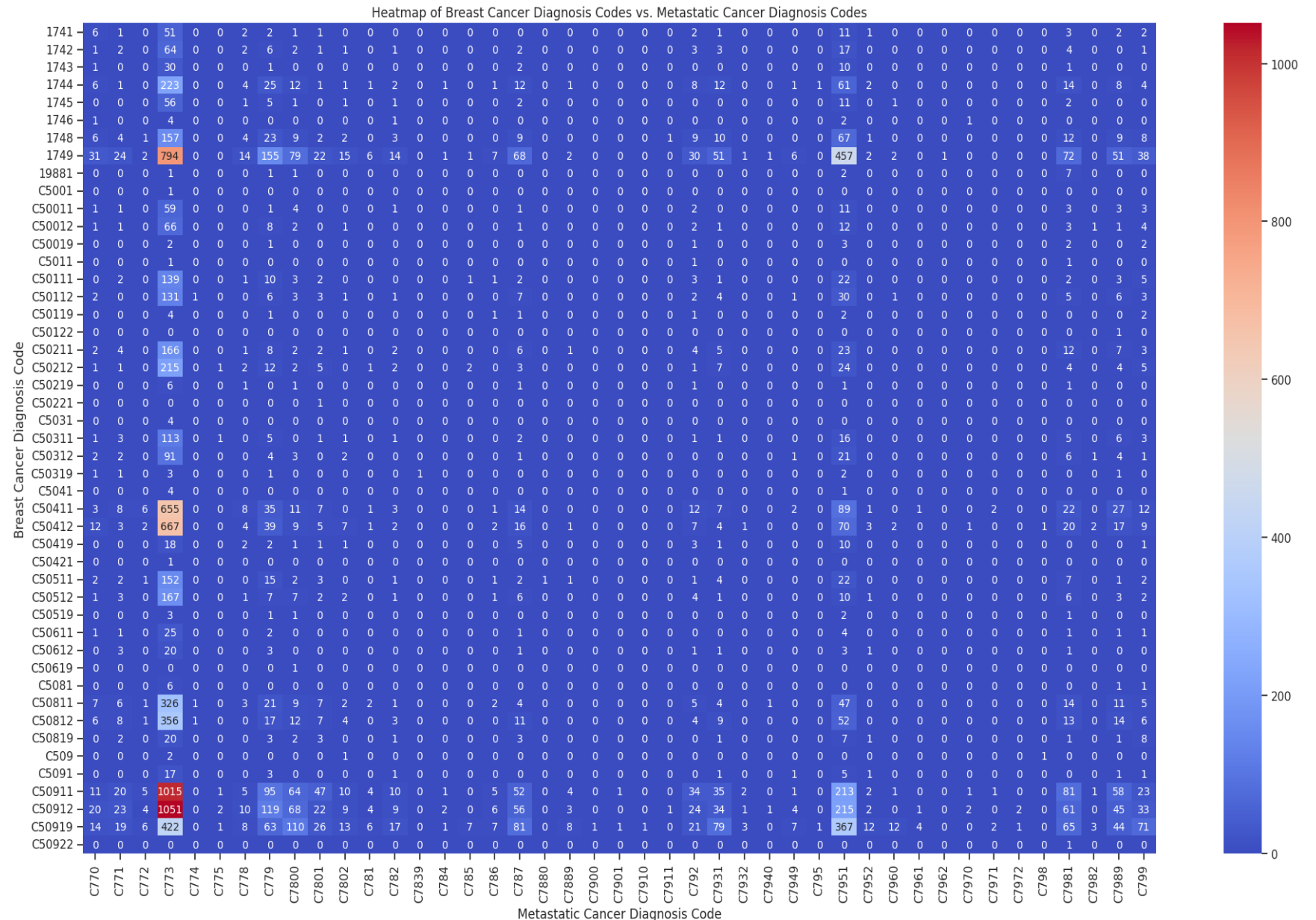Distribution of Metastatic Cancer Diagnosis Codes by Race

**What is the distribution of metastatic cancer diagnosis codes among different races in the dataset?**

# Cont.



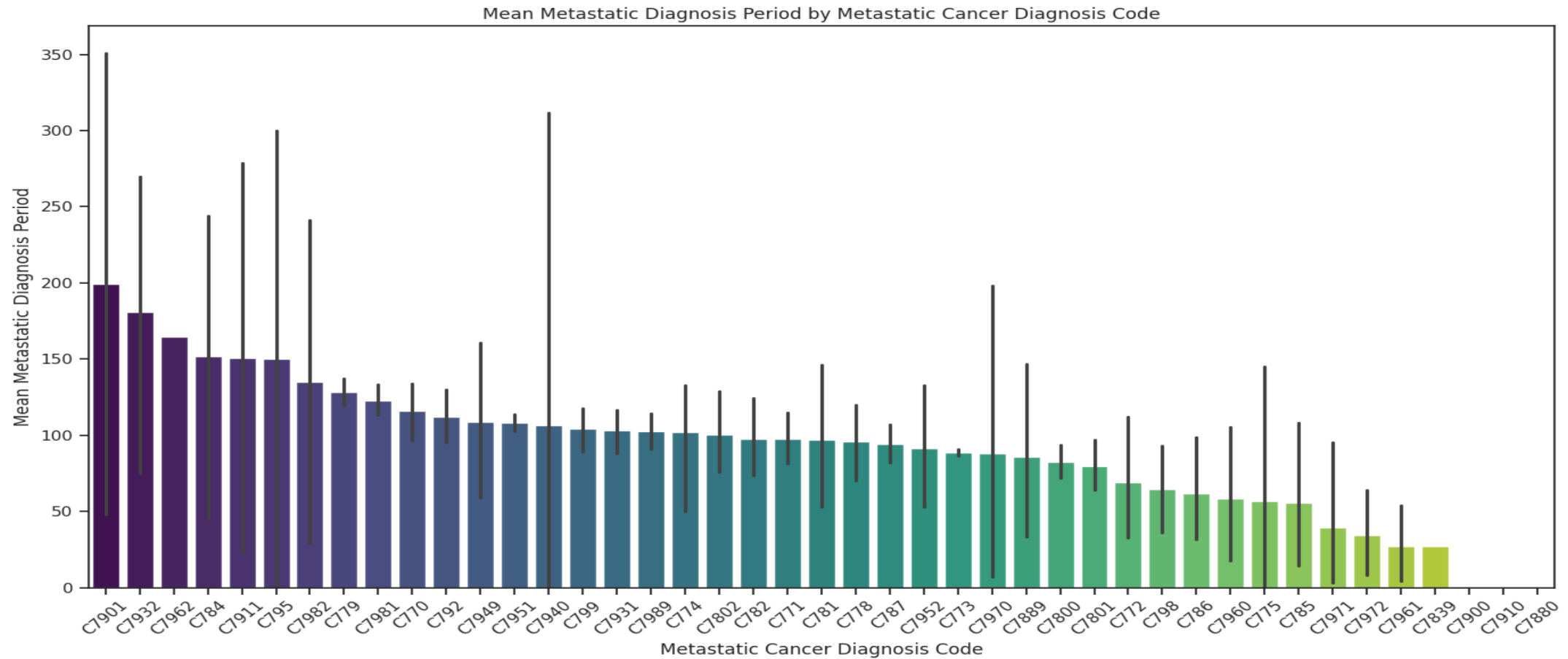Heatmap of Breast Cancer Diagnosis Codes vs. Metastatic Cancer Diagnosis Codes

**What is the co-occurrence pattern between breast cancer diagnosis codes and metastatic cancer diagnosis codes based on the dataset?**

# Cont.



Mean Metastatic Diagnosis Period by Metastatic Cancer Diagnosis Code

**What is the mean metastatic diagnosis period for different metastatic cancer diagnosis codes in the dataset?**

# Word Cloud of Caner Code Description

# THANK YOU