# Healthcare Cost and Utilization Project Code Preparation

**Drexel University LeBow College of Business**
**Fall Term 2023**

## Type 1 Diabetes Table of Contents

## Asthma Table of Contents

# Merging Data

Dataset Importation:
- Import 'dask.dataframe' library for parallel computing with larger datasets.
- Using Dask, three datasets — **Hospital**, **Core**, and **Severity** — are read into memory.

Merging Hospital with Severity Dataset:
- The Hospital and Severity datasets are merged on the 'HOSP_KID' column. It is a 'left join', which means all entries from the Hospital dataset and corresponding entries from the Severity dataset will be included.
- The resultant merged dataset is saved as a **Hospital_Severity.csv** file in the specified output directory (*Dataset/Data Cleaning/Merged*).

Merging Hospital with Core Dataset:

- Similarly, the Hospital and Core datasets are merged on the 'HOSP_KID' column using a 'left join'.
- This merged dataset is also saved as a **Hospital_Core.csv** file in the designated output location (Dataset/Data Cleaning/Merged).

# Diabetes K-Modes Clustering

Dataset Importation:
- Import libraries such as Pandas, Dask, and Numpy and read the hospital datasets using Dask for efficient memory.
- Apply age filters to isolate patients between 0 to 18 years and filter diagnosis codes for type 1 diabetes.

Data Transformation:
- Categorize patients into age groups (*'Infants'*, *'Children'*, and *'Adolescents'*). Extract relevant diagnosis codes and convert Dask dataframes to Pandas dataframes for further processing.

Elbow Method:
- Use the Elbow method to determine the optimal number of clusters by plotting the cost function against the number of clusters.

KModes Algorithm:
- Apply the KModes clustering algorithm using the optimal number of clusters. Ensure that the algorithm is initialized with a consistent random seed for reproducibility.

Export Predicted Cluster Dataset:

- Export the predicted clustering dataframe with assigned clusters to a **HCUP_Clustering.csv** file for further cluster data analysis.

# Diabetes Chi-Square Hypothesis Testing

Dataset Importation:
- Import libraries such as Pandas, and read the cluster datasets **HCUP_Clustering.csv** to dataframe. It includes relevant columns for the analysis, including 'GENDER', 'ALL_PATIENT_SEVERITY','TYPE1_DIABETES_CODES','HOSPITAL_REGION', 'RACE_DESC', 'LOS', and 'AGE_GROUP'.
- To conducting hypothesis tests, first to import scipy.stats library from stats.

Gender and Severity of Diagnosis:
- Create a contingency table for 'GENDER' and 'ALL_PATIENT_SEVERITY' whether the test is between two different categorical variables are associate or not.
- Perform a chi-square test and interpret the results.

Gender and Type 1 Diabetes Codes:
- Create a contingency table for 'GENDER' and 'TYPE1_DIABETES_CODES'.
- Perform a chi-square test and interpret the results.

Hospital Region vs. Ethnicity:
- Create a contingency table for 'HOSPITAL_REGION' and 'RACE_DESC'.
- Perform a chi-square test and interpret the results.

Ethnicity vs. Type 1 Diabetes Codes:
- Create a contingency table for 'RACE_DESC' and 'TYPE1_DIABETES_CODES'.
- Perform a chi-square test and interpret the results.

Length of Stay vs. Age Group:
- Create a contingency table for 'LOS' and 'AGE_GROUP'.
- Perform a chi-square test and interpret the results.

Length of Stay vs. Severity:
- Create a contingency table for 'LOS' and 'ALL_PATIENT_SEVERITY'.
- Perform a chi-square test and interpret the results.

Age Group vs. Severity:
- Create a contingency table for 'AGE_GROUP' and 'ALL_PATIENT_SEVERITY'.
- Perform a chi-square test and interpret the results.

Interpreting the Results:
- These hypothesis tests will output the Chi-square statistic, p-value, degrees of freedom, and a conclusion to accept or reject the null hypothesis based on 5% significance level

# Diabetes Codes: Jupyter Notebook

- Merging data source

```
In [1]:   1  import pandas as pd
          2  import dask.dataframe as dd
          3  import os
          4
          5  # File paths
          6  hospital_file = "C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (Codes and EDA)
          7  core_file = "C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (Codes and EDA)/Dat
          8  severity_file = "C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (Codes and EDA)
          9  merged_severity_output_file = "C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (
         10  merged_core_output_file = "C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (Code
         11
```

**Define dictionary**

```
In [2]:   1  # Define the dtype dictionary
          2  dtype_dict = {
          3       'AGE_NEONATE': 'str', 'AMONTH': 'str', 'ELECTIVE': 'str', 'LOS': 'str', 'PAY1': 'str', 'PRDAY10': 'str',
          4   'PRDAY11': 'str','PRDAY8': 'str', 'PRDAY9': 'str','RACE': 'str', 'TOTCHG': 'str', 'ZIPINC_QRTL': 'str',
          5   'DQTR': 'str', 'PRDAY12':'str', 'PRDAY13':'str', 'PRDAY14':'str', 'PRDAY15':'str', 'PRDAY16':'str', 'PRDAY17':'str',
          6   'PRDAY18':'str', 'PRDAY19':'str', 'PRDAY20':'str', 'PRDAY21':'str', 'PRDAY22':'str', 'PRDAY23':'str',
          7   'PRDAY24':'str', 'PRDAY25':'str', 'PRDAY2':'str', 'PRDAY3':'str', 'PRDAY4':'str', 'PRDAY5':'str',
          8   'PRDAY6':'str', 'TRAN_IN':'str', 'TRAN_OUT':'str', 'FEMALE':'str',  'DISPUNIFORM': 'str', 'DIED':'str'
          9  }
```

**Read data**

```
In [3]:   1  # Read datasets
          2  hospital_df = dd.read_csv(hospital_file)
          3  core_df = dd.read_csv(core_file, dtype=dtype_dict)
          4  severity_df = dd.read_csv(severity_file)
```

**Merging Hospital and Core**

```
In [4]:   1  merged_core_df = hospital_df.merge(core_df, on='HOSP_KID', how='left')
          2
          3  # Check if the file exists
          4  if not os.path.exists(merged_core_output_file):
          5      # If the file does not exist, export the final merged dataframe to a CSV file
          6      merged_core_df.compute().to_csv(merged_core_output_file, index=False)
          7
```

**Merging Hospital and Severity**

```
In [5]:   1  # Merge the dataframes Hospital and Severity
          2  merged_severity_df = hospital_df.merge(severity_df, on='HOSP_KID', how='left')
          3
          4  # Check if the file exists
          5  if not os.path.exists(merged_severity_output_file):
          6      # If the file does not exist, export the final merged dataframe to a CSV file
          7      merged_severity_df.compute().to_csv(merged_severity_output_file, index=False)
          8
```

- Clustering with KModes

```
In [ ]:   1  import pandas as pd
          2  import dask.dataframe as dd
          3  import numpy as np
          4  import matplotlib.pyplot as plt
          5  import seaborn as sns
```

## Load csv file

```
In [ ]:   1  hospital_core_file = "C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (Codes and
          2  hospital_severity_file = "C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (Codes
```

## Define dictionary

```
In [ ]:   1  # Define the dtype dictionary
          2  dtype_dict = {
          3       'AGE_NEONATE': 'str', 'AMONTH': 'str', 'ELECTIVE': 'str', 'LOS': 'str', 'PAY1': 'str', 'PRDAY10': 'str',
          4   'PRDAY11': 'str','PRDAY8': 'str', 'PRDAY9': 'str','RACE': 'str', 'TOTCHG': 'str', 'ZIPINC_QRTL': 'str',
          5   'DQTR': 'str', 'PRDAY12':'str', 'PRDAY13':'str', 'PRDAY14':'str', 'PRDAY15':'str', 'PRDAY16':'str', 'PRDAY17':'str',
          6   'PRDAY18':'str', 'PRDAY19':'str', 'PRDAY20':'str', 'PRDAY21':'str', 'PRDAY22':'str', 'PRDAY23':'str',
          7   'PRDAY24':'str', 'PRDAY25':'str', 'PRDAY2':'str', 'PRDAY3':'str', 'PRDAY4':'str', 'PRDAY5':'str',
          8   'PRDAY6':'str', 'TRAN_IN':'str', 'TRAN_OUT':'str', 'FEMALE':'object',  'DISPUNIFORM': 'str', 'DIED':'str'
          9  }
         10
```

## Read dataset and filter only type 1 diabetes codes tranpose into one columns

```python
# Read hospital_core dataset
hospital_core_df = dd.read_csv(hospital_core_file, dtype=dtype_dict)
# Read severity dataset
hospital_severity_df = dd.read_csv(hospital_severity_file, dtype=dtype_dict)

# Filter for age
age_condition = (hospital_core_df['AGE'] >= 0) & (hospital_core_df['AGE'] <= 18)

# Pre-existing type 1 diabetes diagnosis codes
type1_preexist = ['O24011', 'O24012', 'O23013', 'O24019', 'O2402', 'O2403']

# Vectorized condition for both type1_preexist and 'E10' for columns I10_DX1 to I10_DX5
dx_conditions = False
for column in ['I10_DX1', 'I10_DX2', 'I10_DX3', 'I10_DX4', 'I10_DX5']:
    dx_conditions |= hospital_core_df[column].isin(type1_preexist).fillna(False)
    dx_conditions |= hospital_core_df[column].str.startswith('E10').fillna(False)

# Combine age and diabetes diagnosis conditions
combined_condition = age_condition & dx_conditions

hospital_core_df = hospital_core_df[combined_condition]

# Extract all diagnosis codes related to type 1 diabetes
def extract_diagnosis_codes(row):
    return [code for code in row if pd.notnull(code) and (code.startswith('E10') or code in type1_preexist)]

hospital_core_df['TYPE1_DIABETES_CODES'] = hospital_core_df[['I10_DX1', 'I10_DX2', 'I10_DX3', 'I10_DX4', 'I10_DX5']].apply(e

# Transform dask dataframe to pandas dataframe
hospital_core_df = hospital_core_df.compute()

# Reset the index for a unique index
hospital_core_df.reset_index(drop=True, inplace=True)

#hospital_core_df.head()
```

```python
hospital_core_df.shape[0]
```

```python
# Join Severity to Hospital_Core by 'RECNUM'
hospital_severity_df = hospital_severity_df.compute()

hospital_core_severity_df = hospital_core_df.merge(hospital_severity_df, on='RECNUM', how='inner')
hospital_core_severity_df.head()
```

```python
hospital_core_severity_df.info()
```

```python
# Join Severity to Hospital_Core by 'RECNUM'
hospital_severity_df = hospital_severity_df.compute()

hospital_core_severity_df = hospital_core_df.merge(hospital_severity_df, on='RECNUM', how='inner')
hospital_core_severity_df.head()
```

```python
hospital_core_severity_df.info()
```

```python
# Checking Null values
hospital_core_severity_df.isnull().sum()*100/hospital_core_severity_df.shape[0]
```

```python
def categorize_age(age):
    if age < 1:
        return "Infant"
    elif 1 <= age <= 13:
        return "Children"
    else:
        return "Adolescents"

# Categorization function 'AGE_GROUP'
hospital_core_severity_df['AGE_GROUP'] = hospital_core_severity_df['AGE'].apply(categorize_age)
hospital_core_severity_df['TYPE1_DIABETES_CODES'] = hospital_core_severity_df['TYPE1_DIABETES_CODES'].apply(lambda x: x[0] i


hospital_core_severity_df  = hospital_core_severity_df.drop('AGE',axis = 1)
```

```python
select_diagnosis_columns  = ['AGE_GROUP', 'GENDER', 'HOSPITAL_REGION', 'LOS','ALL_PATIENT_SEVERITY', 'RACE_DESC','TYPE1_DIAB

hospital_core_severity_df = hospital_core_severity_df[select_diagnosis_columns]
hospital_core_severity_df.head()
```

```python
print("Total number of records for clustering: " + str(hospital_core_severity_df.shape))
```

**Copy Data**

```
In [ ]:   1  # Keep a copy of data
          2  hospital_core_severity_df_copy = hospital_core_severity_df.copy()
```

```
In [ ]:   1  hospital_core_severity_df.head()
```

## Find Optimal k by Elbow Method

```
In [ ]:   1  from kmodes.kmodes import KModes
          2  import matplotlib.pyplot as plt
          3
          4  np.random.seed(42)
          5  cost = []
          6  for num_clusters in list(range(1,8)):
          7      kmode = KModes(n_clusters=num_clusters, init ="Cao", n_init = 5, verbose=0)
          8      kmode.fit_predict(hospital_core_severity_df)
          9      cost.append(kmode.cost_)
         10
         11  y = np.array([i for i in range(1,8,1)])
         12
         13  plt.figure(figsize=(10,6))
         14  #plt.plot(y,cost)
         15  plt.plot(y, cost, marker='o')
         16  plt.xlabel('Number of Clusters')
         17  plt.ylabel('Cost')
         18  plt.title('Elbow Method For Optimal k')
         19  plt.grid(False)
         20  plt.show()
         21
         22
```

## Model fitting

```
In [ ]:   1  optimal_k =3 # as per elbow method
```

```
In [ ]:   1  # fit the KModes clustering algorithm with 5 optimal clusters.
          2  np.random.seed(42)
          3  kmd = KModes(n_clusters=optimal_k, init = "Huang", n_init =5 , verbose=0)
          4  kmd_clusters = kmd.fit_predict(hospital_core_severity_df)
```

```
In [ ]:   1  # Combining predicted clusters with the original dataframe
          2  hospital_core_severity_df = hospital_core_severity_df_copy.reset_index()
```

```
In [ ]:   1  kmd_clusters_df = pd.DataFrame(kmd_clusters)
          2  kmd_clusters_df.columns = ['CLUSTERS']
          3  kmd_clusters_df = pd.concat([hospital_core_severity_df, kmd_clusters_df], axis = 1).reset_index()
          4  kmd_clusters_df = kmd_clusters_df.drop(['index', 'level_0'], axis = 1)
          5  # added row numbers
          6  kmd_clusters_df['ROWNUM'] = range(1, len(kmd_clusters_df) + 1)
          7
          8  # Create a mapping from old to new cluster labels
          9  label_mapping = {0: 1, 1: 2, 2: 3}
         10
         11  # start with cluster 1 instead of 0. map the old cluster labels to the new ones
         12  kmd_clusters_df['CLUSTERS'] = kmd_clusters_df['CLUSTERS'].map(label_mapping)
```

**Extract the clusters dataset**

```
In [ ]:   1  import os
          2  output_file = 'C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Production (Codes and EDA)/D
          3
          4  # Check if the file exists
          5  if not os.path.exists(output_file):
          6      # If the file does not exist, create HCUP_Clustering.csv file
          7      kmd_clusters_df.to_csv(output_file, index=False)
          8
          9  kmd_clusters_df.head()
```

## Descriptive Analysis

### Cluster Proportion

```python
# Calculate the proportions
np.random.seed(42)
cluster_size = kmd_clusters_df['CLUSTERS'].value_counts().sort_index()
cluster_size_proportion = (cluster_size / cluster_size.sum()) * 100

# Plotting
plt.figure(figsize=(10, 6))
cluster_size_proportion.plot(kind='bar')

# Adding titles and labels
plt.title('Cluster Size Proportions (%)')
plt.xlabel('Cluster')
plt.ylabel('Proportion (%)')
plt.xticks(rotation=0)

# show percentages above the bars
for i, value in enumerate(cluster_size_proportion):
    plt.text(i, value + 0.5, f"{value:.2f}%", ha='center', va='bottom')

plt.tight_layout()
plt.show()
```

```python
kmd_clusters_df.head()
```

```python
# Heatmap of AGE_GROUP
crosstab = pd.crosstab(kmd_clusters_df['AGE_GROUP'], kmd_clusters_df['CLUSTERS'], normalize='index') * 100

# matrix with percentages
annot_array = crosstab.values
annotations = [[f'{val:.2f}%' for val in row] for row in annot_array]

plt.figure(figsize=(10, 6))
sns.heatmap(crosstab, annot=annotations, cmap="YlGnBu", fmt="")
plt.title("Distribution of Age Group across Clusters")
plt.show()

# Heatmap of GENDER
crosstab = pd.crosstab(kmd_clusters_df['GENDER'], kmd_clusters_df['CLUSTERS'], normalize='index') * 100

# matrix with percentages
annot_array = crosstab.values
annotations = [[f'{val:.2f}%' for val in row] for row in annot_array]

plt.figure(figsize=(10, 6))
sns.heatmap(crosstab, annot=annotations, cmap="YlGnBu", fmt="")
plt.title("Distribution of Gender across Clusters")
plt.show()

# Heatmap of RACE_DESC
crosstab = pd.crosstab(kmd_clusters_df['RACE_DESC'], kmd_clusters_df['CLUSTERS'], normalize='index') * 100

# matrix with percentages
annot_array = crosstab.values
annotations = [[f'{val:.2f}%' for val in row] for row in annot_array]

plt.figure(figsize=(10, 6))
sns.heatmap(crosstab, annot=annotations, cmap="YlGnBu", fmt="")
plt.title("Distribution of Race across Clusters")
plt.show()

# Heatmap of RACE_DESC
crosstab = pd.crosstab(kmd_clusters_df['HOSPITAL_REGION'], kmd_clusters_df['CLUSTERS'], normalize='index') * 100

# matrix with percentages
annot_array = crosstab.values
annotations = [[f'{val:.2f}%' for val in row] for row in annot_array]

plt.figure(figsize=(10, 6))
sns.heatmap(crosstab, annot=annotations, cmap="YlGnBu", fmt="")
plt.title("Distribution of Hospital Region across Clusters")
plt.show()

# Heatmap of Severity
crosstab = pd.crosstab(kmd_clusters_df['ALL_PATIENT_SEVERITY'], kmd_clusters_df['CLUSTERS'], normalize='index') * 100

# matrix with percentages
annot_array = crosstab.values
annotations = [[f'{val:.2f}%' for val in row] for row in annot_array]

plt.figure(figsize=(10, 6))
sns.heatmap(crosstab, annot=annotations, cmap="YlGnBu", fmt="")
plt.title("Distribution of Health Severity across Clusters")
plt.show()

# Heatmap of LOS
kmd_clusters_df['LOS'] = kmd_clusters_df['LOS'].astype(int)
```

```python
64  def categorize_los(los):
65      if los <= 5:
66          return "Short-term"
67      elif 5 < los <= 15:
68          return "Mid-term"
69      else:
70          return "Long-term"
71
72  kmd_clusters_df['LOS_Category'] = kmd_clusters_df['LOS'].apply(categorize_los)
73
74  # crosstab with the LOS_Category and Clusters
75  crosstab = pd.crosstab(kmd_clusters_df['LOS_Category'], kmd_clusters_df['CLUSTERS'])
76  # Convert the crosstab values to proportions
77  crosstab_percentage = (crosstab / crosstab.sum().sum()) * 100
78
79  # Generate annotations with percentage symbol
80  annotations = crosstab_percentage.applymap(lambda x: f"{x:.2f}%")
81
82  plt.figure(figsize=(10, 6))
83  sns.heatmap(crosstab_percentage, annot=annotations, cmap="YlGnBu", fmt="s")
84  plt.title("Distribution of LOS across Clusters in Percentages")
85  plt.xlabel("Cluster")
86  plt.ylabel("LOS Category")
87  plt.show()
88
89
90
91  # Heatmap of TYPE1 DIABETES
92  crosstab = pd.crosstab(kmd_clusters_df['TYPE1_DIABETES_CODES'], kmd_clusters_df['CLUSTERS'], normalize='index') * 100
93
94  # matrix with percentages
95  annot_array = crosstab.values
96  annotations = [[f'{val:.2f}%' for val in row] for row in annot_array]
97
98  plt.figure(figsize=(10, 6))
99  sns.heatmap(crosstab, annot=annotations, cmap="YlGnBu", fmt="")
100 plt.title("Distribution of Type1 Diabetes Diagnosis across Clusters")
101 plt.show()
102
103
104
```

- Chi-Square Testing

```python
In [ ]:  1  import pandas as pd
```

```python
In [ ]:  1  hospital_cluster_data =  pd.read_csv("C:/Users/Lenovo/OneDrive - Drexel University/Fall 2023/BSAN-710/Capstone Project/Produ
         2  hospital_cluster_data.head()
```

## Hypothesis 1: A disparity in gender concerning the severity of the diagnosis of type 1 diabetes.

create the contingency tables for each relationship, performed the chi-square test to determine the significance of associations, and provided the interpretations based on the p-values

```python
In [ ]:  1  import scipy.stats as stats
         2
         3  # Create a contingency table for GENDER vs ALL_PATIENT_SEVERITY
         4  contingency_table_gender_sev = pd.crosstab(hospital_cluster_data['GENDER'], hospital_cluster_data['ALL_PATIENT_SEVERITY'])
         5  print("Contingency Table for Gender and Severity:")
         6  print(contingency_table_gender_sev)
         7
         8  # Conduct the chi-square test
         9  chi2, p, df_1, _ = stats.chi2_contingency(contingency_table_gender_sev)
        10  print(f"\nHypothesis 1 chi-square: {chi2}, p-value: {format(p,'.10f')},  degrees of freedom: {df_1}\n")
        11  print("-" * 50)
        12
        13  # Interpret the result gender vs severity
        14  alpha = 0.05
        15  if p < alpha:
        16      print("Reject the null hypothesis: There is a significant relationship between gender and severity.")
        17  else:
        18      print("Fail to reject the null hypothesis: There is no significant relationship between gender and severity.")
        19
        20  #-------------------------------------------------------------------------------------------------
        21
        22  # Create a contingency table for GENDER vs TYPE1_DIABETES_CODES
        23  contingency_table_gender_dia = pd.crosstab(hospital_cluster_data['GENDER'], hospital_cluster_data['TYPE1_DIABETES_CODES'])
        24  print("\nContingency Table for Gender and Type 1 Diabetes Codes:")
        25  print(contingency_table_gender_dia)
        26
        27  # Conduct the chi-square test
        28  chi2_gender, p_gender, df_2, _ = stats.chi2_contingency(contingency_table_gender_dia)
        29  print(f"\nChi-square for Gender vs. Type 1 Diabetes Codes: {chi2_gender}, p-value: {format(p_gender,'.10f')}, degrees of fre
        30  print("-" * 50)
        31
        32  # Interpret the result for gender vs type 1 diabetes
        33  alpha = 0.05
        34  if p_gender < alpha:
        35      print("Reject the null hypothesis: There is a significant relationship between gender and Type 1 diabetes codes.")
        36  else:
        37      print("Fail to reject the null hypothesis: There is no significant relationship between gender and Type 1 diabetes codes
        38  print("=" * 50)
        39
        40  #-------------------------------------------------------------------------------------------------
        41
```

```
In [ ]:  1  print("\nConclusion:")
         2
         3  print("These findings support Hypothesis 1, gender and severity have correlations with Type 1 diabetes.")
```

## Hypothesis 2: Factors such as the regions of hospitals and ethnicity are found to have a correlation with Type 1 diabetes within the group.

```
In [ ]:   1  # Create a contingency table for HOSPITAL_REGION vs Ethnicity
          2  contingency_table_hosp_race = pd.crosstab(hospital_cluster_data['HOSPITAL_REGION'], hospital_cluster_data['RACE_DESC'])
          3  print("\nContingency Table for Hospital Region vs. Ethnicity:")
          4  print(contingency_table_hosp_race)
          5
          6  # Conduct the chi-square test
          7  chi2, p, df_3, _ = stats.chi2_contingency(contingency_table_hosp_race)
          8  print(f"\nChi-square for Hospital Region vs. Ethnicity: {chi2}, p-value: {format(p,'.10f')}, degrees of freedom: {df_3}\n")
          9  print("-" * 50)
         10
         11  # Interpret the result for hospital region
         12  alpha = 0.05
         13  if p < alpha:
         14      print("Reject the null hypothesis: There is a significant relationship between hospital region and ethnicity.")
         15  else:
         16      print("Fail to reject the null hypothesis: There is no significant relationship between hospital region and Ethnicity.")
         17
         18  print("=" * 50)
         19
         20  #----------------------------------------------------------------------------------------------------
         21  # Create a contingency table for RACE_DESC vs TYPE1_DIABETES_CODES
         22  contingency_table_race_dia = pd.crosstab(hospital_cluster_data['RACE_DESC'], hospital_cluster_data['TYPE1_DIABETES_CODES'])
         23  print("\nContingency Table for Ethnicity vs. Type 1 Diabetes Codes:")
         24  print(contingency_table_race_dia)
         25
         26  # Conduct the chi-square test
         27  chi2_race, p_race, df_4, _ = stats.chi2_contingency(contingency_table_race_dia)
         28  print(f"\nChi-square for Ethnicity vs. Type 1 Diabetes Codes: {chi2_race}, p-value: {format(p_race,'.10f')}, degrees of free
         29  print("-" * 50)
         30
         31  # Interpret the result for ethnicity
         32  if p_race < alpha:
         33      print("Reject the null hypothesis: There is a significant relationship between ethnicity and Type 1 diabetes codes.")
         34  else:
         35      print("Fail to reject the null hypothesis: There is no significant relationship between ethnicity and Type 1 diabetes co
         36
```

```
In [ ]:  1  print("\nConclusion:")
         2
         3  print("These findings support Hypothesis 2, hospital regions and ethnicity have correlations with Type 1 diabetes within the
```

## Hypothesis 3: Distribution of the length of a patient's stay in the hospital and severity of diabetes diagnoses among different age groups.

```python
In [ ]:    1  # Define the bins and their labels
           2  bins = [-1, 0, 7, 14, 30, float('inf')] # the -1 ensures that 0 is included in the first bin
           3  labels = ['0 days', '1-7 days', '8-14 days', '15-30 days', '31+ days']
           4
           5
           6  # Create a new column for the binned LOS
           7  hospital_cluster_data['LOS_Binned'] = pd.cut(hospital_cluster_data['LOS'], bins=bins, labels=labels)
           8
           9  # Create a contingency table for the binned LOS vs AGE_GROUP
          10  contingency_table_los_age = pd.crosstab(hospital_cluster_data['LOS_Binned'], hospital_cluster_data['AGE_GROUP'])
          11  print("\nContingency Table for Length of Stay vs. Age Group:")
          12  print(contingency_table_los_age)
          13
          14  # Conduct the chi-square test
          15  chi2, p, df_5, _ = stats.chi2_contingency(contingency_table_los_age)
          16  print(f"\nChi-square for Length of stay vs. Age Group: {chi2}, p-value: {format(p,'.10f')}, degrees of freedom: {df_5}\n")
          17  print("-" * 50)
          18
          19  # Interpret the result for LOS
          20  alpha = 0.05
          21  if p < alpha:
          22      print("Reject the null hypothesis: There is a significant relationship between length of stay and patient age group.")
          23  else:
          24      print("Fail to reject the null hypothesis: There is no significant relationship between length of stay and patient age g
          25
          26  print("=" * 50)
          27
          28  #-------------------------------------------------------------------------------------------------
          29
          30  # Create a contingency table for Los vs severity
          31  contingency_table_los_sev = pd.crosstab(hospital_cluster_data['LOS_Binned'], hospital_cluster_data['ALL_PATIENT_SEVERITY'])
          32  print("\nContingency Table for Lenght of stay vs. Severity:")
          33  print(contingency_table_los_sev)
          34
          35  # Conduct the chi-square test
          36  chi2_los, p_los, df_6, _ = stats.chi2_contingency(contingency_table_los_sev)
          37  print(f"\nChi-square for Lenght of stay vs. Severity: {chi2_los}, p-value: {format(p_los,'.10f')}, degrees of freedom: {df_6
          38  print("-" * 50)
          39
          40  # Interpret the result for Los and severity
          41  if p_los < alpha:
          42      print("Reject the null hypothesis: There is a significant relationship between lenght of stay and severity.")
          43  else:
          44      print("Fail to reject the null hypothesis: There is no significant relationship between lenght of stay and severity.")
          45
          46  #-------------------------------------------------------------------------------------------------
          47
          48  # Create a contingency table for Age group vs severity
          49  contingency_table_age_sev = pd.crosstab(hospital_cluster_data['AGE_GROUP'], hospital_cluster_data['ALL_PATIENT_SEVERITY'])
          50  print("\nContingency Table for Age Group vs. Severity:")
          51  print(contingency_table_age_sev)
          52
          53  # Conduct the chi-square test
          54  chi2_ag, p_ag, df_7, _ = stats.chi2_contingency(contingency_table_age_sev)
          55  print(f"\nChi-square for Lenght of stay vs. Severity: {chi2_ag}, p-value: {format(p_ag,'.10f')}, degrees of freedom: {df_7}\
          56  print("-" * 50)
          57
          58  # Interpret the result for age group and severity
          59  if p_ag < alpha:
          60      print("Reject the null hypothesis: There is a significant relationship between age group and severity.")
          61  else:
          62      print("Fail to reject the null hypothesis: There is no significant relationship between age group and severity.")
          63
```

## Asthma Initial Setup

Similar to the merging done for Diabetes, the datasets were combined together. Afterward, a basic summary of the dataset was created to determine the spread of the data. In addition, data cleaning and preprocessing were done prior to statistical testing and clustering analysis.

## Asthma Statistical Testing Using Chi-Squared Test

Similar to the statistical testing done for Diabetes, chi-squared test for statistical analysis, examining the relationships between categorical variables. The test generated a p-value and a Chi-Square value. A low p-value ($<0.05$) suggested significant associations, enabling us to make informed decisions about the validity of patterns within our dataset and test our initial hypothesis.

## Asthma K-Means Clustering

Applying K-Means clustering, we segmented the data into distinct groups based on similarities. This unsupervised machine learning technique allowed us to identify inherent patterns and groupings within the dataset.

For better results we created bins for variables like Age and Total Charges. The clustering algorithm divided the entire data for asthmatic patients into 5 clusters.

## Asthma K-Modes Clustering

We also employed K-Modes clustering, specifically designed for categorical data. This approach further enriched our understanding by revealing patterns in non-numeric variables, providing a holistic view of the dataset.

The clustering algorithm divided the entire data for asthmatic patients into 3 clusters.

## Asthma Codes: SQL and Jupyter Notebook

## SQL Queries:

```sql
--Query to create separate table for respiratory diseases based on diagnosis code:
CREATE TABLE capstone2.RESPIRATORY_DISEASES AS

SELECT HOSP_KID
        ,RECNUM
        ,DXCCSR_RSP001
        ,DXCCSR_RSP002
        ,DXCCSR_RSP003
```

```
        ,DXCCSR_RSP004
        ,DXCCSR_RSP005
        ,DXCCSR_RSP006
        ,DXCCSR_RSP007
        ,DXCCSR_RSP008
        ,DXCCSR_RSP009
        ,DXCCSR_RSP010
        ,DXCCSR_RSP011
        ,DXCCSR_RSP012
        ,DXCCSR_RSP013
        ,DXCCSR_RSP014
        ,DXCCSR_RSP015
        ,DXCCSR_RSP016
        ,DXCCSR_RSP017
FROM `capstone - 400517. capstone2.kid_GPRS`


--QUERY FOR asthma cases by Gender
SELECT FEMALE
        ,count(*)
FROM `capstone - 400517. capstone2.table_core`
WHERE Asthma <> '0'
        AND Asthma IS NOT NULL
GROUP BY 1


--QUERY FOR asthma cases by INCOME
SELECT ZIPINC_QRTL
        ,count(*)
FROM `capstone - 400517. capstone2.table_core`
WHERE Asthma <> '0'
GROUP BY 1




--QUERY FOR asthma cases by race
SELECT RACE
        ,count(*)
FROM `capstone - 400517. capstone2.table_core`
--WHERE Asthma <> '0'
GROUP BY 1

--query for asthma and infectious diseases
SELECT Asthma
        ,INF_All
        ,COUNT(*) AS INFANDASTHMA
FROM `capstone - 400517. capstone2.table_core`
WHERE Asthma <> '0'
        AND INF_All <> 0
GROUP BY 1
        ,2

--COMORBIDITY ANALYSIS FOR ASTHMA AND OTHER DISEASES
SELECT RECNUM
FROM `capstone - 400517. capstone2.table_core`
WHERE (
                Blood_All <> 0
                OR MAL_All <> 0
                OR MBD_ALL <> 0
                OR END_ALL <> 0
                OR INF_All <> 0
```

```sql
                    OR DIG_ALL <> 0
                    OR NVS_ALL <> 0
                    OR PRG_All <> 0
                    OR CIR_All <> 0
                    OR SKN_All <> 0
                    OR MUS_All <> 0
                    OR Tumor_All <> 0
                    OR EYE_All <> 0
                    OR EAR_All <> 0
                    )
        AND Asthma <> '0'
```

```sql
--Data cleaning in gender
DELETE
FROM `capstone - 400517. capstone2.table_core`
WHERE FEMALE IN (
                    'A'
                    ,'C'
                    ,'nan'
                    );
```

```sql
--query to get total charges by payer
SELECT
        --RECNUM
        Pay1
        ,sum(TotalChargesInteger)
FROM `capstone - 400517. capstone2.table_core`
WHERE Asthma <> '0'
        AND (
                    Blood_All = 0
                    AND MAL_All = 0
                    AND MBD_ALL = 0
                    AND END_ALL = 0
                    AND INF_All = 0
                    AND DIG_ALL = 0
                    AND NVS_ALL = 0
                    AND PRG_All = 0
                    AND CIR_All = 0
                    AND SKN_All = 0
                    AND MUS_All = 0
                    AND Tumor_All = 0
                    AND EYE_All = 0
                    AND EAR_All = 0
                    )
        AND PAY1 IN (
                    '1'
                    ,'2'
                    )
--and race <> 'nan'
--and Race not like '%nan'
GROUP BY 1
```

```sql
--Query to check count of each endocrine disease
SELECT SUM(CASE
                    WHEN CAST(DXCCSR_END001 AS INT64) = 1
```

```sql
                                THEN 1
                    ELSE 0
                    END) AS Count_1
,SUM(CASE

                    WHEN CAST(DXCCSR_END002 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_2
,SUM(CASE

                    WHEN CAST(DXCCSR_END003 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_3
,SUM(CASE

                    WHEN CAST(DXCCSR_END004 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_4
,SUM(CASE

                    WHEN CAST(DXCCSR_END006 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_6
,SUM(CASE

                    WHEN CAST(DXCCSR_END007 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_7
,SUM(CASE

                    WHEN CAST(DXCCSR_END008 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_8
,SUM(CASE

                    WHEN CAST(DXCCSR_END009 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_9
,SUM(CASE

                    WHEN CAST(DXCCSR_END010 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_10
,SUM(CASE

                    WHEN CAST(DXCCSR_END011 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_11
,SUM(CASE

                    WHEN CAST(DXCCSR_END012 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_12
,SUM(CASE

                    WHEN CAST(DXCCSR_END013 AS INT64) = 1
                                THEN 1
                    ELSE 0
                    END) AS Count_13
,SUM(CASE
```

```sql
                                    WHEN CAST(DXCCSR_END014 AS INT64) = 1
                                            THEN 1
                                    ELSE 0
                                    END) AS Count_14
                ,SUM(CASE

                                    WHEN CAST(DXCCSR_END015 AS INT64) = 1
                                            THEN 1
                                    ELSE 0
                                    END) AS Count_15
                ,SUM(CASE

                                    WHEN CAST(DXCCSR_END016 AS INT64) = 1
                                            THEN 1
                                    ELSE 0
                                    END) AS Count_16
                ,SUM(CASE

                                    WHEN CAST(DXCCSR_END017 AS INT64) = 1
                                            THEN 1
                                    ELSE 0
                                    END) AS Count_17

                ,
FROM `capstone - 400517. capstone2.ENDOCRINE_DISEASES`
--Query to analyze asthma and length of stay
SELECT RECNUM
        ,LOS
FROM `capstone - 400517. capstone2.table_core`
WHERE (
                    Blood_All <> 0
                    OR MAL_All <> 0
                    OR MBD_ALL <> 0
                    OR END_ALL <> 0
                    OR INF_All <> 0
                    OR DIG_ALL <> 0
                    OR NVS_ALL <> 0
                    OR PRG_All <> 0
                    OR CIR_All <> 0
                    OR SKN_All <> 0
                    OR MUS_All <> 0
                    OR Tumor_All <> 0
                    OR EYE_All <> 0
                    OR EAR_All <> 0
                    )
        AND Asthma <> '0'
        AND Asthma IS NOT NULL
--Query to get count of individual infectious diseases
SELECT SUM(CASE
                    WHEN CAST(DXCCSR_INF001 AS INT64) <> 0
                            THEN 1
                    ELSE 0
                    END) AS Count_INF001
                ,SUM(CASE

                    WHEN CAST(DXCCSR_INF002 AS INT64) <> 0
                            THEN 1
                    ELSE 0
                    END) AS Count_INF002
                ,SUM(CASE

                    WHEN CAST(DXCCSR_INF003 AS INT64) <> 0
                            THEN 1
                    ELSE 0
                    END) AS Count_INF003
```

```sql
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF004 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF004
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF005 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF005
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF006 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF006
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF007 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF007
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF008 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF008
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF009 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF009
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF010 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF010
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF011 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF011
        ,SUM(CASE
                        WHEN CAST(DXCCSR_INF012 AS INT64) <> 0
                                THEN 1
                        ELSE 0
                        END) AS Count_INF012
FROM `capstone - 400517. capstone2.Infectious_DISEASES`
--Query to clean total charges col.
--ALTER TABLE `capstone-400517.capstone2.table_core`
--ADD COLUMN TotalChargesInteger INT64;
-- Copy data from the existing column to the new column
UPDATE `capstone - 400517. capstone2.table_core`
SET TotalChargesInteger = CASE
                WHEN TOTCHG = 'nan'
                        THEN NULL -- or set to a default value
                WHEN TOTCHG = 'A'
                        THEN NULL
                WHEN TOTCHG = 'C'
                        THEN NULL
                WHEN TOTCHG = 'TOTCHG'
```

```sql
                    THEN NULL
              ELSE CAST(TOTCHG AS INT64)
              END
WHERE 1 = 1


--Query to get LOS AND TOTAL CHARGE FOR PATIENTS WITH ASTHMA AND OTHER DISEASES
SELECT t1.DXCCSR_RSP009
        ,COUNT(*)
FROM (
        SELECT RECNUM
                ,DXCCSR_Default_DX1
                ,DXCCSR_RSP009
        FROM `capstone - 400517. capstone2.kid_GPRS` A11
        WHERE A11.DXCCSR_Default_DX1 = 'RSP009'
        ) t1
GROUP BY 1


--query to create separate table for CongenitalMalfunction
CREATE TABLE capstone2.CongenitalMalfunction_DISEASES AS

SELECT HOSP_KID
        ,RECNUM
        ,DXCCSR_MAL001
        ,DXCCSR_MAL002
        ,DXCCSR_MAL003
        ,DXCCSR_MAL004
        ,DXCCSR_MAL005
        ,DXCCSR_MAL006
        ,DXCCSR_MAL007
        ,DXCCSR_MAL008
        ,DXCCSR_MAL009
        ,DXCCSR_MAL010
FROM `capstone - 400517. capstone2.kid_GPRS`


--query to make new table for infectious diseases
CREATE TABLE capstone2.Infectious _DISEASES AS

SELECT HOSP_KID
        ,RECNUM
        ,DXCCSR_INF001
        ,DXCCSR_INF002
        ,DXCCSR_INF003
        ,DXCCSR_INF004
        ,DXCCSR_INF005
        ,DXCCSR_INF006
        ,DXCCSR_INF007
        ,DXCCSR_INF008
        ,DXCCSR_INF009
        ,DXCCSR_INF010
        ,DXCCSR_INF011
        ,DXCCSR_INF012
FROM `capstone - 400517. capstone2.kid_GPRS`


--Query to make range for LOS
UPDATE capstone2.table_core
SET LOS_Range = CASE
                WHEN LOS BETWEEN '0'
```

```sql
                                AND '1'
                        THEN '0-1'
                WHEN LOS BETWEEN '2'
                                AND '3'
                        THEN '2-3'
                WHEN LOS BETWEEN '4'
                                AND '9'
                        THEN '4-9'
                WHEN LOS > '10'
                        THEN '10+'
                END
WHERE 1 = 1
```

```sql
--query to make range for total charges
UPDATE capstone2.table_core
SET TotalChargeRange = CASE
                WHEN TotalChargesInteger BETWEEN 0
                                AND 12000
                        THEN 1
                WHEN TotalChargesInteger BETWEEN 12001
                                AND 20000
                        THEN 2
                WHEN TotalChargesInteger BETWEEN 20001
                                AND 40000
                        THEN 3
                WHEN TotalChargesInteger BETWEEN 40001
                                AND 80000
                        THEN 4
                WHEN TotalChargesInteger > 80001
                        THEN 5
                END
WHERE 1 = 1
```

```python
# Python script for calculating Chi-Squared value and p value for asthma by gender
#import library
from scipy.stats import chi2_contingency
import matplotlib.pyplot as plt
import numpy as np

observed_data = [[101427, 1333965], [103168, 1248246]]

# Perform chi-square test
chi2, p, _, _ = chi2_contingency(observed_data)


#print p value
print(f"P-value: {p}")

print(f"Chi-squared value: {chi2}")
```

```python
#Python script to check chi-square and P value for asthma by hospital region
#import library
from scipy.stats import chi2_contingency
import matplotlib.pyplot as plt
import numpy as np

observed_data = [[46993, 463696], [50325, 645687], [81598, 1116237], [25679, 356591]]
```

```
# Perform chi-square test
chi2, p, dof, expected = chi2_contingency(observed_data)
#print p value
print(f"P-value: {p}")
```

```
# Python script for chi-sqared test for income and asthma
#import library
from scipy.stats import chi2_contingency
import matplotlib.pyplot as plt
import numpy as np
observed_data = [[74825, 804022], [48867, 634510], [45054, 617432], [33955, 499721]]
# Perform chi-square test
chi2, p, _, _ = chi2_contingency(observed_data)
#print p value
print(f"P-value: {p}")

print(f"Chi-squared value: {chi2}")
```

```
# Python script for chi-sqared test: asthma with race
#import library
from scipy.stats import chi2_contingency
import matplotlib.pyplot as plt
import numpy as np
observed_data = [[78861, 1227805], [60569, 434845], [38607, 458916], [4366, 94119], [1550,
18515], [10290, 160175]]
# Perform chi-square test
chi2, p, _, _ = chi2_contingency(observed_data)
#print p value
print(f"P-value: {p}")
```

# Code for clustering:

```
# Filtering only for RSP009
results = job.to_dataframe()
results = results[results['DXCCSR_Default_DX1']== 'RSP009']
```

```
# Import numpy to handle NaN values
import numpy as np

# Replace all instances of 'nan' with NaN in the results DataFrame
results.replace('nan', np.nan, inplace=True)
```

```
# Drop NA values from table
results.dropna(inplace=True)
```

```python
results.replace('A', np.nan, inplace=True)
```

```python
# Drop NA values form table
results.dropna(inplace=True)
```

```python
# check if we have NA values among unique values

cols = results.columns
for i in cols:

  print(i, results[i].unique())
```

```python
#Converting all values into numeric so that we can run ML algorithms

import pandas as pd

# Convert all columns to numeric, coercing errors to NaN
results = results.apply(pd.to_numeric, errors='coerce')


# Convert to Int64
results = results.astype('Int64')
```

```python
#Checking if the values converted into numeric

cols = results.columns
for i in cols:

  print(i, results[i].unique())
```

```python
#Converting all values into numeric so that we can run ML algorithms

import pandas as pd

# Convert all columns to numeric, coercing errors to NaN
results = results.apply(pd.to_numeric, errors='coerce')


# Convert to Int64
results = results.astype('Int64')
```

```python
#Checking if the values converted into numeric

cols = results.columns
for i in cols:
```

```python
  print(i, results[i].unique())
```

```python
# Separating the features that  we need for K-Mode algorthm

df_kmode = results[['AGE','FEMALE', 'RACE', 'HOSP_REGION', 'ZIPINC_QRTL', 'LOS', 'TOTCHG',
'PAY1']]
```

```python
# Creating bins for Length of stay column

## labels = ['0-1', '2-3', '4-9', '10 and over']


bins = [-1, 1, 3, 9, float('inf')]  # Using -1 as the lower bound to include 0
labels = [1, 2, 3, 4]

# Create a new column with the binned categories
df_kmode.loc[:, 'LOS'] = pd.cut(df_kmode['LOS'], bins=bins, labels=labels)
```

```python
# Creating bins for "Total charge" feature

## labels = ['0-12000', '12,000-20,000' '20,000-40,000', '40,000-80,000', '80,000+']


bins = [-1, 12000, 20000, 40000, 80000, float('inf')]  # Using -1 as the lower bound to include 0
labels = [1, 2, 3, 4, 5]

# Create a new column with the binned categories
df_kmode.loc[:, 'TOTCHG'] = pd.cut(df_kmode['TOTCHG'], bins=bins, labels=labels)
```

```python
# Creating bins for "AGE" feature of the model

## labels = ['0-1', '1-10' '11-20']


bins = [-1, 1, 10, 21]  # Using -1 as the lower bound to include 0
labels = [1, 2, 3]

# Create a new column with the binned categories
df_kmode.loc[:, 'AGE'] = pd.cut(df_kmode['AGE'], bins=bins, labels=labels)
```

```python
# Install K-Modes ML library

!pip install kmodes
```

```python
# Import necessary libraries from KModes and create clusters
from kmodes.kmodes import KModes
```

```python
import pandas as pd
from sklearn.preprocessing import OneHotEncoder

encoder = OneHotEncoder(sparse=False)
df_encoded = encoder.fit_transform(df_kmode)

km = KModes(n_clusters=3, init='Huang', n_init=5, verbose=1)
```

```python
#Check of the clusters created in vector format

df_encoded
```

```python
# Fit the cluster

clusters = km.fit_predict(df_encoded)
```

```python
# Checking the optimal number of clusters using Elbow method

import matplotlib.pyplot as plt
from kmodes.kmodes import KModes

costs = []
for num_clusters in range(1, 10):
    km = KModes(n_clusters=num_clusters, init='Huang', n_init=5, verbose=0)
    km.fit(df_encoded)
    costs.append(km.cost_)

plt.plot(range(1, 10), costs, marker='o')
plt.xlabel('Number of clusters')
plt.ylabel('Cost')
plt.title('Elbow Method For Optimal Number of Clusters')
plt.show()
```

```python
# Join the cluster results to the original table

df_kmode['Cluster'] = clusters
```

```python
# Using dataframe df_kmode: save to csv file for future analysis and visualization

df_kmode.to_csv('df_kmode_Asthma_primary_3.csv', index=False)
```

```python
#Importing Library for kMeans
import pandas as pd
from sklearn.preprocessing import StandardScaler
from sklearn.cluster import KMeans
import matplotlib.pyplot as plt
```

```python
#Scaling the dataset
scaler = StandardScaler()
scaled_features = scaler.fit_transform(df_asthma)
```

```python
# Elbow Method
inertia = []
for i in range(1, 11):
    kmeans = KMeans(n_clusters=i, random_state=0)
    kmeans.fit(scaled_features)
    inertia.append(kmeans.inertia_)

plt.plot(range(1, 11), inertia)
plt.title('Elbow Method')
plt.xlabel('Number of clusters')
plt.ylabel('Inertia')
plt.show()
```

```python
# Initialize KMeans with k clusters
kmeans = KMeans(n_clusters=5, random_state=0)

# Fit the model
kmeans.fit(scaled_features)

# Get the cluster labels for each data point
labels = kmeans.labels_
```

```python
# Add the cluster labels as a new column in the original DataFrame
df_asthma['Cluster'] = labels


#Export result to CSV
df_asthma.to_csv("/df_cluster_new.csv")
```

--------------------------------------------------------End--------------------------------------------------------