# Data Collection and Preprocessing Phase

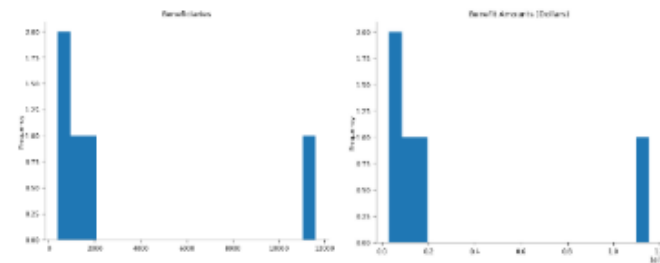| Date | 27 NOVEMBER 2024 |
|---|---|
| Team ID | FACULTY |
| Project Title | Unemployed Insurance Beneficiary Forecasting. |
| Maximum Marks | 6 Marks |

**Preprocessing Template**

The images will be preprocessed by resizing, normalizing, augmenting,Dataset variables will be statistically analyzed to identify patterns and outliers, with Python,employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions, ensuring robust and efficient performance across various computer vision tasks.
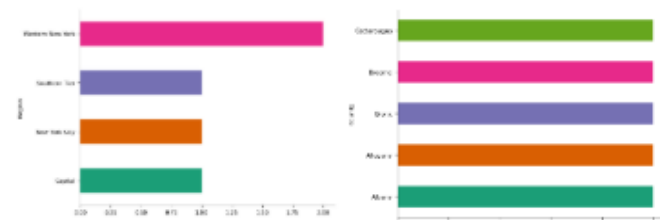
| Section | Description |
|---|---|
| Data Overview | Dimension:<br>13760 rows × 6 columns<br>Descriptive statistics:<br><br>**Descriptive Analysis**<br><br>`df.describe()`<br><br>

| | Beneficiaries | Benefit Amounts (Dollars) | Beneficiaries_diff | Date |
|---|---|---|---|---|
| count | 13759.000000 | 1.375900e+04 | 13759.000000 | 13759 |
| mean | 3858.499891 | 3.847134e+06 | -0.094484 | 2009-11-30 12:10:43.651428352 |
| min | 0.000000 | 0.000000e+00 | -50200.000000 | 2001-01-01 00:00:00 |
| 25% | 600.000000 | 5.700000e+05 | -1200.000000 | 2005-06-01 00:00:00 |
| 50% | 1200.000000 | 1.110000e+06 | 0.000000 | 2009-12-01 00:00:00 |
| 75% | 2800.000000 | 2.720000e+06 | 1500.000000 | 2014-06-01 00:00:00 |
| max | 50700.000000 | 5.681000e+07 | 49000.000000 | 2018-11-01 00:00:00 |
| std | 6557.760758 | 6.878863e+06 | 9431.460815 | NaN |

Distributions

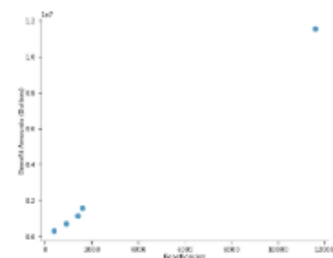## Distributions



## Categorical distributions



## 2-d distributions



## Time series



## Values



## 2-d categorical distributions



## Faceted distributions

| Outliers and Anomalies | -- |
| --- | --- |
| **Data Preprocessing Code Screenshots** | |
| Loading Data | **Read The Dataset**<br><br>`[2] df = pd.read_csv(r"/content/insurance_unemployed_data.csv")`<br><br>`df.head()`<br><br><table><thead><tr><th></th><th>Year</th><th>Month</th><th>Region</th><th>County</th><th>Beneficiaries</th><th>Benefit Amounts (Dollars)</th></tr></thead><tbody><tr><td>0</td><td>2018</td><td>11</td><td>Capital</td><td>Albany</td><td>1600</td><td>1570000</td></tr><tr><td>1</td><td>2018</td><td>11</td><td>Western New York</td><td>Allegany</td><td>400</td><td>300000</td></tr><tr><td>2</td><td>2018</td><td>11</td><td>New York City</td><td>Bronx</td><td>11600</td><td>11530000</td></tr><tr><td>3</td><td>2018</td><td>11</td><td>Southern Tier</td><td>Broome</td><td>1400</td><td>1150000</td></tr><tr><td>4</td><td>2018</td><td>11</td><td>Western New York</td><td>Cattaraugus</td><td>900</td><td>710000</td></tr></tbody></table> |
| Handling Missing Data | **Checking for missing values**<br><br>`[6]  print(df.isnull().sum())`<br><br>```<br>Year                          0<br>Month                         0<br>Region                        0<br>County                        0<br>Beneficiaries                 0<br> Benefit Amounts (Dollars)    0<br>dtype: int64<br>``` |
| Checking Duplicates | **Checking for Duplicates**<br><br>`[7]  df.duplicated().sum()`<br><br>`0`<br><br>`df.info()`<br><br>```<br><class 'pandas.core.frame.DataFrame'><br>RangeIndex: 13760 entries, 0 to 13759<br>Data columns (total 6 columns):<br> #   Column                     Non-Null Count  Dtype<br>---  ------                     --------------  -----<br> 0   Year                       13760 non-null  int64<br> 1   Month                      13760 non-null  int64<br> 2   Region                     13760 non-null  object<br> 3   County                     13760 non-null  object<br> 4   Beneficiaries              13760 non-null  int64<br> 5    Benefit Amounts (Dollars) 13760 non-null  int64<br>dtypes: int64(4), object(2)<br>memory usage: 645.1+ KB<br>```<br><br>`[5] df.shape`<br><br>`(13760, 6)` |

| Feature Engineering | ## Splitting Dataset into Train and Test Sets |
|---|---|
| | ```[9]  df.dropna(inplace=True)```<br><br>```[10]  train_size=int (len(df)*0.8)```<br>```       train,test=df[:train_size],df[train_size:]```<br><br>**create differenced column**<br><br>```[11]  train['Beneficiaries_diff']=train['Beneficiaries'].diff()```<br>```       print(train['Beneficiaries_diff'])```<br><br>```0              NaN```<br>```1          -1200.0```<br>```2          11200.0```<br>```3         -10200.0```<br>```4           -500.0```<br>```            ...```<br>```11003         0.0```<br>```11004       500.0```<br>```11005      6700.0```<br>```11006     -7300.0```<br>```11007      -200.0```<br>```Name: Beneficiaries_diff, Length: 11008, dtype: float64```<br><br>**Distributions**<br><br><br><br> |