# Data Collection and Preprocessing Phase
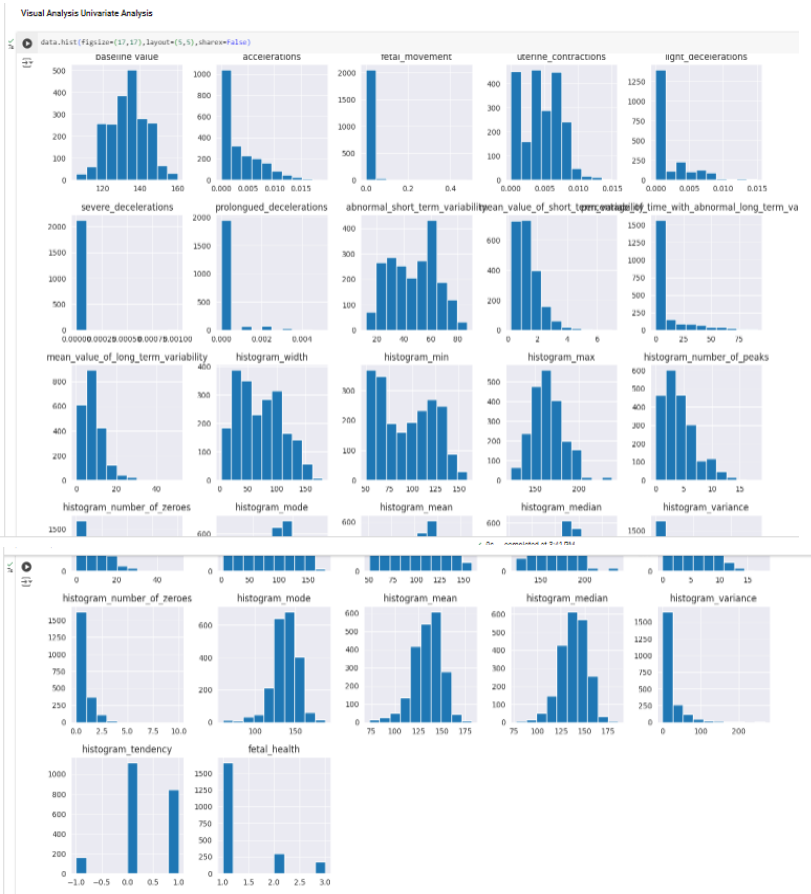
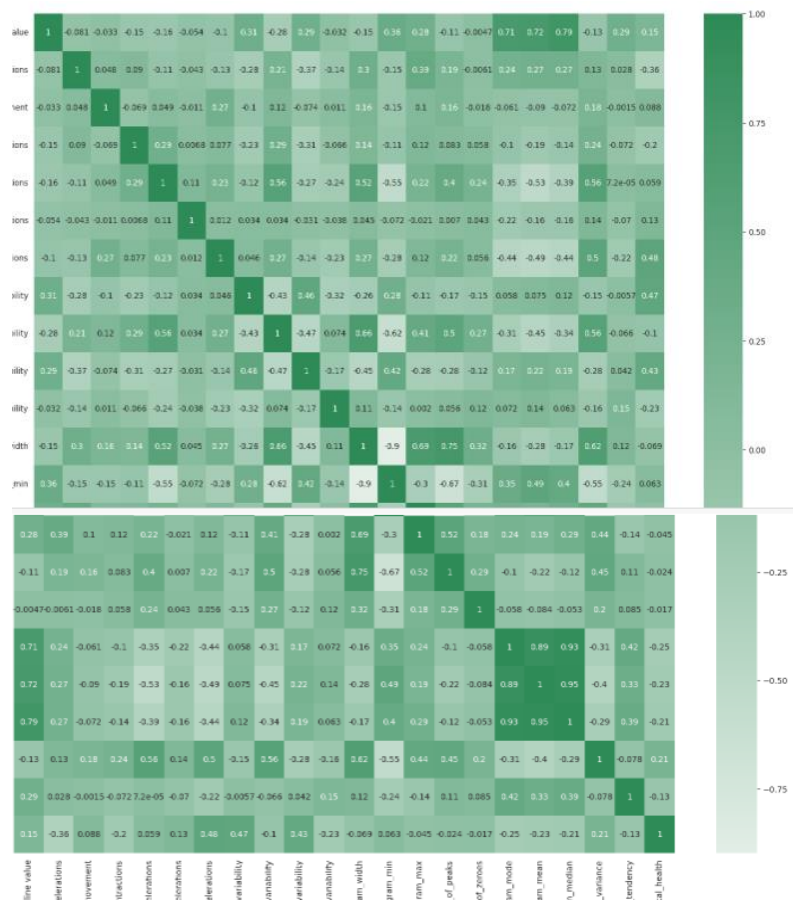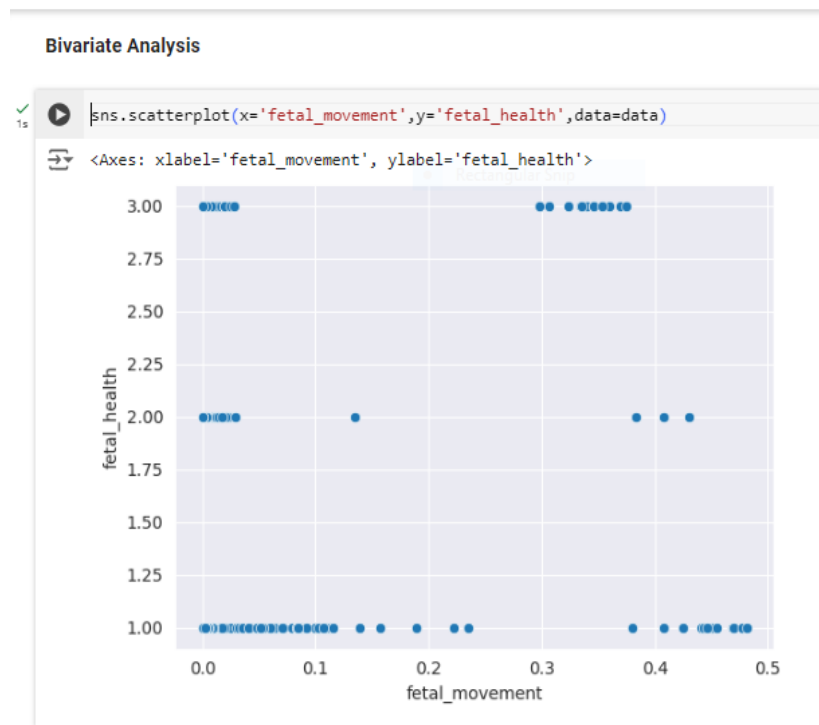| Date | 10 JULY 2024 |
|---|---|
| Team ID | FACULTY |
| Project Title | Fetal AI: Using Machine Learning To Predict And Monitor Fetal Health. |
| Maximum Marks | 6 Marks |

**Data Exploration and Preprocessing Report**

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.

| Section | Description |
|---|---|
| Data Overview | Dimension:<br>2126 rows × 22columns<br>Descriptive statistics:<br><br> |
|  |  |

| Univariate Analysis |  |
|---|---|
| | |

| | |
|---|---|
| Bivariate Analysis | **Bivariate Analysis**<br><br>`sns.scatterplot(x='fetal_movement',y='fetal_health',data=data)`<br><br>`<Axes: xlabel='fetal_movement', ylabel='fetal_health'>`<br><br> |
| Multivariate Analysis |  |

| Outliers and Anomalies | -- |
|---|---|

## Data Preprocessing Code Screenshots

| Loading Data | **Read the Dataset** <br><br> [2] `data=pd.read_csv('/content/fetalhealth.csv')` <br><br> [3] `data.head()` <br><br> (table preview: baseline value, accelerations, fetal_movement, uterine_contractions, light_decelerations, severe_decelerations, prolongued_decelerations) <br> 0: 120.0, 0.000, 0.0, 0.000, 0.000, 0.0, 0.0 <br> 1: 132.0, 0.006, 0.0, 0.006, 0.003, 0.0, 0.0 <br> 2: 133.0, 0.003, 0.0, 0.008, 0.003, 0.0, 0.0 <br> 3: 134.0, 0.003, 0.0, 0.008, 0.003, 0.0, 0.0 <br> 4: 132.0, 0.007, 0.0, 0.008, 0.000, 0.0, 0.0 <br> 5 rows × 22 columns |
|---|---|

| Handling Missing Data | **Data Preparation 1.Handling Missing Values:** <br><br> [5] `data.info()` <br><br> ```<class 'pandas.core.frame.DataFrame'>``` <br> RangeIndex: 2126 entries, 0 to 2125 <br> Data columns (total 22 columns): <br> # Column, Non-Null Count, Dtype <br> 0 baseline value — 2126 non-null float64 <br> 1 accelerations — 2126 non-null float64 <br> 2 fetal_movement — 2126 non-null float64 <br> 3 uterine_contractions — 2126 non-null float64 <br> 4 light_decelerations — 2126 non-null float64 <br> 5 severe_decelerations — 2126 non-null float64 <br> 6 prolongued_decelerations — 2126 non-null float64 <br> 7 abnormal_short_term_variability — 2126 non-null float64 <br> 8 mean_value_of_short_term_variability — 2126 non-null float64 <br> 9 percentage_of_time_with_abnormal_long_term_variability — 2126 non-null float64 <br> 10 mean_value_of_long_term_variability — 2126 non-null float64 <br> 11 histogram_width — 2126 non-null float64 <br> 12 histogram_min — 2126 non-null float64 <br> 13 histogram_max — 2126 non-null float64 <br> 14 histogram_number_of_peaks — 2126 non-null float64 <br> 15 histogram_number_of_zeroes — 2126 non-null float64 <br> 16 histogram_mode — 2126 non-null float64 <br> 17 histogram_mean — 2126 non-null float64 <br> 18 histogram_median — 2126 non-null float64 <br> 19 histogram_variance — 2126 non-null float64 <br> 20 histogram_tendency — 2126 non-null float64 <br> 21 fetal_health — 2126 non-null float64 <br> dtypes: float64(22) <br> memory usage: 365.5 KB <br><br> [6] `data.isnull().sum()` <br> baseline value 0 <br> accelerations 0 <br> fetal_movement 0 <br> uterine_contractions 0 <br> light_decelerations 0 <br> severe_decelerations 0 <br> prolongued_decelerations 0 <br> abnormal_short_term_variability 0 <br> mean_value_of_short_term_variability 0 <br> percentage_of_time_with_abnormal_long_term_variability 0 <br> mean_value_of_long_term_variability 0 <br> histogram_width 0 <br> histogram_min 0 <br> histogram_max 0 <br> histogram_number_of_peaks 0 <br> histogram_number_of_zeroes 0 <br> histogram_mode 0 <br> histogram_mean 0 <br> histogram_median 0 <br> histogram_variance 0 <br> histogram_tendency 0 <br> fetal_health 0 <br> dtype: int64 |
|---|---|

| | |
|---|---|
| Handling Imbalance Data | **2.Handling Imbalance Data**<br><br>[7] `#Evaluating the target and find out if our data is imbalanced or not`<br>`data['fetal_health'].value_counts()`<br><br>fetal_health<br>1.0    1655<br>2.0     295<br>3.0     176<br>Name: count, dtype: int64<br><br>[8] `colours=["#f7b2b0","#8f7198", "#003f5c"]`<br>`sns.countplot(data= data, x="fetal_health",palette=colours)`<br><br>`<Axes: xlabel='fetal_health', ylabel='count'>`<br><br> |
| Feature Engineering | **Feature Selection**<br><br>`data.drop(columns=['histogram_mean'],axis=1,inplace=True)`<br><br>[17] `data.shape`<br><br>`(2126, 21)`<br><br>[18] `data.corr()["fetal_health"].sort_values(ascending=False)`<br><br>fetal_health                                                     1.000000<br>prolongued_decelerations                                         0.484859<br>abnormal_short_term_variability                                  0.471191<br>percentage_of_time_with_abnormal_long_term_variability           0.426146<br>histogram_variance                                               0.206630<br>baseline value                                                   0.148151<br>severe_decelerations                                             0.131934<br>fetal_movement                                                   0.088010<br>histogram_min                                                    0.063175<br>light_decelerations                                              0.058870<br>histogram_number_of_zeroes                                      -0.016682<br>histogram_number_of_peaks                                      -0.023666<br>histogram_max                                                   -0.045265<br>histogram_width                                                -0.068789<br>mean_value_of_short_term_variability                           -0.103382<br>histogram_tendency                                             -0.131976<br>uterine_contractions                                           -0.204894<br>histogram_median                                               -0.205033<br>mean_value_of_long_term_variability                            -0.226797<br>histogram_mode                                                 -0.250412<br>accelerations                                                  -0.364066<br>Name: fetal_health, dtype: float64<br><br>`new_data=data.loc[:,["prolongued_decelerations","abnormal_short_term_variability",`<br>`"percentage_of_time_with_abnormal_long_term_variability"]]`<br><br>[20] `new_data.head()`<br><br><table><tr><th></th><th>prolongued_decelerations</th><th>abnormal_short_term_variability</th><th>percentage_of_time_with_abnormal_long_term_variability</th></tr><tr><td>0</td><td>0.0</td><td>73.0</td><td>43.0</td></tr><tr><td>1</td><td>0.0</td><td>17.0</td><td>0.0</td></tr><tr><td>2</td><td>0.0</td><td>16.0</td><td>0.0</td></tr><tr><td>3</td><td>0.0</td><td>16.0</td><td>0.0</td></tr><tr><td>4</td><td>0.0</td><td>16.0</td><td>0.0</td></tr></table> |

**Scaling Data**

```python
[21] x=data.drop(columns=['fetal_health'])
     y=data["fetal_health"]
     from sklearn.preprocessing import MinMaxScaler
     scale=MinMaxScaler()
     x_scaled=pd.DataFrame(scale.fit_transform(x),columns=x.columns)
     x_scaled.head()
```

| | baseline value | accelerations | fetal_movement | uterine_contractions | light_decelerations | severe_decelerations | prolongued_decelerations | a |
|---|---|---|---|---|---|---|---|---|
| 0 | 0.259259 | 0.000000 | 0.0 | 0.000000 | 0.0 | 0.0 | 0.0 | |
| 1 | 0.481481 | 0.315789 | 0.0 | 0.400000 | 0.2 | 0.0 | 0.0 | |
| 2 | 0.500000 | 0.157895 | 0.0 | 0.533333 | 0.2 | 0.0 | 0.0 | |
| 3 | 0.518519 | 0.157895 | 0.0 | 0.533333 | 0.2 | 0.0 | 0.0 | |
| 4 | 0.481481 | 0.368421 | 0.0 | 0.533333 | 0.0 | 0.0 | 0.0 | |

Next steps: [ Generate code with x_scaled ]  [ View recommended plots ]

```python
[22] data.shape
     (2126, 21)
```

**Splitting data into Train and Test**

```python
[23] from sklearn.metrics import accuracy_score,classification_report,confusion_matrix
```

```python
[24] from sklearn.model_selection import train_test_split
     x_train,x_test,y_train,y_test=train_test_split(x,y,test_size=0.3,random_state=42)
     x_train.shape,x_test.shape
     ((1488, 20), (638, 20))
```

**Applying SMOTE for balancing the Data**

```python
[25] pip install imblearn
     Collecting imblearn
       Downloading imblearn-0.0-py2.py3-none-any.whl (1.9 kB)
     Requirement already satisfied: imbalanced-learn in /usr/local/lib/python3.10/dist-pac
     Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.10/dist-packag
     Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.10/dist-package
     Requirement already satisfied: scikit-learn>=1.0.2 in /usr/local/lib/python3.10/dist-
     Requirement already satisfied: joblib>=1.1.1 in /usr/local/lib/python3.10/dist-packag
     Requirement already satisfied: threadpoolctl>=2.0.0 in /usr/local/lib/python3.10/dist
     Installing collected packages: imblearn
     Successfully installed imblearn-0.0
```

```python
[26] from imblearn.over_sampling import SMOTE
     smote=SMOTE()
```

```python
[27] x_train_smote,y_train_smote=smote.fit_resample(x_train.astype('float'),y_train)
```

```
[27] x_train_smote,y_train_smote=smote.fit_resample(x_train.astype('float'),y_train)
```

```
[28] print(x_train.columns)
```

```
Index(['baseline value', 'accelerations', 'fetal_movement',
       'uterine_contractions', 'light_decelerations', 'severe_decelerations',
       'prolongued_decelerations', 'abnormal_short_term_variability',
       'mean_value_of_short_term_variability',
       'percentage_of_time_with_abnormal_long_term_variability',
       'mean_value_of_long_term_variability', 'histogram_width',
       'histogram_min', 'histogram_max', 'histogram_number_of_peaks',
       'histogram_number_of_zeroes', 'histogram_mode', 'histogram_median',
       'histogram_variance', 'histogram_tendency'],
      dtype='object')
```

```
[29] from collections import Counter
     print("Before SMOTE:",Counter(y_train))
     print("After SMOTE:",Counter(y_train_smote))
```

```
Before SMOTE: Counter({1.0: 1159, 2.0: 194, 3.0: 135})
After SMOTE: Counter({1.0: 1159, 3.0: 1159, 2.0: 1159})
```