# Nompumelelo Ngwenya

# Training a machine learning model using the Employee Attrition data.
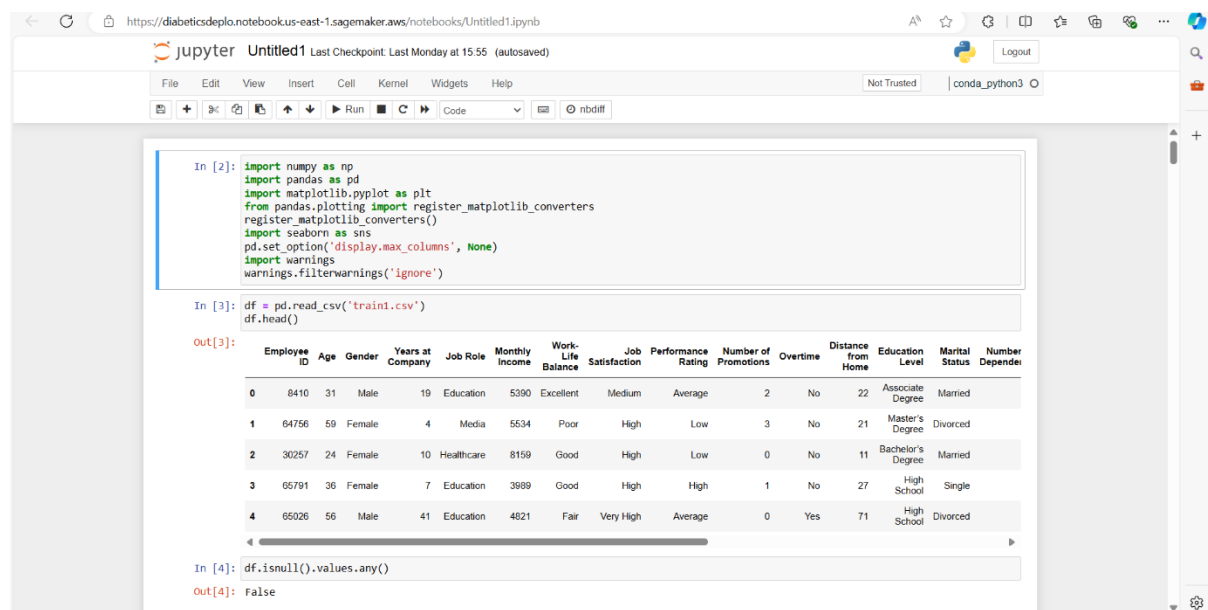
Introduction and Overview:

This document aims to guide stakeholders through the process of deploying a machine learning model trained on the Employees Attrition (Train) dataset on AWS.

The dataset comprises 74,498 samples and It contains detailed information about various aspects of an employee's profile, including demographics, job-related features, and personal circumstances. Each record includes a unique Employee ID and features that influence employee attrition. The aim is to understand the factors contributing to attrition and develop predictive models to identify at-risk employees.

This dataset is ideal for HR analytics, machine learning model development, and demonstrating advanced data analysis techniques. It provides a comprehensive and realistic view of the factors affecting employee retention, making it a valuable resource for researchers and practitioners in the field of human resources and organizational development.

Data source: (Employee Attrition Classification Dataset (kaggle.com))

Importing the necessary python libraries and data extracting from the data source

# Exploratory data analysis(EDA) and data cleaning by checking for any null values

Jupyter Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

File　Edit　View　Insert　Cell　Kernel　Widgets　Help

Not Trusted | conda_python3

In [5]: `df.describe()`

Out[5]:

| | Employee ID | Age | Years at Company | Monthly Income | Number of Promotions | Distance from Home | Number of Dependents | Company Tenure |
|---|---|---|---|---|---|---|---|---|
| count | 59598.000000 | 59598.000000 | 59598.000000 | 59598.000000 | 59598.000000 | 59598.000000 | 59598.000000 | 59598.000000 |
| mean | 37227.118729 | 38.565875 | 15.753901 | 7302.397983 | 0.832578 | 50.007651 | 1.648075 | 55.758415 |
| std | 21519.150028 | 12.079673 | 11.245981 | 2151.457423 | 0.994991 | 28.466459 | 1.555689 | 25.411090 |
| min | 1.000000 | 18.000000 | 1.000000 | 1316.000000 | 0.000000 | 1.000000 | 0.000000 | 2.000000 |
| 25% | 18580.250000 | 28.000000 | 7.000000 | 5658.000000 | 0.000000 | 25.000000 | 0.000000 | 36.000000 |
| 50% | 37209.500000 | 39.000000 | 13.000000 | 7354.000000 | 1.000000 | 50.000000 | 1.000000 | 56.000000 |
| 75% | 55876.750000 | 49.000000 | 23.000000 | 8880.000000 | 2.000000 | 75.000000 | 3.000000 | 76.000000 |
| max | 74498.000000 | 59.000000 | 51.000000 | 16149.000000 | 4.000000 | 99.000000 | 6.000000 | 128.000000 |

In [6]: `df.isna().sum()`

Out[6]:
```
Employee ID             0
Age                     0
Gender                  0
Years at Company        0
Job Role                0
Monthly Income          0
Work-Life Balance       0
Job Satisfaction        0
Performance Rating      0
Number of Promotions    0
Overtime                0
Distance from Home      0
Education Level         0
Marital Status          0
Number of Dependents    0
Job Level               0
Company Size            0
Company Tenure          0
```

---

Jupyter Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

File　Edit　View　Insert　Cell　Kernel　Widgets　Help

Not Trusted | conda_python3

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 75% | 55876.750000 | 49.000000 | 23.000000 | 8880.000000 | 2.000000 | 75.000000 | 3.000000 | 76.000000 |
| max | 74498.000000 | 59.000000 | 51.000000 | 16149.000000 | 4.000000 | 99.000000 | 6.000000 | 128.000000 |

In [6]: `df.isna().sum()`

Out[6]:
```
Employee ID                 0
Age                         0
Gender                      0
Years at Company            0
Job Role                    0
Monthly Income              0
Work-Life Balance           0
Job Satisfaction            0
Performance Rating          0
Number of Promotions        0
Overtime                    0
Distance from Home          0
Education Level             0
Marital Status              0
Number of Dependents        0
Job Level                   0
Company Size                0
Company Tenure              0
Remote Work                 0
Leadership Opportunities    0
Innovation Opportunities    0
Company Reputation          0
Employee Recognition        0
Attrition                   0
dtype: int64
```

In [7]: `df.shape`

Out[7]: `(59598, 24)`

```
dtype: int64
```

In [7]: `df.shape`

Out[7]: (59598, 24)

In [8]: `df.dtypes`

Out[8]:
```
Employee ID                int64
Age                        int64
Gender                     object
Years at Company           int64
Job Role                   object
Monthly Income             int64
Work-Life Balance          object
Job Satisfaction           object
Performance Rating         object
Number of Promotions       int64
Overtime                   object
Distance from Home         int64
Education Level            object
Marital Status             object
Number of Dependents       int64
Job Level                  object
Company Size               object
Company Tenure             int64
Remote Work                object
Leadership Opportunities   object
Innovation Opportunities   object
Company Reputation         object
Employee Recognition       object
Attrition                  object
dtype: object
```

In [9]: `#Get the information about the datasets`
`df.info()`

In [9]: `#Get the information about the datasets`
`df.info()`

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 59598 entries, 0 to 59597
Data columns (total 24 columns):
 #   Column                    Non-Null Count  Dtype
---  ------                    --------------  -----
 0   Employee ID               59598 non-null  int64
 1   Age                       59598 non-null  int64
 2   Gender                    59598 non-null  object
 3   Years at Company          59598 non-null  int64
 4   Job Role                  59598 non-null  object
 5   Monthly Income            59598 non-null  int64
 6   Work-Life Balance         59598 non-null  object
 7   Job Satisfaction          59598 non-null  object
 8   Performance Rating        59598 non-null  object
 9   Number of Promotions      59598 non-null  int64
 10  Overtime                  59598 non-null  object
 11  Distance from Home        59598 non-null  int64
 12  Education Level           59598 non-null  object
 13  Marital Status            59598 non-null  object
 14  Number of Dependents      59598 non-null  int64
 15  Job Level                 59598 non-null  object
 16  Company Size              59598 non-null  object
 17  Company Tenure            59598 non-null  int64
 18  Remote Work               59598 non-null  object
 19  Leadership Opportunities  59598 non-null  object
 20  Innovation Opportunities  59598 non-null  object
 21  Company Reputation        59598 non-null  object
 22  Employee Recognition      59598 non-null  object
 23  Attrition                 59598 non-null  object
dtypes: int64(8), object(16)
memory usage: 10.9+ MB
```

In [10]: `#Get the count of the number of Employee that stayed or left the company`
`df['Attrition'].value_counts()`

Jupyter  Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

Code

```
dtypes: int64(8), object(16)
memory usage: 10.9+ MB
```

In [10]: #Get the count of the number of Employee that stayed or left the company
df['Attrition'].value_counts()

Out[10]: Attrition
Stayed    31260
Left      28338
Name: count, dtype: int64

In [11]: #Get all the data types and their unique values
for column in df.columns:
    if df[column].dtype == object:
        print(str(column)+ ' : '+ str(df[column].unique()))
        print(df[column].value_counts())
        print('_____')

```
Gender : ['Male' 'Female']
Gender
Male      32739
Female    26859
Name: count, dtype: int64

Job Role : ['Education' 'Media' 'Healthcare' 'Technology' 'Finance']
Job Role
Technology    15507
Healthcare    13642
Education     12490
Media          9574
Finance        8385
Name: count, dtype: int64

Work-Life Balance : ['Excellent' 'Poor' 'Good' 'Fair']
Work-Life Balance
Good    22528
Fair    18046
```
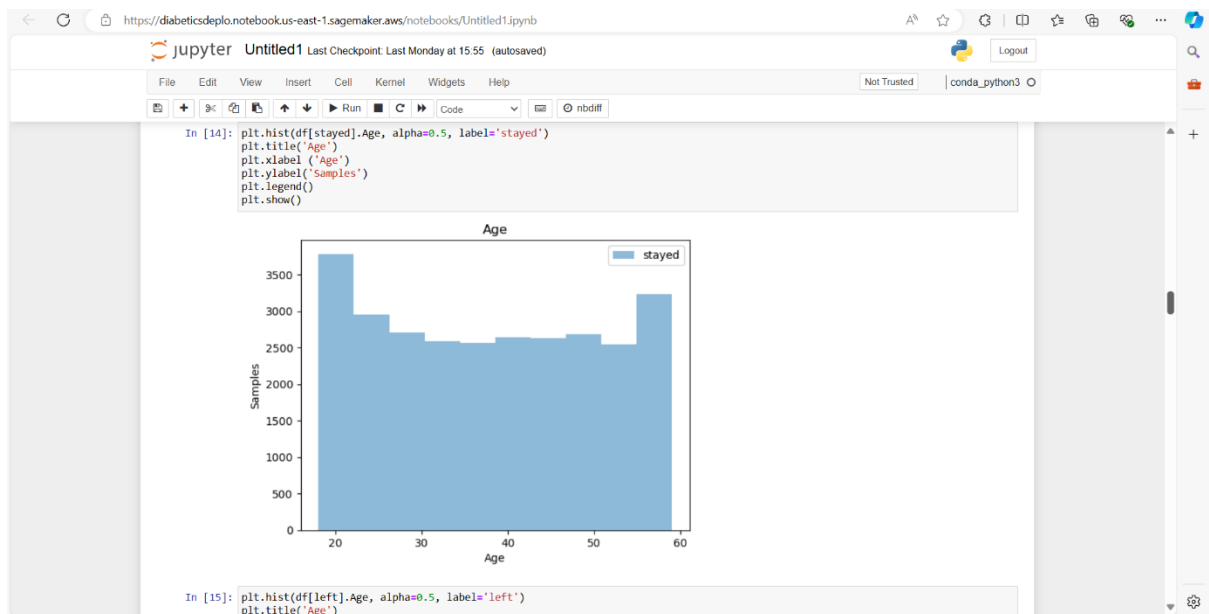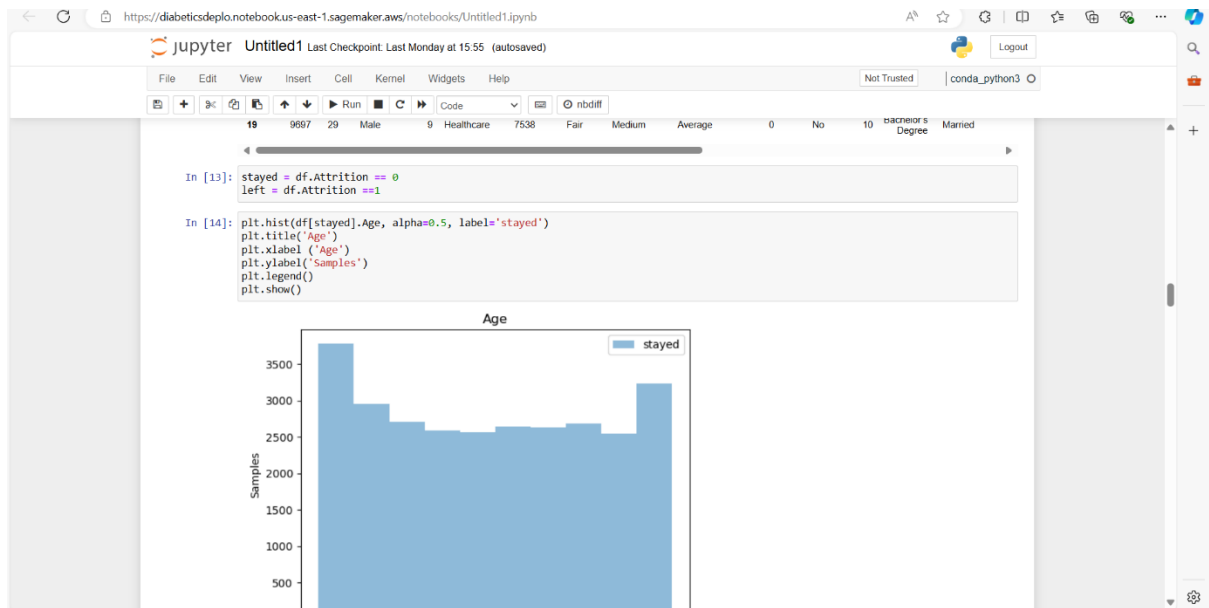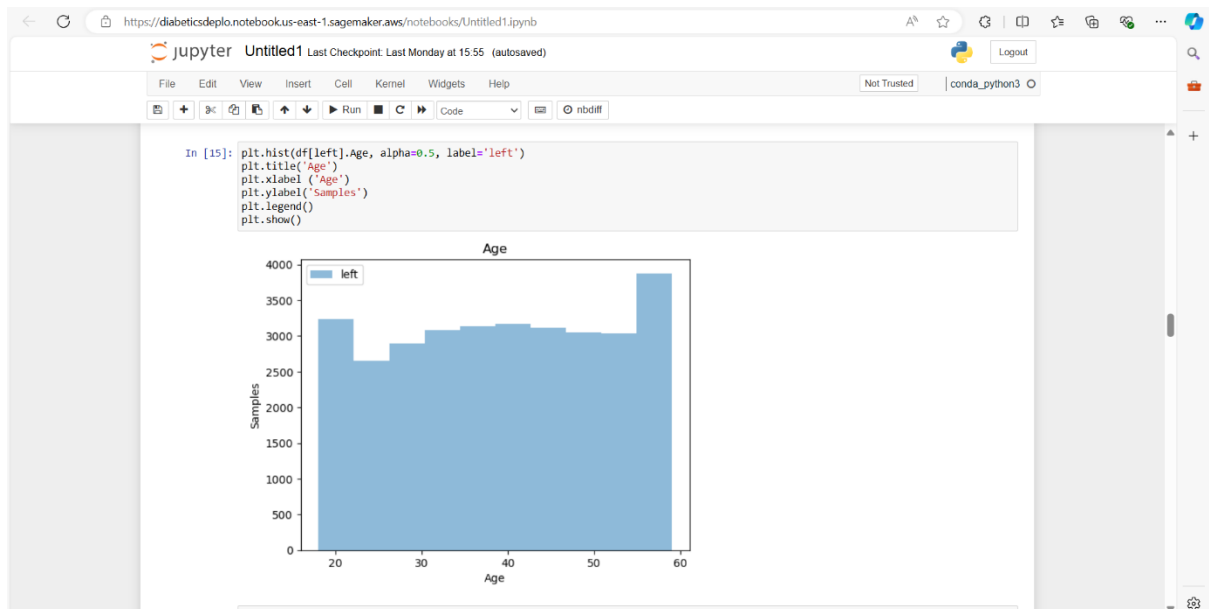
Encoding the attrition. Employees that stayed as 1, employees who left as 0
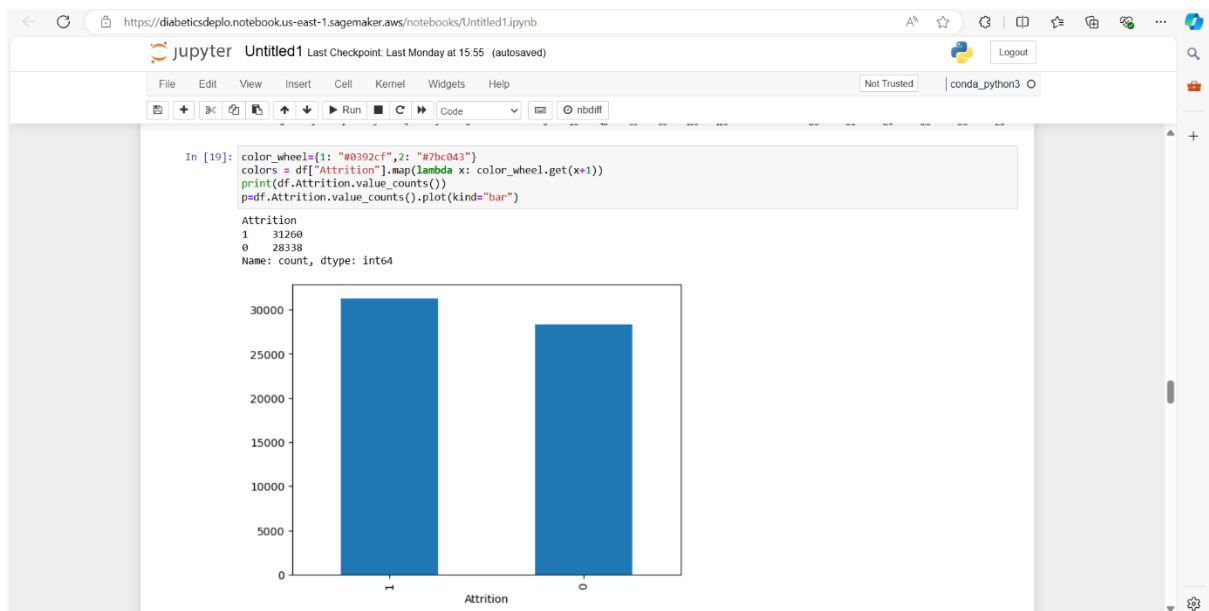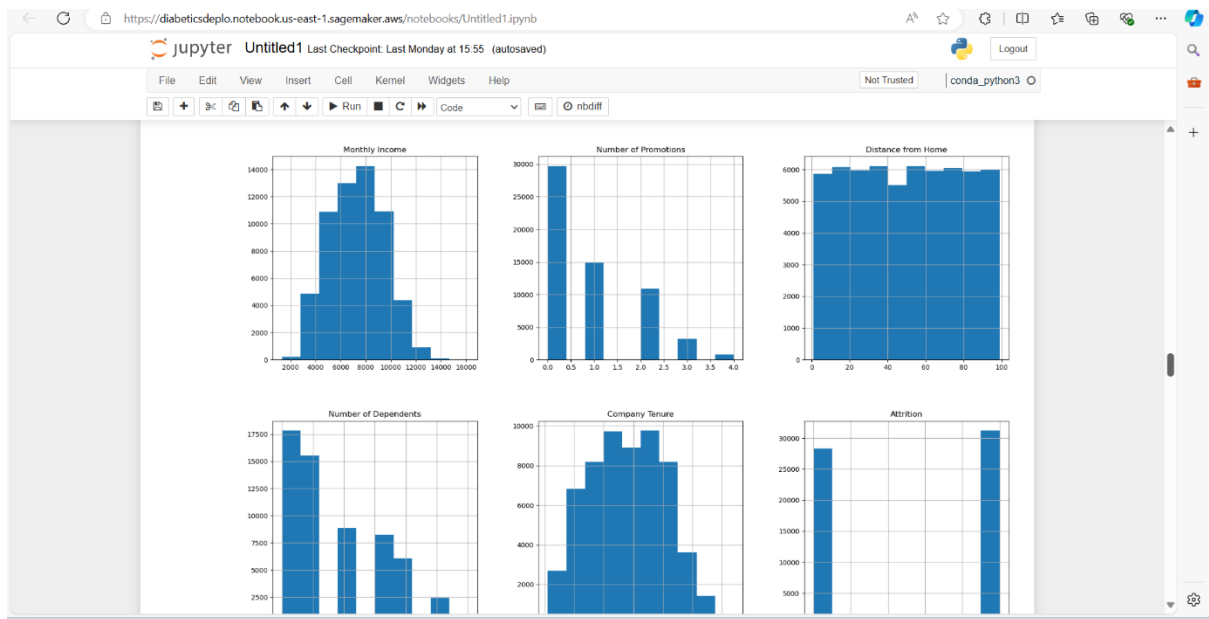
Jupyter  Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

Code

In [12]: #lets convert this attrition to label
from sklearn.preprocessing import LabelEncoder
le = LabelEncoder()
df['Attrition']= le.fit_transform(df['Attrition'])
df.head(20)

Out[12]:

| | Employee ID | Age | Gender | Years at Company | Job Role | Monthly Income | Work-Life Balance | Job Satisfaction | Performance Rating | Number of Promotions | Overtime | Distance from Home | Education Level | Marital Status | Numb Depend |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 8410 | 31 | Male | 19 | Education | 5390 | Excellent | Medium | Average | 2 | No | 22 | Associate Degree | Married | |
| 1 | 64756 | 59 | Female | 4 | Media | 5534 | Poor | High | Low | 3 | No | 21 | Master's Degree | Divorced | |
| 2 | 30257 | 24 | Female | 10 | Healthcare | 8159 | Good | High | Low | 0 | No | 11 | Bachelor's Degree | Married | |
| 3 | 65791 | 36 | Female | 7 | Education | 3989 | Good | High | High | 1 | No | 27 | High School | Single | |
| 4 | 65026 | 56 | Male | 41 | Education | 4821 | Fair | Very High | Average | 0 | Yes | 71 | High School | Divorced | |
| 5 | 24368 | 38 | Female | 3 | Technology | 9977 | Fair | High | Below Average | 3 | No | 37 | Bachelor's Degree | Married | |
| 6 | 64970 | 47 | Male | 23 | Education | 3681 | Fair | High | High | 1 | Yes | 75 | High School | Divorced | |
| 7 | 36999 | 48 | Male | 16 | Finance | 11223 | Excellent | Very High | High | 2 | No | 5 | Master's Degree | Married | |
| 8 | 32714 | 57 | Male | 44 | Education | 3773 | Good | Medium | High | 1 | Yes | 39 | High School | Married | |
| 9 | 15944 | 24 | Female | 1 | Healthcare | 7319 | Poor | High | Average | 1 | Yes | 57 | PhD | Single | |
| 10 | 29972 | 30 | Female | 12 | Education | 5443 | Good | High | Average | 1 | No | 51 | High School | Single | |
| 11 | 9063 | 29 | Female | 6 | Healthcare | 8950 | Poor | Medium | Low | 2 | No | 26 | Master's Degree | Single | |

jupyter Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

Logout

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

Not Trusted    conda_python3

| 19 | 9697 | 29 | Male | 9 | Healthcare | 7538 | Fair | Medium | Average | 0 | No | 10 | Bachelor's Degree | Married |

In [13]:
```python
stayed = df.Attrition == 0
left = df.Attrition ==1
```

In [14]:
```python
plt.hist(df[stayed].Age, alpha=0.5, label='stayed')
plt.title('Age')
plt.xlabel ('Age')
plt.ylabel('Samples')
plt.legend()
plt.show()
```



---

jupyter Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

Logout

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

Not Trusted    conda_python3

In [14]:
```python
plt.hist(df[stayed].Age, alpha=0.5, label='stayed')
plt.title('Age')
plt.xlabel ('Age')
plt.ylabel('Samples')
plt.legend()
plt.show()
```



In [15]:
```python
plt.hist(df[left].Age, alpha=0.5, label='left')
plt.title('Age')
```

Jupyter **Untitled1** Last Checkpoint: Last Monday at 15:55 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | conda_python3

Code

```
In [16]: le.classes_
```

```
Out[16]: array(['Left', 'Stayed'], dtype=object)
```

```
In [17]: #Attrition{stayed=1, left=1}
         Attrition_employee=le.classes_
         print(Attrition_employee)
```

```
['Left' 'Stayed']
```

```
In [18]: #plotting the distributions

         p = df.hist(figsize=(20,20))
```



---

Jupyter **Untitled1** Last Checkpoint: Last Monday at 15:55 (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted | conda_python3

Code

```
In [15]: plt.hist(df[left].Age, alpha=0.5, label='left')
         plt.title('Age')
         plt.xlabel ('Age')
         plt.ylabel('Samples')
         plt.legend()
         plt.show()
```

Encoding the columns to prepare for training.

```python
In [20]: #Define the columns to be label encoded
labels_cols = ['Gender','Job Role','Overtime','Education Level','Marital Status','Company Size','Remote Work',
               'Leadership Opportunities', 'Innovation Opportunities','Work-Life Balance', 'Job Satisfaction','Performance Rating'
               'Company Reputation','Job Level', 'Employee Recognition']

#initialize label encoders
label_encoders = {col: LabelEncoder() for col in labels_cols}

#Apply the label Encoding
for col in labels_cols:
    df[col] = label_encoders[col].fit_transform(df[col])
```

```python
In [21]: df.head(10)
```

Out[21]:

| | Employee ID | Age | Gender | Years at Company | Job Role | Monthly Income | Work-Life Balance | Job Satisfaction | Performance Rating | Number of Promotions | Overtime | Distance from Home | Education Level | Marital Status | Number of Dependents | J Le |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 8410 | 31 | 1 | 19 | 0 | 5390 | 0 | 2 | 0 | 2 | 0 | 22 | 0 | 1 | 0 | |
| 1 | 64756 | 59 | 0 | 4 | 3 | 5534 | 3 | 0 | 3 | 3 | 0 | 21 | 3 | 0 | 3 | |
| 2 | 30257 | 24 | 0 | 10 | 2 | 8159 | 2 | 0 | 3 | 0 | 0 | 11 | 1 | 1 | 3 | |
| 3 | 65791 | 36 | 0 | 7 | 0 | 3989 | 2 | 0 | 2 | 1 | 0 | 27 | 2 | 2 | 2 | |
| 4 | 65026 | 56 | 1 | 41 | 0 | 4821 | 1 | 3 | 0 | 0 | 1 | 71 | 2 | 0 | 0 | |
| 5 | 24368 | 38 | 0 | 3 | 4 | 9977 | 1 | 0 | 1 | 3 | 0 | 37 | 1 | 1 | 0 | |
| 6 | 64970 | 47 | 1 | 23 | 0 | 3681 | 1 | 0 | 2 | 1 | 1 | 75 | 2 | 0 | 3 | |
| 7 | 36999 | 48 | 1 | 16 | 1 | 11223 | 0 | 3 | 2 | 2 | 0 | 5 | 3 | 1 | 4 | |
| 8 | 32714 | 57 | 1 | 44 | 0 | 3773 | 2 | 2 | 2 | 1 | 1 | 39 | 2 | 1 | 4 | |
| 9 | 15944 | 24 | 0 | 1 | 2 | 7319 | 3 | 0 | 0 | 1 | 1 | 57 | 4 | 2 | 4 | |

Splitting data as X and Y, to prepare for training.

```python
In [22]: x=df.iloc[:, :-1]
```

```python
In [23]: x.head()
```

Out[23]:

| | Employee ID | Age | Gender | Years at Company | Job Role | Monthly Income | Work-Life Balance | Job Satisfaction | Performance Rating | Number of Promotions | Overtime | Distance from Home | Education Level | Marital Status | Number of Dependents | J Le |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 8410 | 31 | 1 | 19 | 0 | 5390 | 0 | 2 | 0 | 2 | 0 | 22 | 0 | 1 | 0 | |
| 1 | 64756 | 59 | 0 | 4 | 3 | 5534 | 3 | 0 | 3 | 3 | 0 | 21 | 3 | 0 | 3 | |
| 2 | 30257 | 24 | 0 | 10 | 2 | 8159 | 2 | 0 | 3 | 0 | 0 | 11 | 1 | 1 | 3 | |
| 3 | 65791 | 36 | 0 | 7 | 0 | 3989 | 2 | 0 | 2 | 1 | 0 | 27 | 2 | 2 | 2 | |
| 4 | 65026 | 56 | 1 | 41 | 0 | 4821 | 1 | 3 | 0 | 0 | 1 | 71 | 2 | 0 | 0 | |

```python
In [24]: y=df.iloc[:,-1]
```

```python
In [25]: y.head()
```

```
Out[25]: 0    1
1    1
2    1
3    1
4    1
Name: Attrition, dtype: int64
```

```python
In [26]: from sklearn.model_selection import train_test_split
```

```python
In [27]: xTrain, xTest, yTrain, yTest = train_test_split(x,y, test_size = 0.2)
```

```python
In [28]: xTrain.head(5)
```

jupyter  Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

Logout

File   Edit   View   Insert   Cell   Kernel   Widgets   Help

Not Trusted | conda_python3 ○

```python
In [30]: trainDF=xTrain.join(yTrain)
         trainDF.head(5)
```

Out[30]:

| | Employee ID | Age | Gender | Years at Company | Job Role | Monthly Income | Work-Life Balance | Job Satisfaction | Performance Rating | Number of Promotions | Overtime | Distance from Home | Education Level | Marital Status | Number of Dependents |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 48722 | 73255 | 29 | 1 | 8 | 4 | 7870 | 1 | 3 | 0 | 3 | 0 | 96 | 3 | 2 | 4 |
| 1768 | 29901 | 26 | 0 | 8 | 0 | 3685 | 2 | 0 | 0 | 0 | 0 | 57 | 2 | 2 | 0 |
| 1516 | 4409 | 39 | 0 | 25 | 3 | 5980 | 0 | 2 | 0 | 1 | 0 | 31 | 3 | 0 | 1 |
| 23326 | 53341 | 28 | 1 | 14 | 1 | 7857 | 2 | 2 | 0 | 0 | 1 | 54 | 1 | 2 | 1 |
| 7246 | 56311 | 55 | 1 | 6 | 2 | 7678 | 3 | 3 | 0 | 3 | 0 | 55 | 2 | 0 | 1 |

```python
In [31]: testDF = xTest.join(yTest)
         testDF.head(5)
```

Out[31]:

| | Employee ID | Age | Gender | Years at Company | Job Role | Monthly Income | Work-Life Balance | Job Satisfaction | Performance Rating | Number of Promotions | Overtime | Distance from Home | Education Level | Marital Status | Number of Dependents |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 22101 | 15182 | 39 | 0 | 12 | 1 | 12931 | 3 | 0 | 1 | 0 | 1 | 99 | 0 | 1 | 1 |
| 17205 | 62820 | 55 | 0 | 25 | 2 | 6623 | 0 | 3 | 3 | 3 | 0 | 41 | 2 | 1 | 1 |
| 49282 | 17777 | 33 | 1 | 1 | 0 | 4264 | 1 | 0 | 2 | 1 | 0 | 54 | 1 | 1 | 0 |
| 26225 | 1824 | 45 | 0 | 35 | 1 | 7434 | 2 | 3 | 3 | 1 | 0 | 22 | 0 | 2 | 1 |
| 8047 | 70265 | 51 | 0 | 11 | 3 | 6240 | 3 | 0 | 0 | 0 | 1 | 84 | 0 | 1 | 2 |

```python
In [32]: column = ['Attrition',
                   'Age','Gender',
                   'Years at Company',
                   'Job Role','Marital Status',
```

---

jupyter  Untitled1 Last Checkpoint: Last Monday at 15:55 (autosaved)

Logout

File   Edit   View   Insert   Cell   Kernel   Widgets   Help

Not Trusted | conda_python3 ○

```python
In [32]: column = ['Attrition',
                   'Age','Gender',
                   'Years at Company',
                   'Job Role','Marital Status',
                   'Education Level',
                   'Job Level',
                   'Number of Dependents',
                   'Monthly Income',
                   'Work-Life Balance',
                   'Job Satisfaction','Overtime',
                   'Distance from Home','Company Size',
                   'Company Tenure','Remote Work',
                   'Leadership Opportunities',
                   'Innovation Opportunities',
                   'Company Reputation',
                   'Employee Recognition',
                  ]
```

```python
In [33]: trainDF= trainDF[column]
         trainDF.head(10)
```

Out[33]:

| | Attrition | Age | Gender | Years at Company | Job Role | Marital Status | Education Level | Job Level | Number of Dependents | Monthly Income | Work-Life Balance | Job Satisfaction | Overtime | Distance from Home | Company Size | Company Tenure |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 48722 | 0 | 29 | 1 | 8 | 4 | 2 | 3 | 0 | 4 | 7870 | 1 | 3 | 0 | 96 | 2 | 37 |
| 1768 | 0 | 26 | 0 | 8 | 0 | 2 | 2 | 0 | 0 | 3685 | 2 | 0 | 0 | 57 | 1 | 65 |
| 1516 | 1 | 39 | 0 | 25 | 3 | 0 | 3 | 0 | 1 | 5980 | 0 | 2 | 0 | 31 | 1 | 67 |
| 23326 | 0 | 28 | 1 | 14 | 1 | 2 | 1 | 0 | 1 | 7857 | 2 | 2 | 1 | 54 | 1 | 50 |
| 7246 | 0 | 55 | 1 | 6 | 2 | 0 | 2 | 0 | 1 | 7678 | 3 | 3 | 0 | 55 | 2 | 37 |
| 3645 | 0 | 25 | 1 | 13 | 3 | 1 | 1 | 0 | 2 | 5385 | 0 | 3 | 0 | 69 | 1 | 90 |
| 11086 | 1 | 38 | 0 | 27 | 2 | 1 | 1 | 0 | 4 | 7074 | 0 | 0 | 1 | 83 | 1 | 86 |
| 33008 | 1 | 32 | 1 | 3 | 2 | 1 | 3 | 0 | 4 | 10464 | 2 | 3 | 0 | 55 | 2 | 69 |

Jupyter **Untitled1** Last Checkpoint: Last Monday at 15:55 (autosaved)

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

Not Trusted | conda_python3

```
In [26]: from sklearn.model_selection import train_test_split
```

```
In [27]: xTrain, xTest, yTrain, yTest = train_test_split(x,y, test_size = 0.2)
```

```
In [28]: xTrain.head(5)
```

Out[28]:

| | Employee ID | Age | Gender | Years at Company | Job Role | Monthly Income | Work-Life Balance | Job Satisfaction | Performance Rating | Number of Promotions | Overtime | Distance from Home | Education Level | Marital Status | Number of Dependents |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 48722 | 73255 | 29 | 1 | 8 | 4 | 7870 | 1 | 3 | 0 | 3 | 0 | 96 | 3 | 2 | 4 |
| 1768 | 29901 | 26 | 0 | 8 | 0 | 3685 | 2 | 0 | 0 | 0 | 0 | 57 | 2 | 2 | 0 |
| 1516 | 4409 | 39 | 0 | 25 | 3 | 5980 | 0 | 2 | 0 | 1 | 0 | 31 | 3 | 0 | 1 |
| 23326 | 53341 | 28 | 1 | 14 | 1 | 7857 | 2 | 2 | 0 | 0 | 1 | 54 | 1 | 2 | 1 |
| 7246 | 56311 | 55 | 1 | 6 | 2 | 7678 | 3 | 3 | 0 | 3 | 0 | 55 | 2 | 0 | 1 |

```
In [29]: yTrain.head(5)
```

Out[29]:
```
48722    0
1768     0
1516     1
23326    0
7246     0
Name: Attrition, dtype: int64
```

```
In [30]: trainDF=xTrain.join(yTrain)
         trainDF.head(5)
```

Out[30]:

| | Employee ID | Age | Gender | Years at Company | Job Role | Monthly Income | Work-Life Balance | Job Satisfaction | Performance Rating | Number of Promotions | Overtime | Distance from Home | Education Level | Marital Status | Number of Dependents |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 48722 | 73255 | 29 | 1 | 8 | 4 | 7870 | 1 | 3 | 0 | 3 | 0 | 96 | 3 | 2 | 4 |

---

Jupyter **Untitled1** Last Checkpoint: Last Monday at 15:55 (autosaved)

File  Edit  View  Insert  Cell  Kernel  Widgets  Help

Not Trusted | conda_python3

```
In [34]: testDF = testDF [column[1:]]
         testDF.head()
```

Out[34]:

| | Age | Gender | Years at Company | Job Role | Marital Status | Education Level | Job Level | Number of Dependents | Monthly Income | Work-Life Balance | Job Satisfaction | Overtime | Distance from Home | Company Size | Company Tenure | Remote Work |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 22101 | 39 | 0 | 12 | 1 | 1 | 0 | 0 | 1 | 12931 | 3 | 0 | 1 | 99 | 2 | 30 | 0 |
| 17205 | 55 | 0 | 25 | 2 | 1 | 2 | 1 | 1 | 6623 | 0 | 3 | 0 | 41 | 1 | 102 | 0 |
| 49282 | 33 | 1 | 1 | 0 | 1 | 1 | 2 | 0 | 4264 | 1 | 0 | 0 | 54 | 1 | 13 | 1 |
| 26225 | 45 | 0 | 35 | 1 | 2 | 0 | 1 | 1 | 7434 | 2 | 3 | 0 | 22 | 1 | 92 | 0 |
| 8047 | 51 | 0 | 11 | 3 | 1 | 0 | 1 | 2 | 6240 | 3 | 0 | 1 | 84 | 0 | 44 | 1 |

```
In [35]: trainDF.to_csv('traineddataattritions.csv',index=False, index_label='Row',header=False, columns=column)
```

```
In [36]: testDF.head()
```

Out[36]:

| | Age | Gender | Years at Company | Job Role | Marital Status | Education Level | Job Level | Number of Dependents | Monthly Income | Work-Life Balance | Job Satisfaction | Overtime | Distance from Home | Company Size | Company Tenure | Remote Work |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 22101 | 39 | 0 | 12 | 1 | 1 | 0 | 0 | 1 | 12931 | 3 | 0 | 1 | 99 | 2 | 30 | 0 |
| 17205 | 55 | 0 | 25 | 2 | 1 | 2 | 1 | 1 | 6623 | 0 | 3 | 0 | 41 | 1 | 102 | 0 |
| 49282 | 33 | 1 | 1 | 0 | 1 | 1 | 2 | 0 | 4264 | 1 | 0 | 0 | 54 | 1 | 13 | 1 |
| 26225 | 45 | 0 | 35 | 1 | 2 | 0 | 1 | 1 | 7434 | 2 | 3 | 0 | 22 | 1 | 92 | 0 |
| 8047 | 51 | 0 | 11 | 3 | 1 | 0 | 1 | 2 | 6240 | 3 | 0 | 1 | 84 | 0 | 44 | 1 |

```
In [37]: trainDF.to_csv('testddataattritions.csv',index=False, index_label='Row',header=False, columns=column)
```

```
In [38]: import boto3 #this package is to integrate with s3 bucket or other cloude service//
```

Creating an s3 storage and moving the trained and test data to the s3 object storage

```
In [37]: trainDF.to_csv('testddataattritions.csv',index=False, index_label='Row',header=False, columns=column)

In [38]: import boto3 #this package is to integrate with s3 bucket or other clouce service//
         import re #this package is to folow a strict pattern to save your work/regular expresession//

In [39]: bucketNM = 'diabeticsinstancebucket'
         TrainFile = r'attritiondata/traineddataattritions/traineddataattritions.csv'
         TestFile = r'attritiondata/testddataattritions/testddataattritions.csv'
         ValFile = r'attritiondata/Val/Val.csv'
         ModelFolder = r'attritiondata/model/'

In [40]: s3ModelOutput = r's3://{0}/{1}'.format(bucketNM,ModelFolder)
         s3Train = r's3://{0}/{1}'.format(bucketNM,TrainFile)
         s3Test = r's3://{0}/{1}'.format(bucketNM,TestFile)
         s3Val = r's3://{0}/{1}'.format(bucketNM,ValFile)

In [41]: s3ModelOutput

Out[41]: 's3://diabeticsinstancebucket/attritiondata/model/'

In [42]: with open('traineddataattritions.csv','rb') as f:
             boto3.Session().resource('s3').Bucket(bucketNM).Object(TrainFile).upload_fileobj(f)

In [43]: import sagemaker
         from sagemaker import get_execution_role

         sagemaker.config INFO - Not applying SDK defaults from location: /etc/xdg/sagemaker/config.yaml
         sagemaker.config INFO - Not applying SDK defaults from location: /home/ec2-user/.config/sagemaker/config.yaml

In [44]: sagemakerSess=sagemaker.Session()
         role=get_execution_role()
```

Setting up the hyperparameter for tuning/optimization



```
In [44]: sagemakerSess=sagemaker.Session()
         role=get_execution_role()

In [45]: sagemakerSess.boto_region_name

Out[45]: 'us-east-1'

In [46]: ECRdockercontainer=sagemaker.amazon.amazon_estimator.get_image_uri(sagemakerSess.boto_region_name,'linear-learner','latest')

         The method get_image_uri has been renamed in sagemaker>=2.
         See: https://sagemaker.readthedocs.io/en/stable/v2.html for details.
         Defaulting to the only supported framework/algorithm version: 1. Ignoring framework/algorithm version: latest.

In [47]: LogisticModel=sagemaker.estimator.Estimator(image_uri=ECRdockercontainer,
                                         role=role,
                                         train_instance_count=1,
                                         train_instance_type='ml.m4.xlarge',
                                         output_path=s3ModelOutput,
                                         sagemaker_session=sagemakerSess,
                                         base_job_name = 'Logistic-Demo-v1'
                                         )

         train_instance_count has been renamed in sagemaker>=2.
         See: https://sagemaker.readthedocs.io/en/stable/v2.html for details.
         train_instance_type has been renamed in sagemaker>=2.
         See: https://sagemaker.readthedocs.io/en/stable/v2.html for details.

In [48]: LogisticModel.set_hyperparameters(predictor_type='binary_classifier', mini_batch_size=100)

In [49]: LogisticModel.hyperparameters()

Out[49]: {'predictor_type': 'binary_classifier', 'mini_batch_size': 100}

In [50]: trainConfig=sagemaker.session.s3_input(s3_data=s3Train,content_type='text/csv')
```
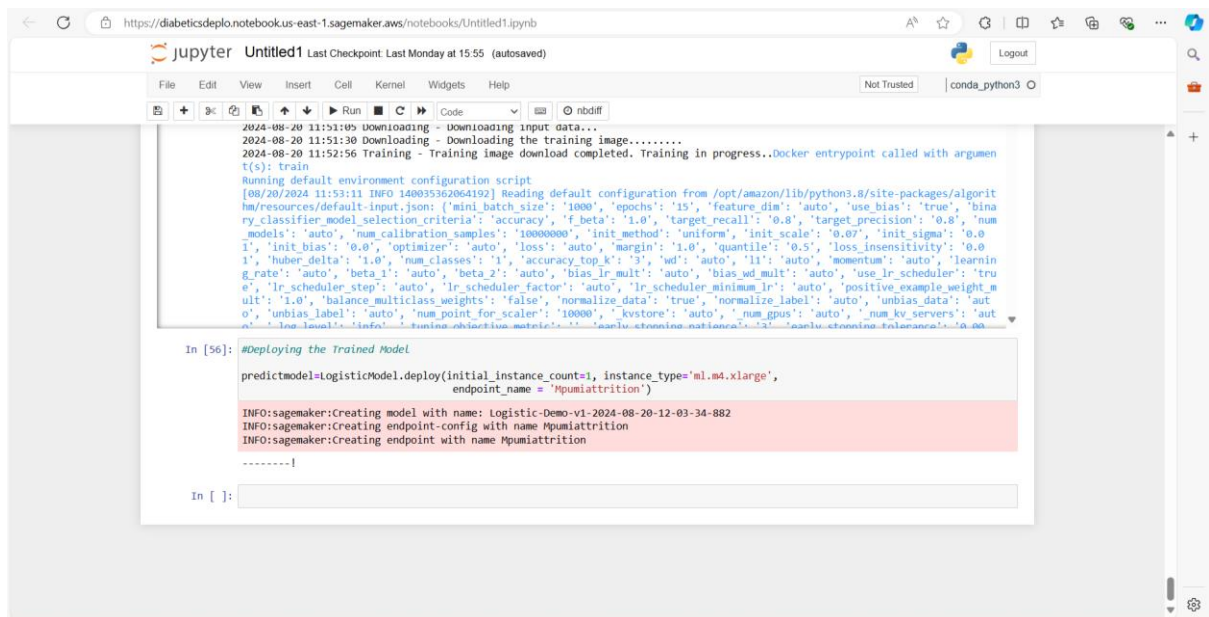
Building and Training the model on sagemaker

Deploying the model after training



The deployed model on AWS sagemaker
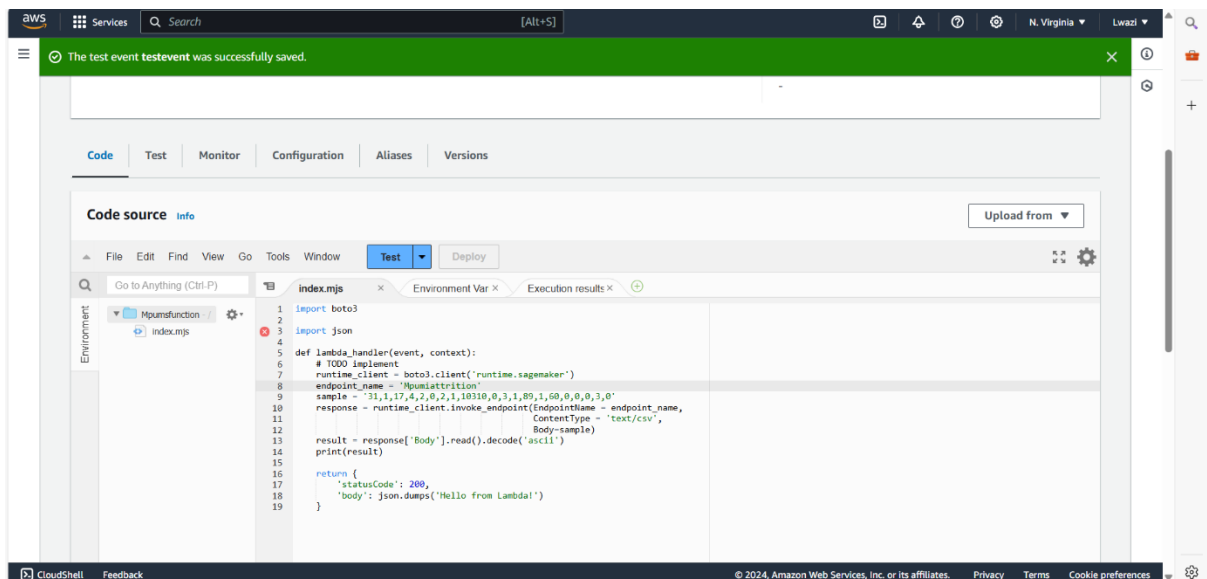
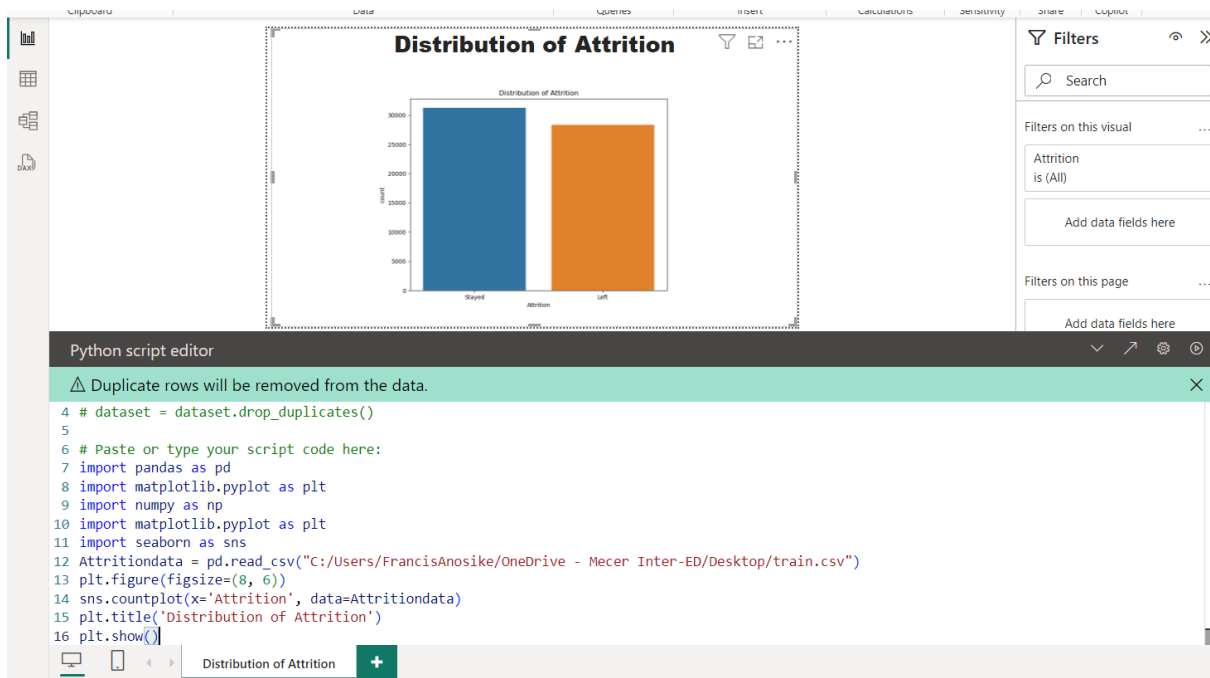The metrics and monitoring of resources

The end point link:

https://runtime.sagemaker.us-east-1.amazonaws.com/endpoints/Mpumiattrition/invocations
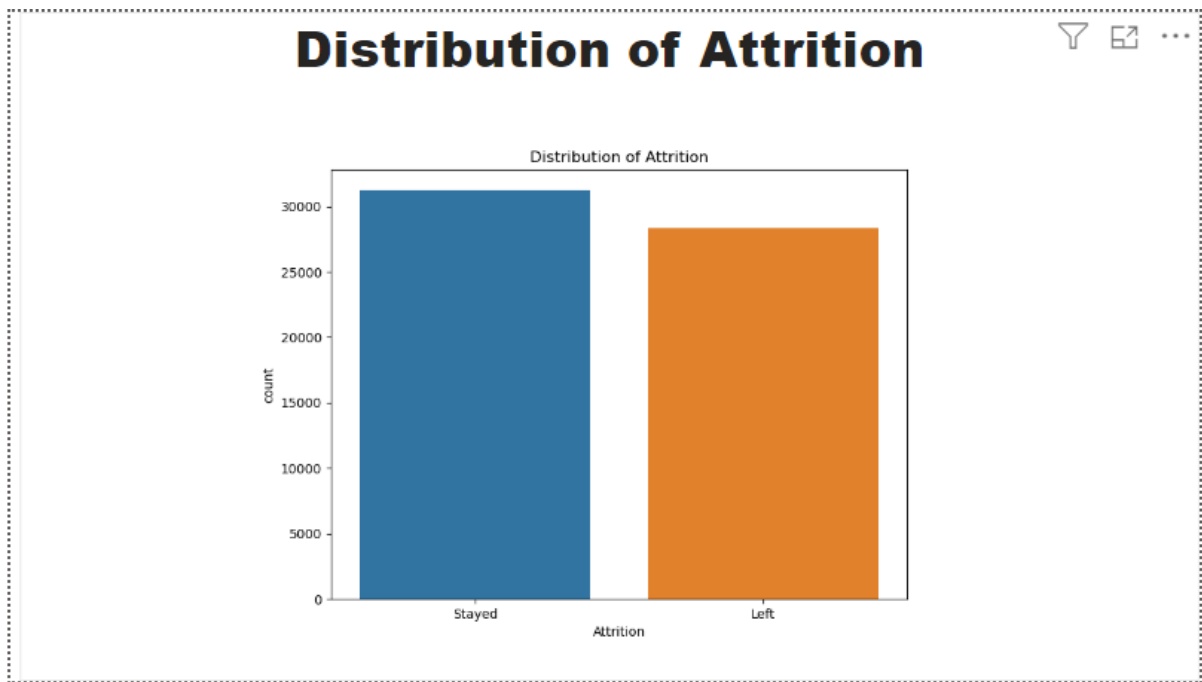
Testing the model on the lambda service on AWS:

Visualizing the dataset using python scripts embedded in PowerBI

The employee attrition visual

# Distribution of Attrition



Distribution of Attrition

Employee Attrition by job role



```python
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
Attritiondata = pd.read_csv("C:/Users/FrancisAnosike/OneDrive - Mecer Inter-ED/Desktop/train.csv")
# Plot Attrition by Job Role
plt.figure(figsize=(12, 8))
sns.countplot(x='Job Role', hue='Attrition', data=Attritiondata)
plt.title('Attrition by Job Role')
plt.xticks(rotation=45)
plt.show()
```

# Attrition by Job Role



Employee attrition by job satisfaction.



```python
# Paste or type your script code here:
import pandas as pd
import matplotlib.pyplot as plt
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
Attritiondata = pd.read_csv("C:/Users/FrancisAnosike/OneDrive - Mecer Inter-ED/Desktop/train.csv")

# Plot Attrition by Job Satisfaction
plt.figure(figsize=(8, 6))
sns.countplot(x='Job Satisfaction', hue='Attrition', data=Attritiondata)
plt.title('Attrition by Job Satisfaction')
plt.show()
```
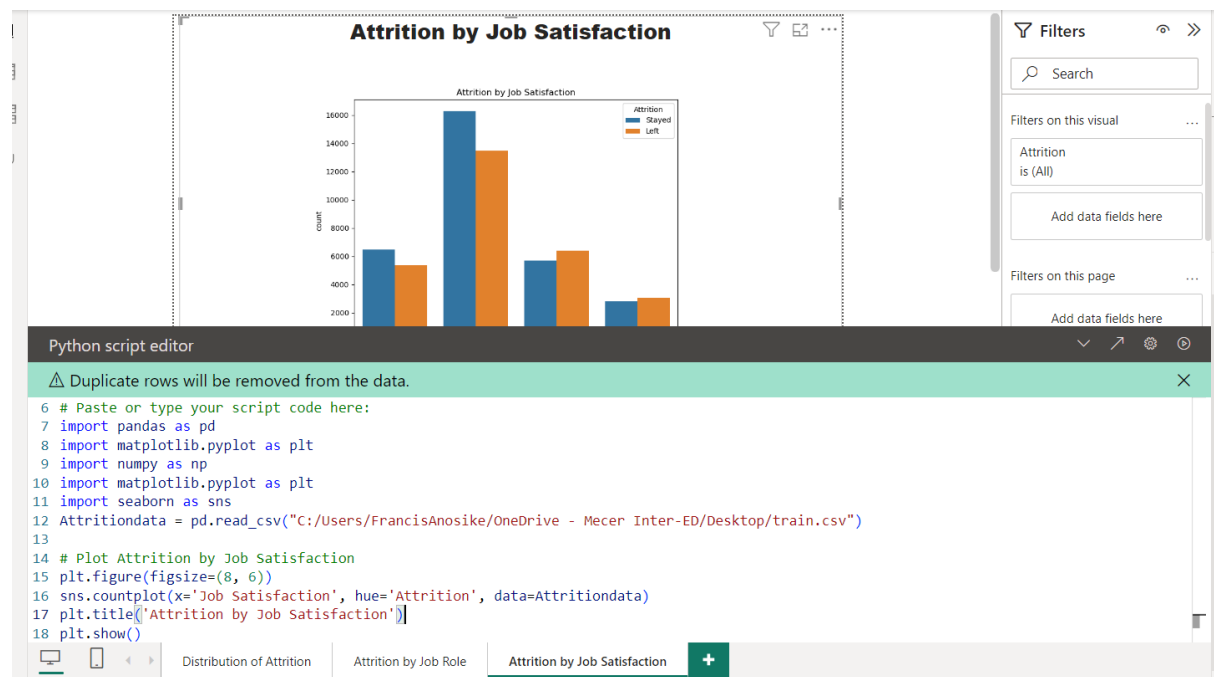
Attrition by Job Satisfaction

Employee attrition by work-life balance.



```
6  # Paste or type your script code here:
7  import pandas as pd
8  import matplotlib.pyplot as plt
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 Attritiondata = pd.read_csv("C:/Users/FrancisAnosike/OneDrive - Mecer Inter-ED/Desktop/train.csv")
13
14 plt.figure(figsize=(8, 6))
15 sns.countplot(x='Work-Life Balance', hue='Attrition', data=Attritiondata)
16 plt.title('Attrition by Work-Life Balance')
17 plt.show()
18
```
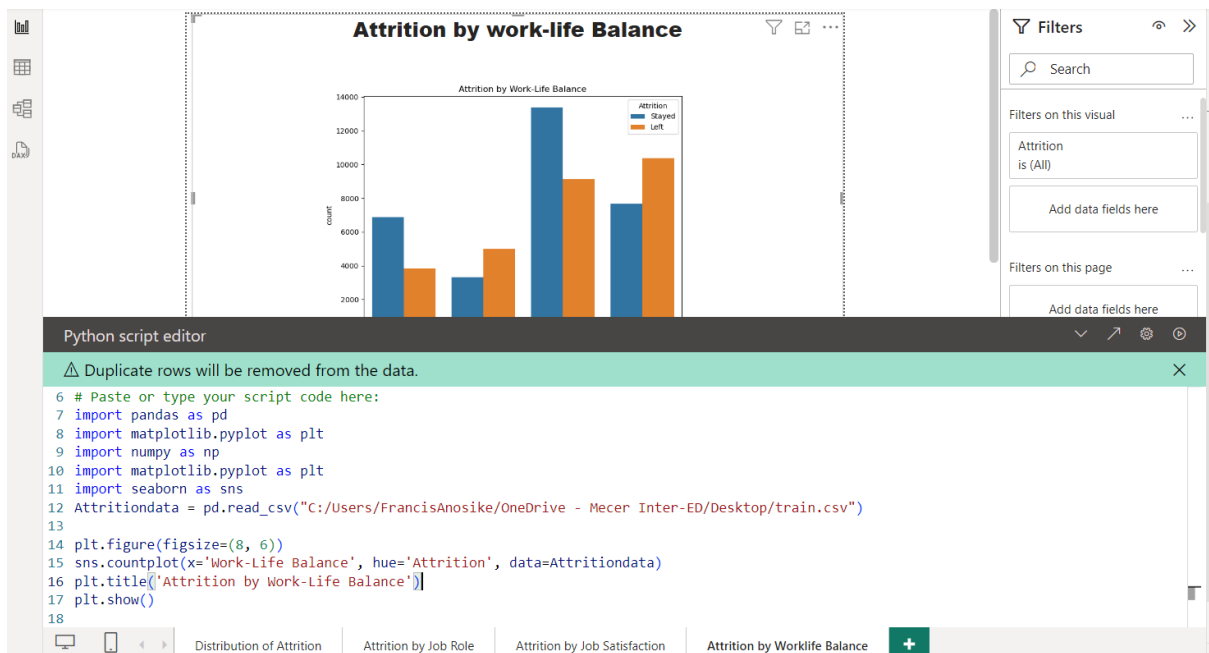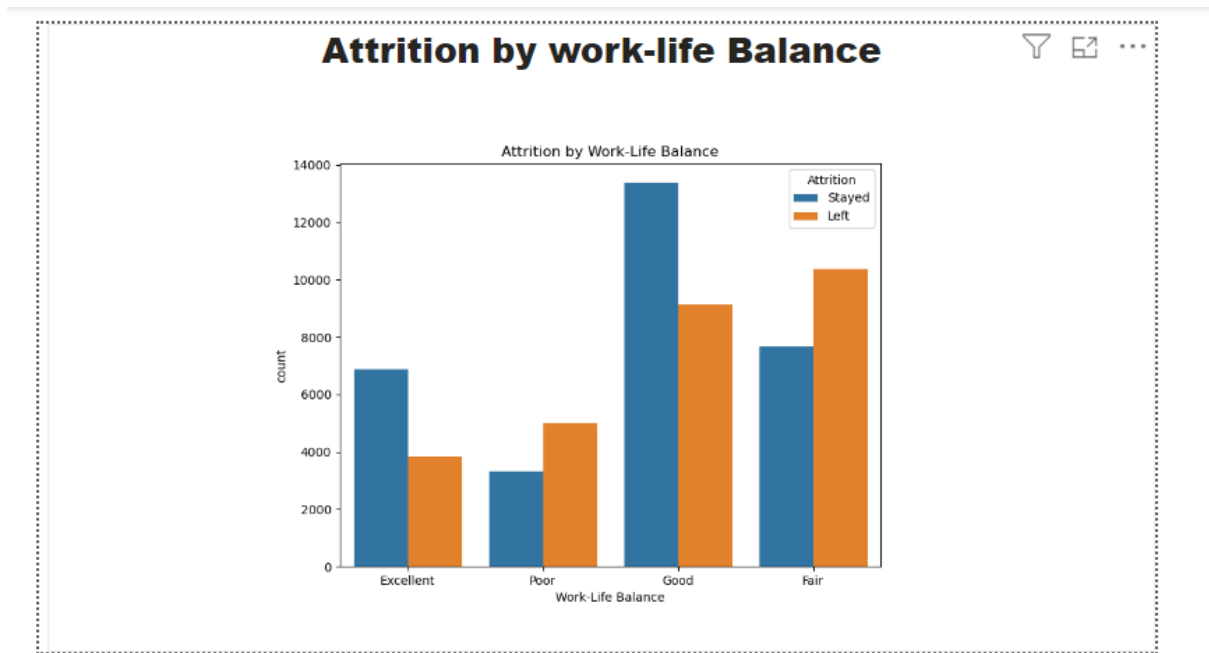
Employee attrition by marital status.



```python
4  # dataset = dataset.drop_duplicates()
5
6  # Paste or type your script code here:
7  import pandas as pd
8  import matplotlib.pyplot as plt
9  import numpy as np
10 import matplotlib.pyplot as plt
11 import seaborn as sns
12 Attritiondata = pd.read_csv("C:/Users/FrancisAnosike/OneDrive - Mecer Inter-ED/Desktop/train.csv")
13 plt.figure(figsize=(8, 6))
14 sns.countplot(x='Marital Status', hue='Attrition', data=Attritiondata)
15 plt.title('Attrition by Marital Status')
16 plt.show()
```
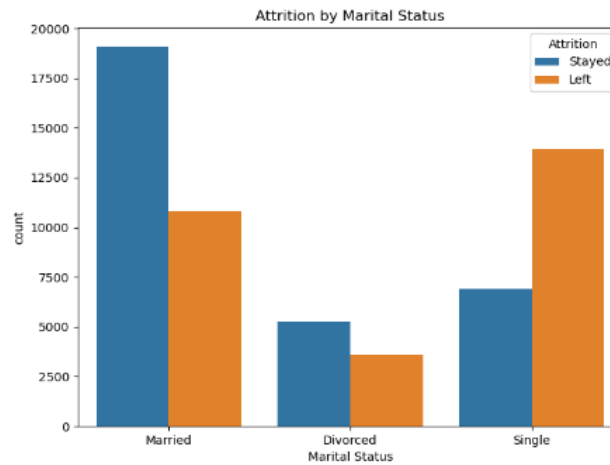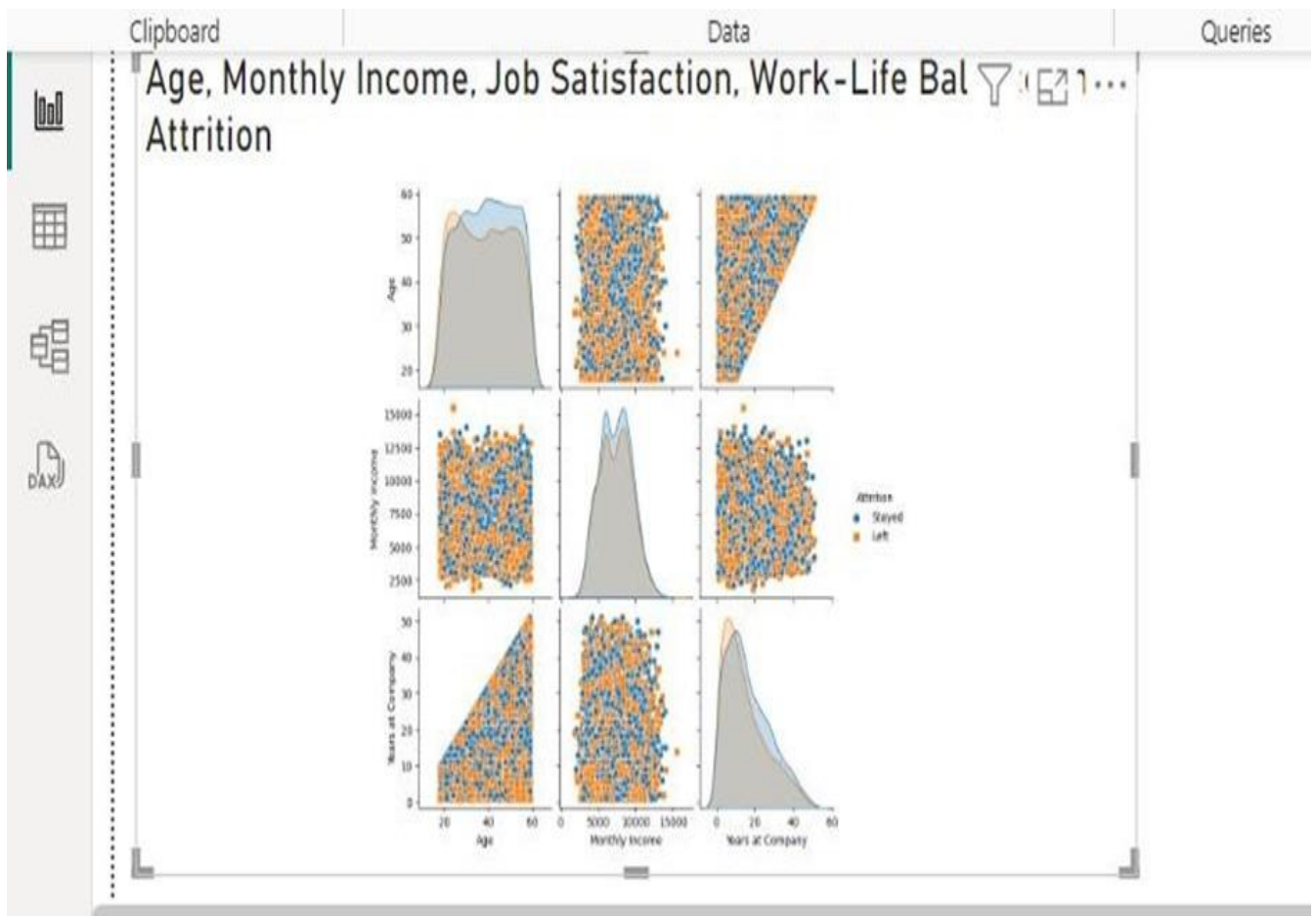
# Attrition by Marital Status



Attrition by Marital Status

```python
1  import statsmodels.api as sm
2  import numpy as np
3  import pandas as pd
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  from sklearn.preprocessing import LabelEncoder
7
8  df = pd.read_csv(r"C:\Users\Tesserai\Desktop\Train Project\train.csv")
9
10 # Pair plot for selected features
11 selected_features = ['Age', 'Monthly Income', 'Years at Company', 'Job Satisfaction', 'Work-Life Balance', 'Attrition']
12 sns.pairplot(df[selected_features], hue='Attrition', diag_kind='kde', markers=["o", "s"])
13 plt.show()
```

Age, Monthly Income, Job Satisfaction, Work-Life Bal ▽ ⟨[⟩ 1 ...
Attrition

## Monitoring and Maintenance:

Once deployed, the model's performance needs to be monitored regularly to ensure that it continues to provide accurate predictions over time. This may involve updating the model with new data or retraining it periodically to maintain its accuracy.

## Deployment Environment used for the deployment of the model:

Frameworks and Libraries:

NumPy: For numerical computations and array manipulations.

Pandas: For data manipulation and analysis, particularly useful for handling datasets like the Iris dataset.

Scikit-learn: For machine learning algorithms and model training. It includes logistic regression and utilities for model evaluation.

Development Tools:

IDEs: PowerBI, VS Code, or Jupyter notebook for coding and testing.

Version Control: Git for managing code versions.

Package Management: pip or conda for managing Python packages and dependencies.

Security Considerations:

Access Control: Restrict access to the deployed model and its endpoints. Implement role-based access control (RBAC) to ensure only authorized personnel can interact with the model.

Encryption: Use encryption mechanisms (e.g., HTTPS/TLS) to secure data transmission between clients and the deployed model, preventing eavesdropping and data tampering.

Input Validation: Validate input data to prevent injection attacks and ensure that only expected and sanitized data is processed by the model.