# Libraries ¶

```
In [2]:   1  # Suppressing all warnings
          2  import warnings
          3  warnings.filterwarnings("ignore")
          4
          5  # Multiprocessing
          6  import multiprocessing
          7
          8  # Threading
          9  import threading
         10
         11  # Time
         12  import time
         13
         14  # Data handling
         15  import pandas as pd
         16  import numpy as np
         17
         18  # Plotting
         19  import matplotlib.pyplot as plt
         20  import seaborn as sns
         21
         22  # Oversampling
         23  from imblearn.over_sampling import ADASYN
         24
         25  # Metrics
         26  from sklearn import metrics
         27
         28  # Preprocessing
         29  from imblearn.over_sampling import ADASYN
         30  from sklearn.model_selection import train_test_split
         31  from sklearn.decomposition import PCA
         32
         33  # Models
         34  from sklearn.pipeline import Pipeline
         35  from sklearn.svm import SVC
         36  from sklearn.neighbors import KNeighborsClassifier
         37  from sklearn.tree import DecisionTreeClassifier
         38  from sklearn.linear_model import LogisticRegression
         39  from xgboost import XGBClassifier
         40  from sklearn.ensemble import AdaBoostClassifier
         41  # from catboost import CatBoostClassifier
         42
         43  # Evaluation metrics
         44
         45  # Cross validation
         46  from sklearn.model_selection import RandomizedSearchCV, GridSearchCV, train_test_split
         47  from skopt import BayesSearchCV
         48  from sklearn.metrics import classification_report
```

```
In [2]:  1  def multi_time():
         2      print("Starting thread...")
         3      start = time.time()
         4
         5      while True:
         6          elapsed = int(time.time() - start)
         7          print(f"Time elapsed: {elapsed} seconds", end = "\r")
         8          time.sleep(1)
```

## Loading CSV file onto Pandas DataFrame

```
In [74]:  1  df = pd.read_csv("player_stats.csv", na_values="None")
          2  df = df[["X", "Y", "xG", "h_a", "situation", "shotType", "lastAction", "result"]]
          3  # df = df.loc[:50000, :].copy()
```

## Check for missing data and percentage

```
In [75]:  1  print("Percentage of data missing from each column\n")
          2  for col in df.columns:
          3      if len(df[df[col].isnull() == True]):
          4          print(f"{col}: {len(df[df[col].isnull() == True])/len(df)*100}%",
          5                  end = "")
          6          print(f" or {len(df[df[col].isnull() == True])} rows")
          7      else:
          8          print(f"{col}: {0}%")
```

```
Percentage of data missing from each column

X: 0%
Y: 0%
xG: 0%
h_a: 0%
situation: 0%
shotType: 0%
lastAction: 9.228053426520521% or 10246 rows
result: 0%
```

## Dropping rows with missing values

```
In [76]:  1  df = df[df["lastAction"].isnull() == False].reset_index(drop =True).copy()
          2  df = df.drop(df[df["result"]=="OwnGoal"].index)
          3  df = df[df.situation != 'Penalty']
```

## Encoding categorical columns

- h_a: One-Hot Encoder
- situation: One-Hot Encoder
- shotType: One-Hot Encoder
- lastAction: One-Hot Encoder
- result: One-Hot Encoder

```
In [77]:  1  print(f"\x1b[32mUnique values in:\x1b[0m {df['h_a'].value_counts()}")
          2  print(f"\x1b[32mUnique values in:\x1b[0m {df['situation'].value_counts()}")
          3  print(f"\x1b[32mUnique values in:\x1b[0m {df['shotType'].value_counts()}")
          4  print(f"\x1b[32mUnique values in:\x1b[0m {df['lastAction'].value_counts()}")
          5  print(f"\x1b[32mUnique values in:\x1b[0m {df['result'].value_counts()}")
```

Unique values in: h     54108
a     44615
Name: h_a, dtype: int64
Unique values in: OpenPlay        82980
FromCorner       8514
SetPiece         3918
DirectFreekick    3311
Name: situation, dtype: int64
Unique values in: RightFoot        48933
LeftFoot        32717
Head            16725
OtherBodyPart     348
Name: shotType, dtype: int64
Unique values in: Pass             41618
Cross           15324
Aerial           7000
Chipped          6630
TakeOn           6616
Rebound          4705
Throughball       3515
Standard         3311
BallRecovery      3081
HeadPass         2644
BallTouch        1919
LayOff            707
Dispossessed      547
Tackle            221
CornerAwarded     194
BlockedPass       171
Foul              159
Goal               88
Interception       65
End                64
OffsidePass        23
GoodSkill          23
Challenge          18
Card               18
SubstitutionOn     16
Clearance          13
OffsideProvoked     9
KeeperPickup        7
Save                7
FormationChange     5
Start               3
ChanceMissed        1
ShieldBallOpp       1
Name: lastAction, dtype: int64
Unique values in: MissedShots     34019
SavedShot       25463
BlockedShot     24329
Goal            12795
ShotOnPost       2117
Name: result, dtype: int64

```python
In [78]:    1  remapping = {"Goal": 1,
            2              "MissedShots": 0,
            3              "SavedShot": 0,
            4              "BlockedShot": 0,
            5              "ShotOnPost": 0,
            6              }
            7
            8  lastAction = {"Cross": ["Aerial", "Cross", "Chipped"] ,
            9              "Pass":["Pass", "Throughball", "HeadPass", "BallTouch", "TakeOn"],
            10             "Dispossessed":["Dispossessed", "Tackle", "BlockedPass", "Interception",
            11                             "Challenge", "Clearance", "BallRecovery"],
            12             "Other":["Standard", "Rebound", "LayOff", "CornerAwarded",
            13                     "Foul", "Goal", "End", "OffsidePass",
            14             "GoodSkill", "Card", "SubstitutionOn", "OffsideProvoked",
            15                     "Save", "KeeperPickup",
            16             "FormationChange", "Start", "ChanceMissed" , "ShieldBallOpp"]}
            17
            18 df["result"] = df["result"].map(remapping)
            19 for k, v in lastAction.items():
            20     for i in v:
            21         df["lastAction"].replace(to_replace = i, value = k, inplace = True)
            22
            23 df_encoded = pd.get_dummies(df,
            24                             drop_first = True)
            25 df_encoded = df_encoded.loc[:50000, :]
```

```python
In [79]:    1  X = df_encoded.drop(["result"], axis = 1).copy()
            2  Y = df_encoded["result"].copy()
```

```python
In [80]:    1  X_train, X_test, y_train, y_test_1 = train_test_split(X, Y,
            2                                              train_size = 0.8,
            3                                              stratify = Y)
```

## Step 2: Grid Search

```python
In [81]:  1  pipe = Pipeline([("classifier", LogisticRegression())])
          2
          3  param_grid = [
          4      {
          5          'classifier': [LogisticRegression(random_state = 42)],
          6          'classifier__solver': ['saga', 'lbfgs'],
          7          'classifier__C': [0.1, 1],
          8      },
          9      {
         10          'classifier': [SVC(random_state = 42)],
         11          'classifier__C':[0.1, 1],
         12          'classifier__degree': [1, 2],
         13          'classifier__gamma': [0.1, 1],
         14      },
         15      {
         16          'classifier': [XGBClassifier(random_state = 42)],
         17          'classifier__n_estimators': [10, 1000],
         18          'classifier__max_depth': [10, 50],
         19      },
         20      {
         21          'classifier': [AdaBoostClassifier()],
         22          'classifier__base_estimator': [DecisionTreeClassifier(random_state = 42)],
         23          'classifier__learning_rate': [0.01, 0.1, 1],
         24          'classifier__base_estimator__criterion': ['gini', 'entropy'],
         25          'classifier__base_estimator__max_depth': [1, 5, 9],
         26      }
         27  ]
         28
         29  bs = GridSearchCV(pipe, param_grid, cv = 5, verbose = 3)
         30  bs.fit(X_train, y_train)
```

```
[CV 5/5] END classifier=LogisticRegression(random_state=42), classifier__C=1, classifie
r__solver=saga;, score=0.876 total time=   0.1s
[CV 1/5] END classifier=LogisticRegression(random_state=42), classifier__C=1, classifie
r__solver=lbfgs;, score=0.879 total time=   0.1s
[CV 2/5] END classifier=LogisticRegression(random_state=42), classifier__C=1, classifie
r__solver=lbfgs;, score=0.879 total time=   0.1s
[CV 3/5] END classifier=LogisticRegression(random_state=42), classifier__C=1, classifie
r__solver=lbfgs;, score=0.875 total time=   0.1s
[CV 4/5] END classifier=LogisticRegression(random_state=42), classifier__C=1, classifie
r__solver=lbfgs;, score=0.877 total time=   0.1s
[CV 5/5] END classifier=LogisticRegression(random_state=42), classifier__C=1, classifie
r__solver=lbfgs;, score=0.876 total time=   0.1s
[CV 1/5] END classifier=SVC(random_state=42), classifier__C=0.1, classifier__degree=1,
classifier__gamma=0.1;, score=0.867 total time=   37.1s
[CV 2/5] END classifier=SVC(random_state=42), classifier__C=0.1, classifier__degree=1,
classifier__gamma=0.1;, score=0.867 total time=   37.5s
[CV 3/5] END classifier=SVC(random_state=42), classifier__C=0.1, classifier__degree=1,
classifier__gamma=0.1;, score=0.867 total time=   38.7s
[CV 4/5] END classifier=SVC(random_state=42), classifier__C=0.1, classifier__degree=1,
classifier__gamma=0.1;, score=0.867 total time=   37.6s
```

```python
In [82]:  1  bs.best_params_
```

```
Out[82]: {'classifier': AdaBoostClassifier(base_estimator=DecisionTreeClassifier(max_depth=1,
                                                                                  random_state=42),
                     learning_rate=1),
          'classifier__base_estimator': DecisionTreeClassifier(max_depth=1, random_state=42),
          'classifier__base_estimator__criterion': 'gini',
          'classifier__base_estimator__max_depth': 1,
          'classifier__learning_rate': 1}
```

```
In [83]:   1  y_preds_1 = bs.predict(X_test)
           2  acc = metrics.accuracy_score(y_test_1, y_preds_1)
```

```
In [84]:   1  cr_1 = classification_report(y_test_1, y_preds_1)
           2  print(classification_report(y_test_1, y_preds_1))
```

```
              precision    recall  f1-score   support

           0       0.89      0.98      0.93      8502
           1       0.60      0.20      0.30      1307

    accuracy                           0.88      9809
   macro avg       0.74      0.59      0.62      9809
weighted avg       0.85      0.88      0.85      9809
```

# Recall isn't good for minority class

This can attributed to the heavy class imbalance in the dataset. We can try to balance this by first filtering data and then using ADASYN

## Higher cutoff for number of shots per player

In [85]:
```python
def preprocessing(df, num_rows = 50000):

    print("\x1b[32mStarting preprocessing...\x1b[31m")

    df = df[["X", "Y", "xG", "h_a", "situation", "shotType", "lastAction", "result"]]
    df = df[df["lastAction"].isnull() == False].reset_index(drop =True).copy()
    df = df.drop(df[df["result"]=="OwnGoal"].index)
    df = df[df.situation != 'Penalty']
    remapping = {"Goal": 1,
                "MissedShots": 0,
                "SavedShot": 0,
                "BlockedShot": 0,
                "ShotOnPost": 0,
                }

    lastAction = {"Cross": ["Aerial", "Cross", "Chipped"] ,
                    "Pass":["Pass", "Throughball", "HeadPass", "BallTouch", "TakeOn"],
                    "Dispossessed":["Dispossessed", "Tackle", "BlockedPass",
                                        "Interception", "Challenge", "Clearance",
                                        "BallRecovery"],
                    "Other":["Standard", "Rebound", "LayOff", "CornerAwarded",
                                "Foul", "Goal", "End", "OffsidePass",
                            "GoodSkill", "Card", "SubstitutionOn", "OffsideProvoked",
                                    "Save", "KeeperPickup", "FormationChange",
                                    "Start", "ChanceMissed" , "ShieldBallOpp"]}

    print("Remapping target column")
    df["result"] = df["result"].map(remapping)

    print("Remapping lastAction column")
    for k, v in lastAction.items():
        for i in v:
            df["lastAction"].replace(to_replace = i, value = k, inplace = True)

    print("One hot encoding remaining columns")
    df_encoded = pd.get_dummies(df,
                                drop_first = True)
    if num_rows != None:
        print(f"Using only {num_rows} rows")
        df_encoded = df_encoded.loc[:num_rows, :]
    else:
        print(f"Using all rows")

    print("Preprocessing done\x1b[0m")
    return df_encoded
```

```python
In [86]:  1  def search(X_train, y_train, verbosity = 1, use_svc = True):
          2
          3      print("\x1b[32m Starting search...\x1b[31m")
          4
          5      print("Building base pipeline")
          6      pipe = Pipeline([("classifier", LogisticRegression())])
          7
          8      print("Building parameter grid")
          9      if use_svc:
         10          print("Using SVC")
         11          param_grid = [
         12      {
         13          'classifier': [LogisticRegression(random_state = 42)],
         14          'classifier__solver': ['saga', 'lbfgs'],
         15          'classifier__C': [0.1, 1],
         16      },
         17      {
         18          'classifier': [SVC(random_state = 42)],
         19          'classifier__C':[0.1, 1],
         20          'classifier__degree': [1, 2],
         21          'classifier__gamma': [0.1, 1],
         22      },
         23      {
         24          'classifier': [XGBClassifier(random_state = 42)],
         25          'classifier__n_estimators': [10, 1000],
         26          'classifier__max_depth': [10, 50],
         27      },
         28      {
         29          'classifier': [AdaBoostClassifier()],
         30          'classifier__base_estimator': [DecisionTreeClassifier(random_state = 42)],
         31          'classifier__learning_rate': [0.01, 0.1, 1],
         32          'classifier__base_estimator__criterion': ['gini', 'entropy'],
         33          'classifier__base_estimator__max_depth': [1, 5, 9],
         34      }
         35  ]
         36      else:
         37          print("Replacing SVC with KNN")
         38          param_grid = [
         39              {
         40                  'classifier': [LogisticRegression(random_state = 42)],
         41                  'classifier__solver': ['saga', 'lbfgs'],
         42                  'classifier__C': [0.1, 1],
         43              },
         44              {
         45                  'classifier': [KNeighborsClassifier()],
         46                  'classifier__n_neighbors':[3, 4, 5],
         47                  'classifier__p': [1, 2],
         48              },
         49              {
         50                  'classifier': [XGBClassifier(random_state = 42)],
         51                  'classifier__n_estimators': [10, 1000],
         52                  'classifier__max_depth': [10, 50],
         53              },
         54              {
         55                  'classifier': [AdaBoostClassifier()],
         56                  'classifier__base_estimator': [DecisionTreeClassifier(random_state = 42
         57                  'classifier__learning_rate': [0.01, 0.1, 1],
         58                  'classifier__base_estimator__criterion': ['gini', 'entropy'],
         59                  'classifier__base_estimator__max_depth': [1, 5, 9],
         60              }
         61          ]
         62      print("Creating Grid Search")
         63      bs = GridSearchCV(pipe, param_grid, cv = 5, verbose = verbosity)
```

```python
64     print("Fitting train data onto Grid Search")
65     bs.fit(X_train, y_train)
66
67     print(f"\x1b[32mBest Parameters:\x1b[0m\n {bs.best_params_}")
68
69     return bs
```

In [87]:
```python
1  def test(best_model, X_test, y_test):
2      y_pred = best_model.predict(X_test)
3      acc = metrics.accuracy_score(y_test, y_pred)
4
5      print(f"\x1b[32mClassification report on test accuracy:\x1b[0m\n {classification_re
6
7      return y_pred, classification_report(y_test, y_pred)
```

```
In [88]:   1  df2 = pd.read_csv("player_stats.csv",  na_values="None")
           2
           3  df3 = df2.copy()
           4
           5  print("Before filtering, df shape", df3.shape)
           6
           7  df3_grouped = df3.groupby(["player", "result"]).size()
           8  df3['shot_count'] = df3.groupby('player')['player'].transform('size')
           9  df4 = df3.groupby(["player", 'shot_count', "result"])[['shot_count']].count()
          10  df4 = df4.rename(columns = {'shot_count':'count'})
          11  df5 = df4.sort_values(by = df4.index.names[1], ascending = False)
          12  df5["result"] = df5.index.get_level_values(2)
          13  df5.index = df5.index.droplevel(2)
          14  df5['shot_count'] = df5.index.get_level_values(1)
          15  df5.index= df5.index.droplevel(1)
          16  df5 = df5[df5['shot_count']>40]
          17  df3 = df3[df3["player"].isin(df5.index)]
          18  print("After filtering, df shape", df3.shape)
          19
          20  df_encoded = preprocessing(df3, num_rows = None)
          21  X = df_encoded.drop(["result"], axis = 1).copy()
          22  Y = df_encoded["result"].copy()
          23
          24  X_train, X_test, y_train, y_test_2 = train_test_split(X, Y,
          25                                                train_size = 0.8,
          26                                                stratify = Y)
          27
          28  best_model_2 = search(X_train, y_train, verbosity = 3,
          29                        use_svc = True)
          30
          31  y_preds_2, cr_2 = test(best_model_2, X_test, y_test_2)
          32
          33
```

```
1.6s
[CV 3/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=0.1;, score=0.880 total time=  1
1.7s
[CV 4/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=0.1;, score=0.878 total time=  1
1.8s
[CV 5/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=0.1;, score=0.881 total time=  1
2.0s
[CV 1/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=1;, score=0.878 total time=  11.9
s
[CV 2/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase estimator  max depth=5, classifier  learning rate=1;, score=0.879 total time=  11.6
```

## Metrics didn't increase as much.

Let's try oversampling using ADASYN and replace KNN with SVC

```python
In [29]:   1  df2 = pd.read_csv("player_stats.csv",  na_values="None")
           2
           3  df3 = df2.copy()
           4
           5  print("Before filtering, df shape", df3.shape)
           6
           7  df3_grouped = df3.groupby(["player", "result"]).size()
           8  df3['shot_count'] = df3.groupby('player')['player'].transform('size')
           9  df4 = df3.groupby(["player", 'shot_count', "result"])[['shot_count']].count()
          10  df4 = df4.rename(columns = {'shot_count':'count'})
          11  df5 = df4.sort_values(by = df4.index.names[1], ascending = False)
          12  df5["result"] = df5.index.get_level_values(2)
          13  df5.index = df5.index.droplevel(2)
          14  df5['shot_count'] = df5.index.get_level_values(1)
          15  df5.index= df5.index.droplevel(1)
          16  df5 = df5[df5['shot_count']>40]
          17  df3 = df3[df3["player"].isin(df5.index)]
          18  print("After filtering, df shape", df3.shape)
          19
          20  df_encoded = preprocessing(df3, num_rows = 50000)
          21  df5_g = df_encoded[df_encoded["result"]==1]
          22  df5_ng = df_encoded[df_encoded["result"]==0]
          23
          24  print("\x1b[32mBefore ADASYN\x1b[0m")
          25  print(f"Ratio of goals to not goals: {len(df5_g)/len(df5_ng)}")
          26
          27  print("\x1b[31mOversampling using ADASYN\x1b[0m")
          28  X = df_encoded.drop(["result"], axis = 1).copy()
          29  Y = df_encoded["result"].copy()
          30
          31  X_resampled, y_resampled = ADASYN().fit_resample(X, Y)
          32  print("\x1b[32mAfter ADASYN\x1b[0m")
          33  print(f"Number of samples added by ADASYN {len(y_resampled) - len(Y)}, or {(len(y_resar
          34  print(f"Ratio of goals to not goals: {len(y_resampled[y_resampled == 1])/len(y_resample
          35
          36  X_train, X_test, y_train, y_test_3 = train_test_split(X_resampled, y_resampled,
          37                                                        train_size = 0.8,
          38                                                        stratify = y_resampled)
          39
          40  best_model_3 = search(X_train, y_train, verbosity = 3, use_svc = False)
          41
          42  y_preds_3, cr_3 = test(best_model_3, X_test, y_test_3)
```

```
ase_estimator__max_depth=5, classifier__learning_rate=0.1;, score=0.787 total time=  1
3.6s
[CV 3/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=0.1;, score=0.792 total time=  1
3.7s
[CV 4/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=0.1;, score=0.831 total time=  1
4.3s

[CV 5/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=0.1;, score=0.768 total time=  1
4.5s
[CV 1/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=5, classifier__learning_rate=1;, score=0.817 total time=  14.5
s
[CV 2/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
```

# Entire dataset since it can run faster

```
In [93]:    1  df2 = pd.read_csv("player_stats.csv",  na_values="None")
            2
            3  df3 = df2.copy()
            4
            5  print("Before filtering, df shape", df3.shape)
            6
            7  df3_grouped = df3.groupby(["player", "result"]).size()
            8  df3['shot_count'] = df3.groupby('player')['player'].transform('size')
            9  df4 = df3.groupby(["player", 'shot_count', "result"])[['shot_count']].count()
           10  df4 = df4.rename(columns = {'shot_count':'count'})
           11  df5 = df4.sort_values(by = df4.index.names[1], ascending = False)
           12  df5["result"] = df5.index.get_level_values(2)
           13  df5.index = df5.index.droplevel(2)
           14  df5['shot_count'] = df5.index.get_level_values(1)
           15  df5.index= df5.index.droplevel(1)
           16  df5 = df5[df5['shot_count']>40]
           17  df3 = df3[df3["player"].isin(df5.index)]
           18  print("After filtering, df shape", df3.shape)
           19
           20  df_encoded = preprocessing(df3, num_rows = None)
           21  df5_g = df_encoded[df_encoded["result"]==1]
           22  df5_ng = df_encoded[df_encoded["result"]==0]
           23
           24  print("\x1b[32mBefore ADASYN\x1b[0m")
           25  print(f"Ratio of goals to not goals: {len(df5_g)/len(df5_ng)}")
           26
           27  print("\x1b[31mOversampling using ADASYN\x1b[0m")
           28  X = df_encoded.drop(["result"], axis = 1).copy()
           29  Y = df_encoded["result"].copy()
           30
           31  X_resampled, y_resampled = ADASYN().fit_resample(X, Y)
           32  print("\x1b[32mAfter ADASYN\x1b[0m")
           33  print(f"Number of samples added by ADASYN {len(y_resampled) - len(Y)}, or {(len(y_resar
           34  print(f"Ratio of goals to not goals: {len(y_resampled[y_resampled == 1])/len(y_resampl
           35
           36  X_train, X_test, y_train, y_test_4 = train_test_split(X_resampled, y_resampled,
           37                                                         train_size = 0.8,
           38                                                         stratify = y_resampled)
           39
           40  best_model_4 = search(X_train, y_train, verbosity = 3, use_svc = False)
           41
           42  y_preds_4, cr_4 = test(best_model_4, X_test, y_test_4)
```

```
assifier(random_state=42), classifier__base_estimator__criterion=gini, classifier__base
_estimator__max_depth=9, classifier__learning_rate=1;, score=0.894 total time=  39.2s
[CV 1/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=1, classifier__learning_rate=0.01;, score=0.674 total time=
7.5s
[CV 2/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=1, classifier__learning_rate=0.01;, score=0.680 total time=
7.7s
[CV 3/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=1, classifier__learning_rate=0.01;, score=0.677 total time=
7.7s
[CV 4/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
ase_estimator__max_depth=1, classifier__learning_rate=0.01;, score=0.681 total time=
8.0s
[CV 5/5] END classifier=AdaBoostClassifier(), classifier__base_estimator=DecisionTreeCl
assifier(random_state=42), classifier__base_estimator__criterion=entropy, classifier__b
```

# Classification Reports

```
In [96]:   1  print("\x1b[32mBaseline evaluation metrics:\x1b[0m")
           2  print("\x1b[31mNumber of rows = 50000")
           3  print("List of models:\nLogistic Regression\nSVC\nXgboost\nAdaBoost\x1b[0m")
           4  print(cr_1)
           5
           6  print("\x1b[36mMinority class recall and precision was low.\x1b[0m")
           7  print("\x1b[36mChoosing top 75% percentile of players with respect to number of shots\x
           8
           9  print("\x1b[32mEvaluation metrics after choosing cutoff of number of shots for each pla
          10  print("\x1b[31mNumber of rows = 50000")
          11  print("List of models:\nLogistic Regression\nSVC\nXgboost\nAdaBoost\x1b[0m")
          12  print(cr_2)
          13
          14  print("\x1b[36mMetrics did not change by much\x1b[0m")
          15
          16  print("\x1b[36mOversampling using ADASYN to balance classes and replacing SVC with KNN
          17  print("\x1b[32mEvaluation metrics after oversampling and replacing SVC with KNN:\x1b[0m
          18  print("\x1b[31mNumber of rows = 50000")
          19  print("List of models:\nLogistic Regression\nKNN\nXgboost\nAdaBoost\x1b[0m")
          20  print(cr_3)
          21
          22  print("\x1b[36mMetrics increased by a significant amount with respect to minority class
          23  print("\x1b[36mRunning gridsearch again with entire dataset\x1b[0m")
          24  print("\x1b[31mNumber of rows = entire dataset")
          25  print("List of models:\nLogistic Regression\nKNN\nXgboost\nAdaBoost\x1b[0m")
          26  print(cr_4)
```

Baseline evaluation metrics:
Number of rows = 50000
List of models:
Logistic Regression
SVC
Xgboost
AdaBoost

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.89      | 0.98   | 0.93     | 8502    |
| 1            | 0.60      | 0.20   | 0.30     | 1307    |
|              |           |        |          |         |
| accuracy     |           |        | 0.88     | 9809    |
| macro avg    | 0.74      | 0.59   | 0.62     | 9809    |
| weighted avg | 0.85      | 0.88   | 0.85     | 9809    |

Minority class recall and precision was low.
Choosing top 75% percentile of players with respect to number of shots
Evaluation metrics after choosing cutoff of number of shots for each player(top 75% percentile):
Number of rows = 50000
List of models:
Logistic Regression
SVC
Xgboost
AdaBoost

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.89      | 0.98   | 0.93     | 16711   |
| 1            | 0.60      | 0.22   | 0.32     | 2498    |
|              |           |        |          |         |
| accuracy     |           |        | 0.88     | 19209   |
| macro avg    | 0.75      | 0.60   | 0.63     | 19209   |
| weighted avg | 0.86      | 0.88   | 0.85     | 19209   |

Metrics did not change by much
Oversampling using ADASYN to balance classes and replacing SVC with KNN
Evaluation metrics after oversampling and replacing SVC with KNN:
Number of rows = 50000
List of models:
Logistic Regression
KNN
Xgboost
AdaBoost

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| 0            | 0.88      | 0.93   | 0.90     | 8498    |
| 1            | 0.92      | 0.87   | 0.90     | 8494    |
|              |           |        |          |         |
| accuracy     |           |        | 0.90     | 16992   |
| macro avg    | 0.90      | 0.90   | 0.90     | 16992   |
| weighted avg | 0.90      | 0.90   | 0.90     | 16992   |

Metrics increased by a significant amount with respect to minority class and grid search ran computationally much faster.
Running gridsearch again with entire dataset
Number of rows = entire dataset
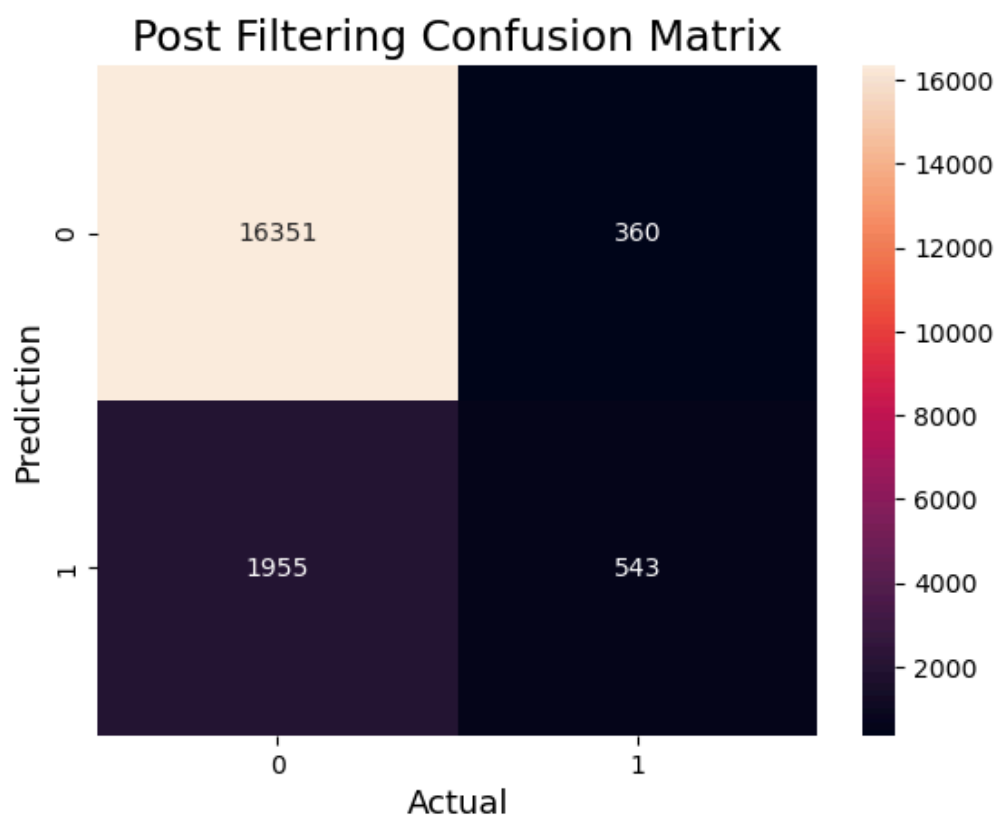List of models:
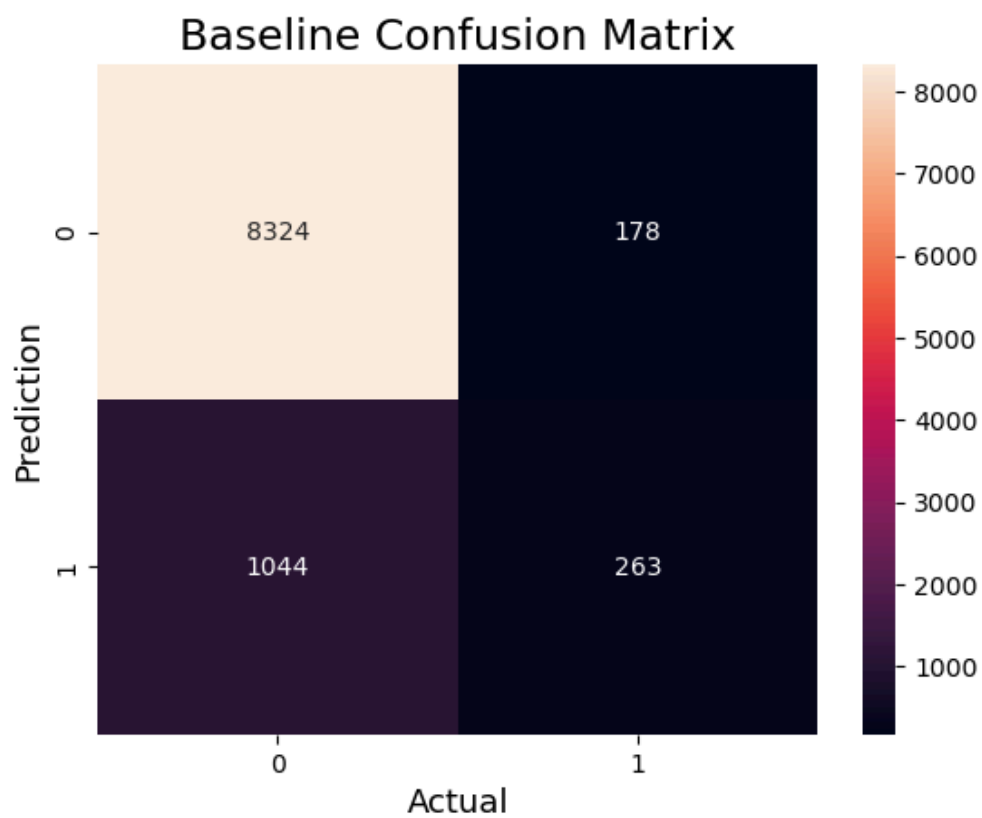Logistic Regression
KNN
Xgboost
AdaBoost

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|

```
           0       0.89      0.93      0.91     16711
           1       0.93      0.89      0.91     16977

    accuracy                           0.91     33688
   macro avg       0.91      0.91      0.91     33688
weighted avg       0.91      0.91      0.91     33688
```
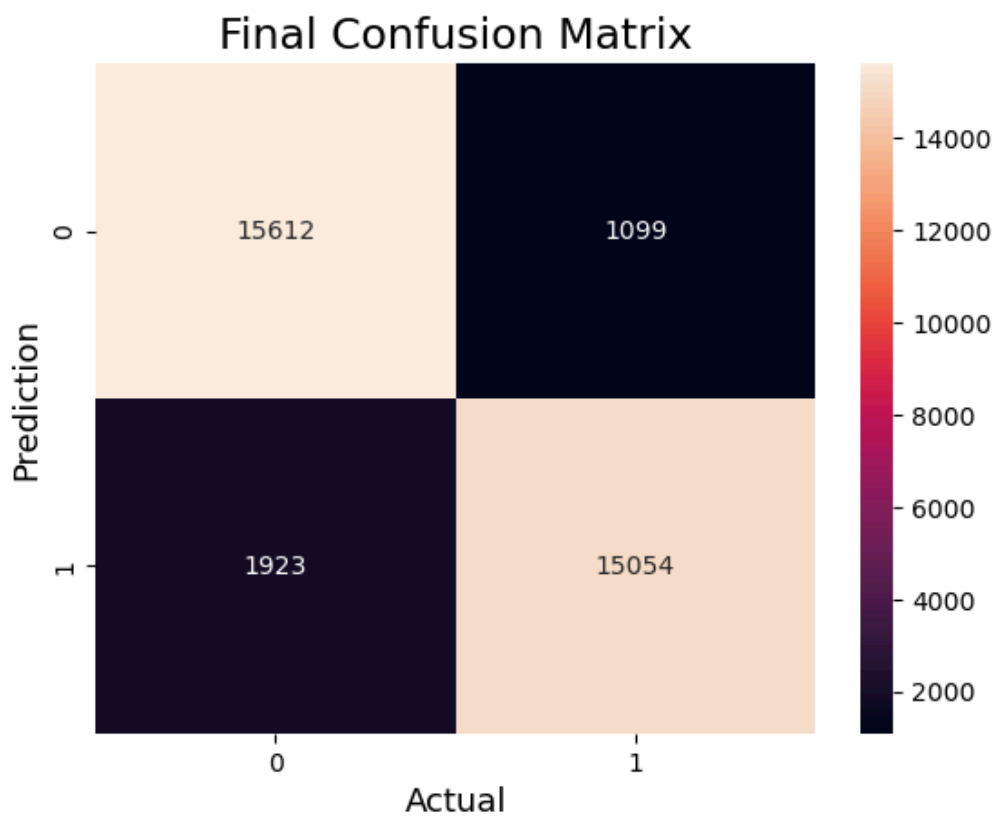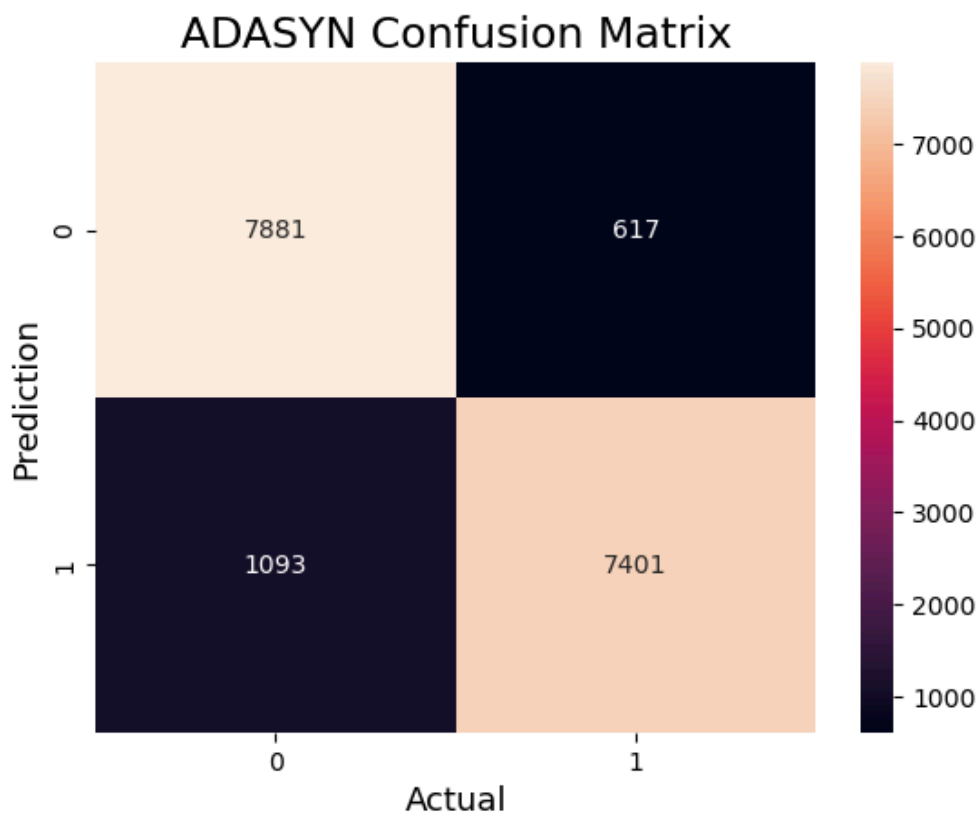
# Confusion Matrices

```python
In [97]:    1  cm_1 = metrics.confusion_matrix(y_test_1, y_preds_1)
            2  cm_2 = metrics.confusion_matrix(y_test_2, y_preds_2)
            3  cm_3 = metrics.confusion_matrix(y_test_3, y_preds_3)
            4  cm_4 = metrics.confusion_matrix(y_test_4, y_preds_4)
            5
            6  sns.heatmap(cm_1,
            7              annot=True,
            8              fmt='g',
            9              xticklabels=range(2),
           10              yticklabels=range(2))
           11  plt.ylabel('Prediction',fontsize=13)
           12  plt.xlabel('Actual',fontsize=13)
           13  plt.title('Baseline Confusion Matrix',fontsize=17)
           14  plt.show()
           15
           16  sns.heatmap(cm_2,
           17              annot=True,
           18              fmt='g',
           19              xticklabels=range(2),
           20              yticklabels=range(2))
           21  plt.ylabel('Prediction',fontsize=13)
           22  plt.xlabel('Actual',fontsize=13)
           23  plt.title('Post Filtering Confusion Matrix',fontsize=17)
           24  plt.show()
           25
           26  sns.heatmap(cm_3,
           27              annot=True,
           28              fmt='g',
           29              xticklabels=range(2),
           30              yticklabels=range(2))
           31  plt.ylabel('Prediction',fontsize=13)
           32  plt.xlabel('Actual',fontsize=13)
           33  plt.title('ADASYN Confusion Matrix',fontsize=17)
           34  plt.show()
           35
           36  sns.heatmap(cm_4,
           37              annot=True,
           38              fmt='g',
           39              xticklabels=range(2),
           40              yticklabels=range(2))
           41  plt.ylabel('Prediction',fontsize=13)
           42  plt.xlabel('Actual',fontsize=13)
           43  plt.title('Final Confusion Matrix',fontsize=17)
           44  plt.show()
```

## Baseline Confusion Matrix

|  | Actual 0 | Actual 1 |
|---|---|---|
| **Prediction 0** | 8324 | 178 |
| **Prediction 1** | 1044 | 263 |

## Post Filtering Confusion Matrix

|  | Actual 0 | Actual 1 |
|---|---|---|
| **Prediction 0** | 16351 | 360 |
| **Prediction 1** | 1955 | 543 |

ADASYN Confusion Matrix



Final Confusion Matrix

**AUC_ROC Curves**

```python
In [98]:  1  fpr_1, tpr_1, _ = metrics.roc_curve(y_test_1, y_preds_1)
          2  auc_1 = metrics.roc_auc_score(y_test_1, y_preds_1)
          3
          4  plt.plot(fpr_1, tpr_1, label='Baseline ROC\nAUC = '+ str(auc_1))
          5  # plt.plot([0, 1], [0, 1], '--', color='grey', label='Worst case')
          6  plt.plot()
          7  plt.xlabel('False Positive Rate')
          8  plt.ylabel('True Positive Rate')
          9  plt.title('ROC Curve')
         10  # plt.legend();
         11
         12  fpr_2, tpr_2, _ = metrics.roc_curve(y_test_2, y_preds_2)
         13  auc_2 = metrics.roc_auc_score(y_test_2, y_preds_2)
         14
         15  plt.plot(fpr_2, tpr_2, label='Filtered ROC\nAUC = '+ str(auc_2))
         16  # plt.plot([0, 1], [0, 1], '--', color='grey', label='Worst case')
         17  plt.plot()
         18  plt.xlabel('False Positive Rate')
         19  plt.ylabel('True Positive Rate')
         20  plt.title('ROC Curve')
         21  # plt.legend();
         22
         23  fpr_3, tpr_3, _ = metrics.roc_curve(y_test_3, y_preds_3)
         24  auc_3 = metrics.roc_auc_score(y_test_3, y_preds_3)
         25
         26  plt.plot(fpr_3, tpr_3, label='ADASYN ROC\nAUC = '+ str(auc_3))
         27  # plt.plot([0, 1], [0, 1], '--', color='grey', label='Worst case')
         28  plt.plot()
         29  plt.xlabel('False Positive Rate')
         30  plt.ylabel('True Positive Rate')
         31  plt.title('ROC Curve')
         32  # plt.legend();
         33
         34  fpr_4, tpr_4, _ = metrics.roc_curve(y_test_4, y_preds_4)
         35  auc_4 = metrics.roc_auc_score(y_test_4, y_preds_4)
         36
         37  plt.plot(fpr_4, tpr_4, label='Final ROC\nAUC = '+ str(auc_4))
         38  plt.plot([0, 1], [0, 1], '--', color='grey', label='Worst case')
         39  plt.plot()
         40  plt.xlabel('False Positive Rate')
         41  plt.ylabel('True Positive Rate')
         42  plt.title('ROC Curve')
         43  plt.legend(bbox_to_anchor = (1.02, 1.02));
```
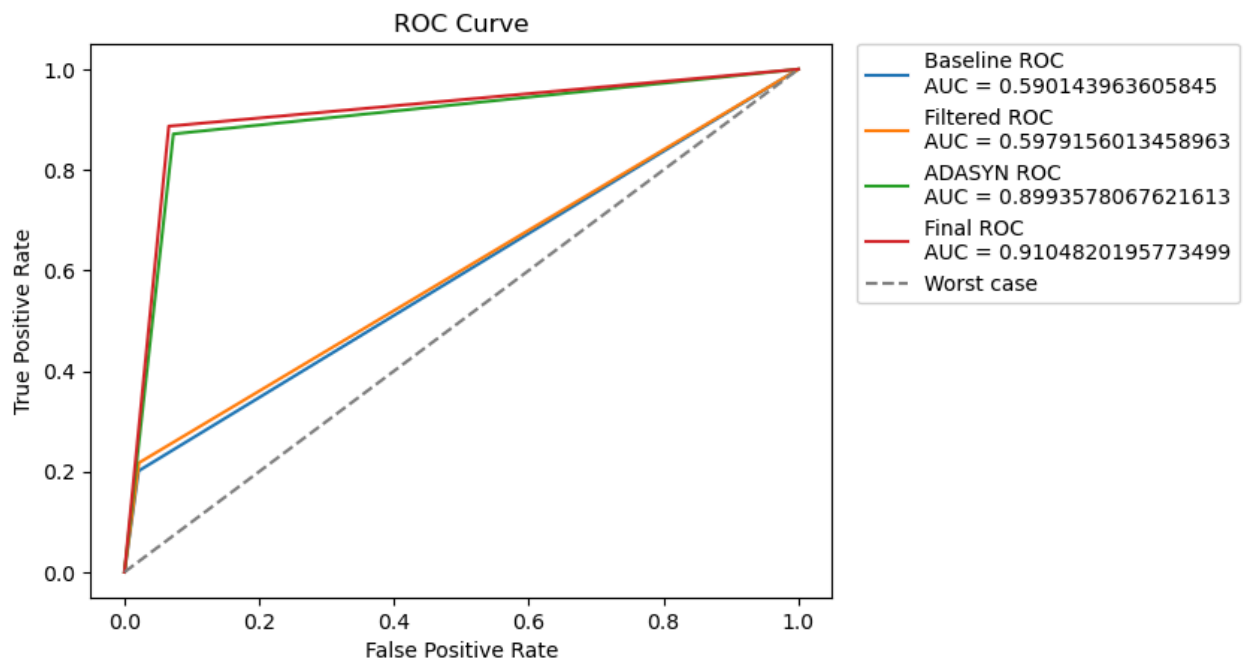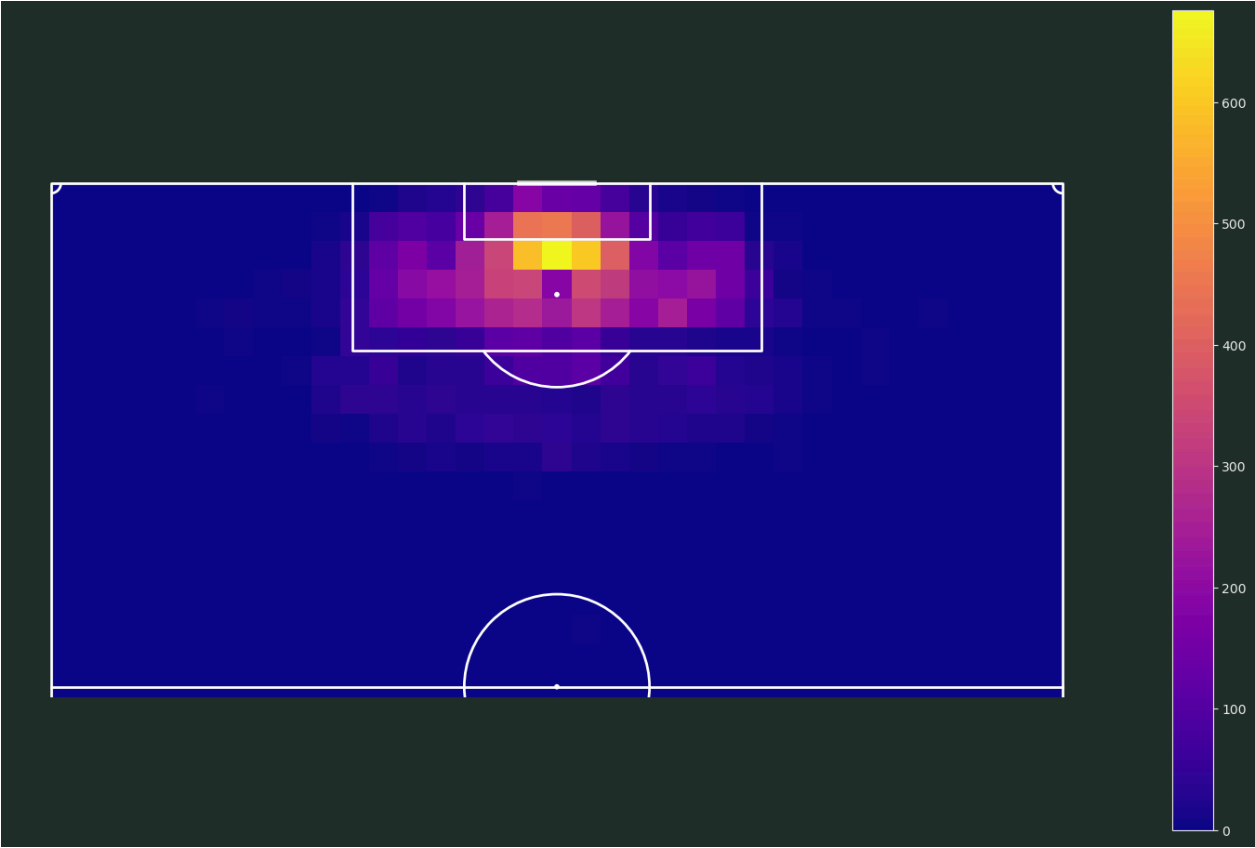
ROC Curve

Baseline ROC
AUC = 0.590143963605845
Filtered ROC
AUC = 0.5979156013458963
ADASYN ROC
AUC = 0.8993578067621613
Final ROC
AUC = 0.9104820195773499
Worst case

## Heatmap

## Predicted Label Heatmap

```python
In [94]:    1  # from highlight_text import fig_text
            2  import matplotlib.patheffects as path_effects
            3
            4  from matplotlib.colors import LinearSegmentedColormap
            5  from scipy.ndimage import gaussian_filter
            6
            7  from mplsoccer import Pitch, VerticalPitch, FontManager, Sbopen
            8  from sklearn.model_selection import train_test_split
            9
           10  mean_with_zero = lambda x: np.sum(x) if not np.sum(np.sum(x)) else 0
           11  pred_heatmap = pd.concat([pd.DataFrame(X_test), pd.Series(y_preds_4,index=X_test.index
           12
           13  pred_heatmap = pred_heatmap.rename(columns = {0 : 'Result'})
           14
           15  plt.figure(dpi=4800)
           16  path_eff = [path_effects.Stroke(linewidth=0.5, foreground='black'),
           17              path_effects.Normal()]
           18  pitch = VerticalPitch(pad_bottom = 1, half = True, goal_type = 'line', goal_alpha = 0.8
           19                        pitch_type = 'custom', pitch_length = 99.5, pitch_width = 100,
           20                        line_zorder=2, line_color='white', corner_arcs=True, pitch_color=
           21
           22  fig, ax = pitch.draw(figsize=(15,15))
           23  fig.set_facecolor('#22312b')
           24
           25  bin_statistic = pitch.bin_statistic(pred_heatmap.X * 100, pred_heatmap.Y* 100,
           26                                      values = pred_heatmap.Result, statistic='sum', bins
           27
           28  bin_statistic['statistic'] = gaussian_filter(bin_statistic['statistic'], -10)
           29
           30  pcm = pitch.heatmap(bin_statistic, ax=ax, cmap='plasma')
           31  cbar = fig.colorbar(pcm, ax=ax, shrink=0.6)
           32
           33  cbar.outline.set_edgecolor('#efefef')
           34  cbar.ax.yaxis.set_tick_params(color='#efefef')
           35
           36  ticks = plt.setp(plt.getp(cbar.ax.axes, 'yticklabels'), color='#efefef')
           37  # labels = pitch.label_heatmap(bin_statistic, color='black', fontsize=10,
           38  #                              ax=ax, ha='center', va='center',
           39  #                              str_format='{:.2f}%', path_effects=path_eff)
```
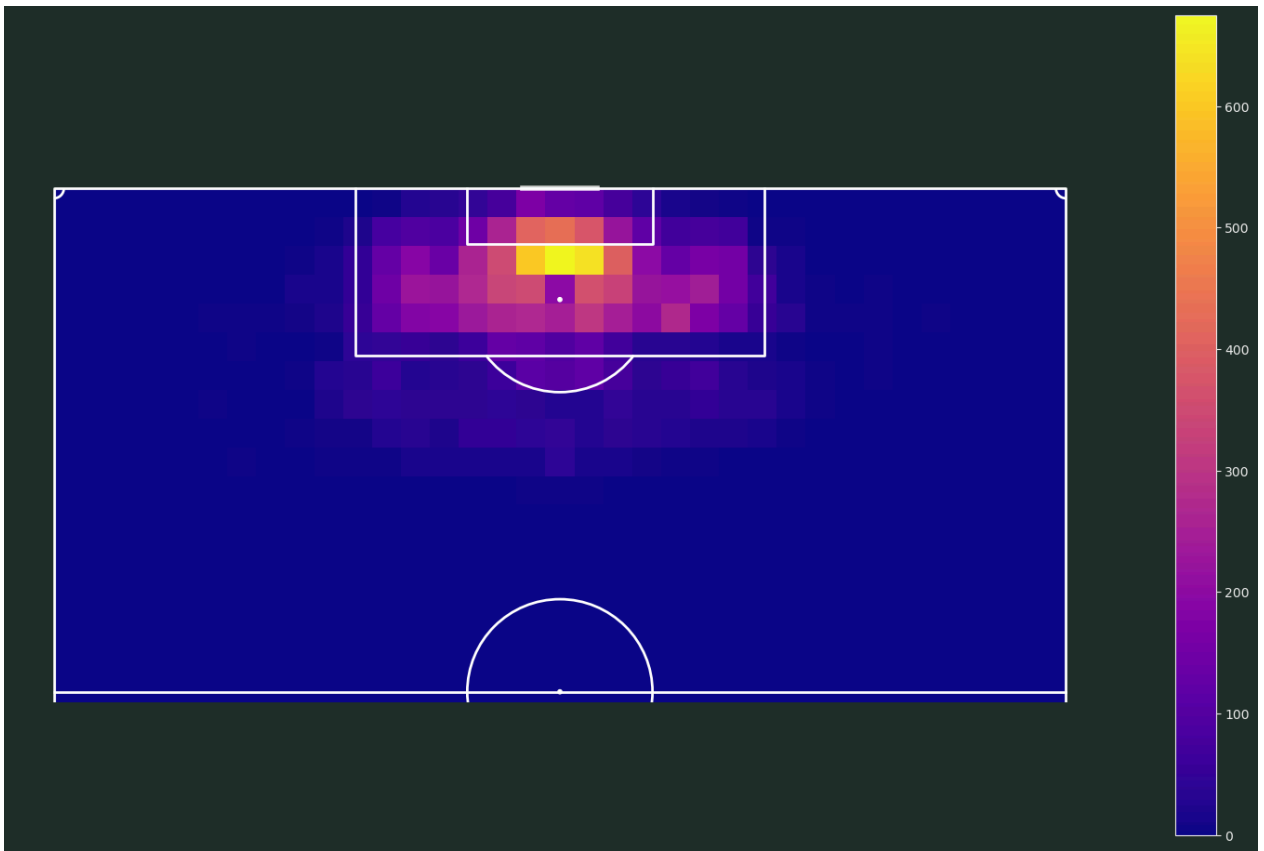
```
<Figure size 30720x23040 with 0 Axes>
```

## True Label Heatmap

```
In [95]:   1  # from highlight_text import fig_text
           2  import matplotlib.patheffects as path_effects
           3
           4  from matplotlib.colors import LinearSegmentedColormap
           5  from scipy.ndimage import gaussian_filter
           6
           7  from mplsoccer import Pitch, VerticalPitch, FontManager, Sbopen
           8  from sklearn.model_selection import train_test_split
           9
          10  mean_with_zero = lambda x: np.sum(x) if not np.sum(np.sum(x)) else 0
          11  pred_heatmap = pd.concat([pd.DataFrame(X_test), pd.Series(y_test_4,index=X_test.index)
          12
          13  # pred_heatmap = pred_heatmap.rename(columns = {0 : 'Result'})
          14
          15  plt.figure(dpi=4800)
          16  path_eff = [path_effects.Stroke(linewidth=0.5, foreground='black'),
          17              path_effects.Normal()]
          18  pitch = VerticalPitch(pad_bottom = 1, half = True, goal_type = 'line', goal_alpha = 0.8
          19                        pitch_type = 'custom', pitch_length = 99.5, pitch_width = 100,
          20                        line_zorder=2, line_color='white', corner_arcs=True, pitch_color=
          21
          22  fig, ax = pitch.draw(figsize=(15,15))
          23  fig.set_facecolor('#22312b')
          24
          25  bin_statistic = pitch.bin_statistic(pred_heatmap.X * 100, pred_heatmap.Y* 100,
          26                                      values = pred_heatmap.result, statistic='sum', bins
          27
          28  bin_statistic['statistic'] = gaussian_filter(bin_statistic['statistic'], -10)
          29
          30  pcm = pitch.heatmap(bin_statistic, ax=ax, cmap='plasma')
          31  cbar = fig.colorbar(pcm, ax=ax, shrink=0.6)
          32
          33  cbar.outline.set_edgecolor('#efefef')
          34  cbar.ax.yaxis.set_tick_params(color='#efefef')
          35
          36  ticks = plt.setp(plt.getp(cbar.ax.axes, 'yticklabels'), color='#efefef')
          37  # labels = pitch.label_heatmap(bin_statistic, color='black', fontsize=10,
          38  #                              ax=ax, ha='center', va='center',
          39  #                              str_format='{:.2f}%', path_effects=path_eff)
```

<Figure size 30720x23040 with 0 Axes>