

```
In [28]: from polars import read_csv, col, Int64, min_horizontal, concat, Series, len as pl_len
from sklearn.model_selection import train_test_split, RandomizedSearchCV
from sklearn.metrics import log_loss, roc_auc_score, brier_score_loss
from xgboost import XGBClassifier
from numpy import corrcoef
from pathlib import Path
import pickle
```

Load data

```
In [2]: df = read_csv("ML_TAKES_ENCODED.csv")
```

Generate Features

Strike Zone

Turn Strikes and Balls from categorical to numerical feature

```
In [3]: df = df.with_columns((col("PITCHCALL") == "StrikeCalled").cast(Int64).alias("IS_STRIKE"))
```

Strike Zone Features

Determine if in zone

```
In [4]: df = df.with_columns(
    ((col("BOT_ZONE") <= col("PLATELOCHEIGHT"))
     & (col("PLATELOCHEIGHT") <= col("TOP_ZONE"))).cast(Int64).alias("IN_ZONE"))
```

Determine how close to the edge of the zone the ball was

```
In [5]: df = df.with_columns(
    min_horizontal(
        (col("PLATELOCHEIGHT") - col("BOT_ZONE")).abs(),
        (col("PLATELOCHEIGHT") - col("TOP_ZONE")).abs()
    ).alias("NEAR_VERT_EDGE")
)
```

```
In [6]: PLATE_CENTER_WIDTH_FT = 0.708
```

```
In [7]: df = df.with_columns((col("PLATELOCSIDE") - PLATE_CENTER_WIDTH_FT).abs().alias("NEAR_HORZ_EDGE"))
```

```
In [8]: df = df.with_columns(PLATE_LOC_HEIGHT.col("PLATELOCHEIGHT"))
df = df.with_columns(PLATE_LOC_HEIGHT.col("PLATELOCHEIGHT"))
```

Pitch Movement Features

Approach Angle

```
In [9]: df = df.with_columns(VERT_APPROACH.col("VERTAPPRANGLE"))
df = df.with_columns(HORZ_APPROACH.col("HORZAPPRANGLE"))
```

Spread and Break

TODO: Need to test if these features are relevant or just noise

```
In [10]: df = df.with_columns(INDUCED_VERT_BREAK.col("INDUCEDVERTBREAK"))
df = df.with_columns(HORZ_BREAK.col("HORZBREAK"))
df = df.with_columns(REL_SPEED.col("RELSPEED"))
```

Final Features

```
In [11]: features = [
    "PLATE_LOC_HEIGHT",
    "PLATE_LOC_SIDE",
    "TOP_ZONE",
    "BOT_ZONE",
    "IN_ZONE",
    "NEAR_VERT_EDGE",
    "NEAR_HORZ_EDGE",
    "BALLS",
    "STRIKES",
    "VERT_APPROACH",
    "HORZ_APPROACH",
    "INDUCED_VERT_BREAK",
    "HORZ_BREAK",
    "REL_SPEED"
]
```

Remove nulls from critical features

```
In [12]: df = df.drop_nulls(subset=features + ["IS_STRIKE"])
```

Train Model

```
In [13]: print(f"Rows: {len(df)}")
print(f"Model: Strike rate: {df.select('IS_STRIKE').mean().item():.3f}")
print(f"Features: {len(features)}")
print(f"Years: {df.select('GAME_YEAR').unique().to_series().to_list()}")
print(f"Unique Catchers: {df.select('CATCHER_ID').n_unique()}"
```

Rows: 1109138

Strike rate: 0.315

Features: 14

Years: [2022, 2023, 2021]

Unique catchers: 164

Test Brier: 0.0465

==== Feature Importance ==

STRIKES: 0.0334

BOT_ZONE: 0.0074

BALLS: 0.0049

VERT_APPROACH: 0.0042

HORZ_APPROACH: 0.0041

HORZ_BREAK: 0.0028

TOP_ZONE: 0.0024

REL_SPEED: 0.0020

INDUCED_VERT_BREAK: 0.0012

PLATE_LOC_HEIGHT: 0.0350

PLATE_LOC_SIDE: 0.1571

NEAR_HORZ_EDGE: 0.1373

NEAR_VERT_EDGE: 0.0830

IN_ZONE: 0.0530

TOP_ZONE: 0.0530

PLATE_CENTER_WIDTH_FT: 0.708

PLATE_CENTER_HEIGHT_FT: 0.0708

PLATE_CENTER_SIDE_FT: 0.0229

NEAR_VERT_EDGE: 0.0702

NEAR_HORZ_EDGE: 0.0702

INDUCED_VERT_BREAK: 0.0571

HORZ_BREAK: 0.0571

PLATELOCHEIGHT: 0.0571

PLATELOCSIDE: 0.0571

PLATELOCWIDTH: 0.0571

PLATELOCHEIGHT: 0.0571

PLATELOCWIDTH: 0.0571

PLATELOCSIDE: 0.0571