



AeroProphet



# Team 2.2

Hailee Schuele



hschuele@berkeley.edu

Landon Yurica



lyurica@berkeley.edu

Sreeram Ravinoothala



sreeram@berkeley.edu

Nick Johnson



nickjohnson@berkeley.edu

# Business Case

**Objective:** Predict departure delays greater than 15 minutes, 2 hours before takeoff

## Benefits:

- For Airlines :
  - Proactively address issues
  - Improve efficiency
- For Passengers:
  - Adapt plans
  - Make informed decisions

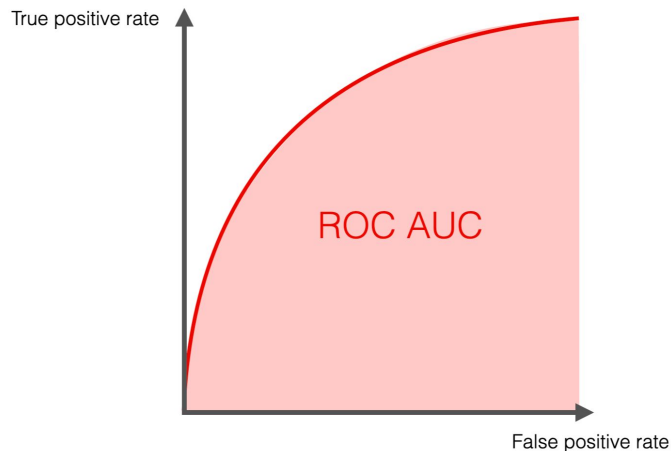


# Metric Selection

**Primary Metric:**

$$\text{F1 Score} = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$$

**Secondary Metric:**



# Final Result

## **Out of Sample Performance:**

Primary Metric:

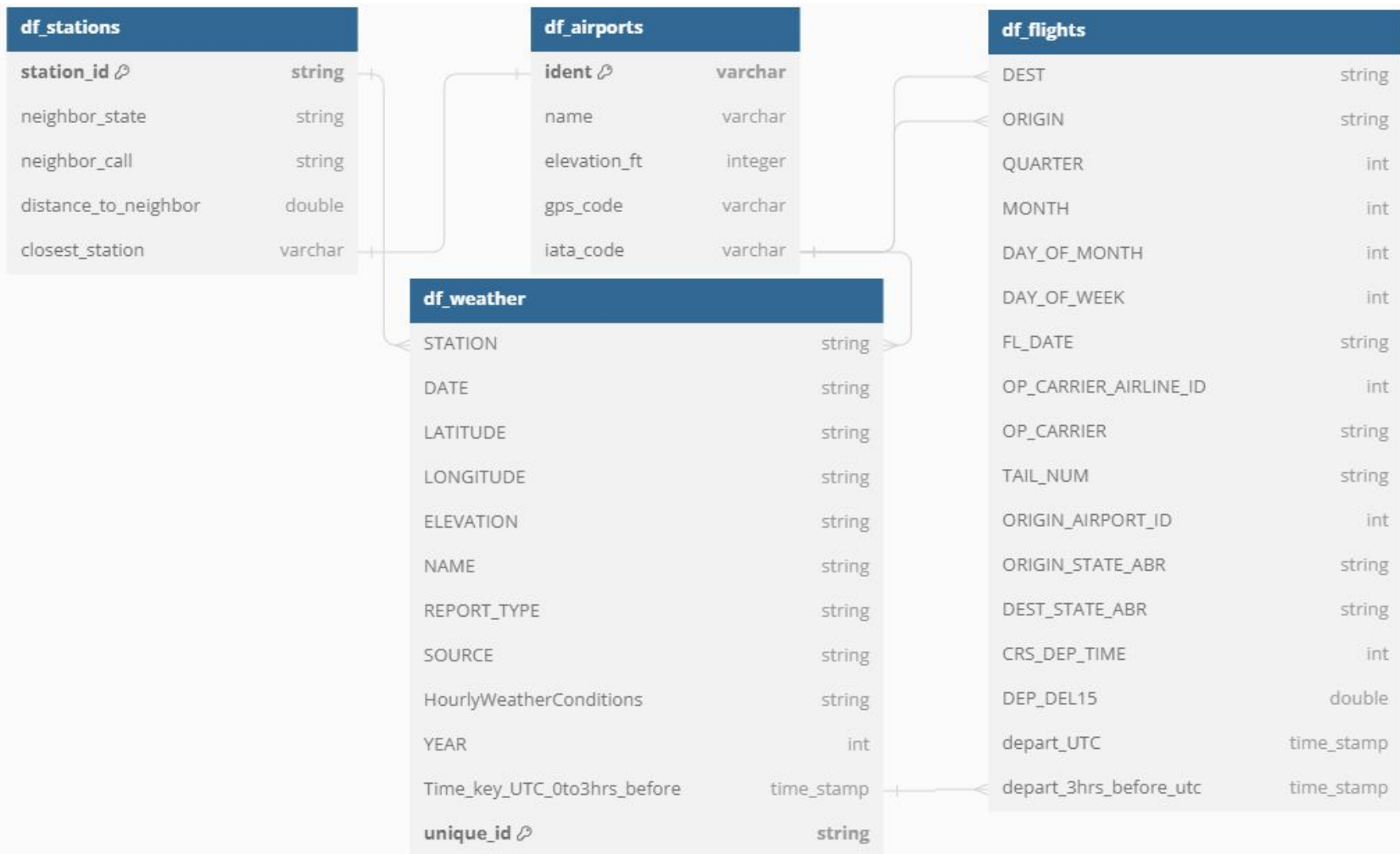
F1 Score = 0.803

Secondary Metric:

AUC = 0.734

# Datasets

- 2015-2021 US Flight records
- 2015-2021 International Weather records
- Weather Station location data
- Airport location data



# Data Join Details

Table	Rows	Columns	Memory (GB)	Run Time (HH:MM:SS)
df_flights	74,177,433	109	2.93	None
df_weather	898,983,399	124	35.05	None
df_stations	5,004,169	12	1.3	None
df_airports	57,421	12	0.01	None
df_FSW	72,515,921	292	43.5	02:13:00
df_FSW_Clean	64,457,088	87	64.78	06:18:00

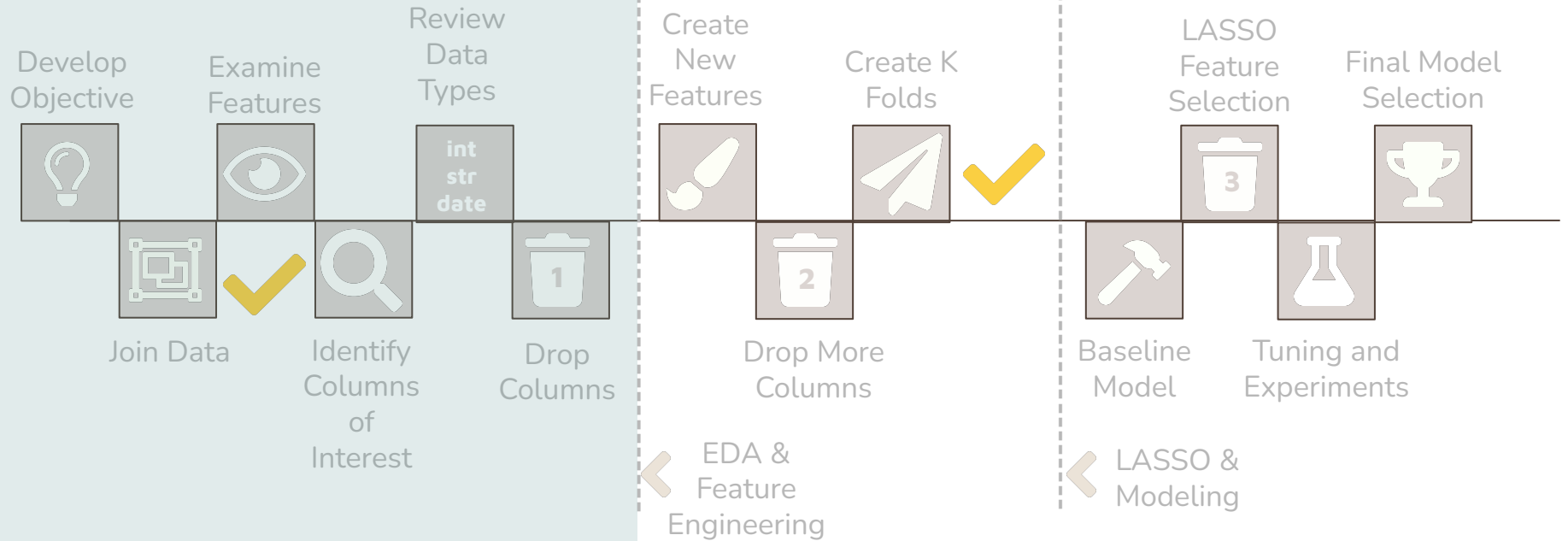
Cluster Details: DBR 13.3 LTS ML, Spark 3.4.1, Scala 2.12, Standaard\_DS3\_v2, 14GB, 4 Cores, Standard\_DS3\_v2, 84GB, 24 Cores, 1-6 workers



# Data Cleaning

- Dropped columns containing:
  - Substantial missing data
  - Potential data leakage
  - Redundant information
- Dropped duplicate rows
- Mean imputation

# Data Join & Cleaning



# Feature Engineering

## Existing Engineered Features

### **previous\_flight\_delay**

- Delay time of the previous flight



### **state\_to\_region**

- Convert state abbreviations to geographic region



### **Plane\_Delays\_last\_24h**

- Sum of departure delays for each plane (based on tail number) within the last 24 hours



### **airport\_delays\_in\_previous\_24\_hours**

- The number of flights delayed at the airport in the last 24 hours



### **days\_to\_nearest\_holiday**

- The number of days to the nearest holiday



## Newly Engineered or Added

→ Added min & max hourly weather values

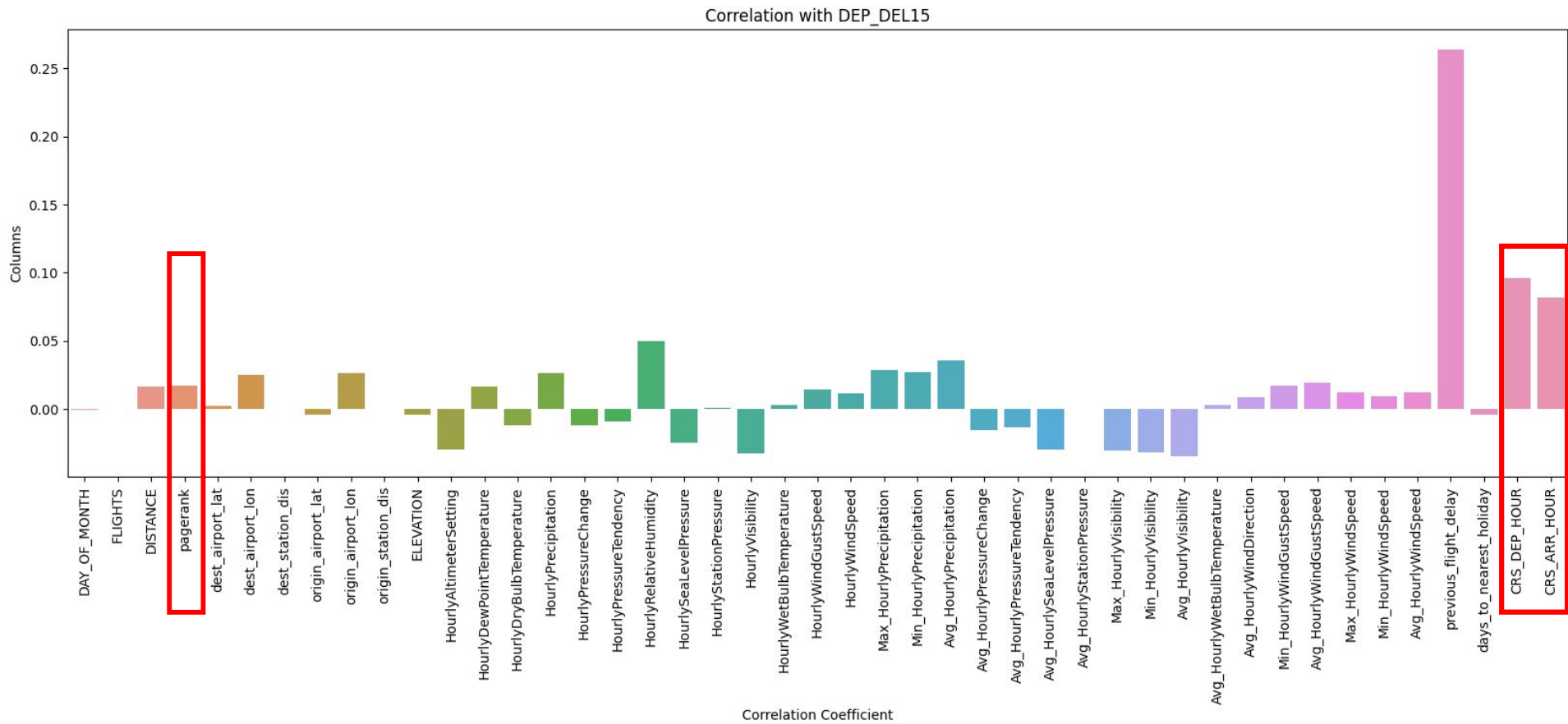
→ Airport Page Rank

→ Airport Size

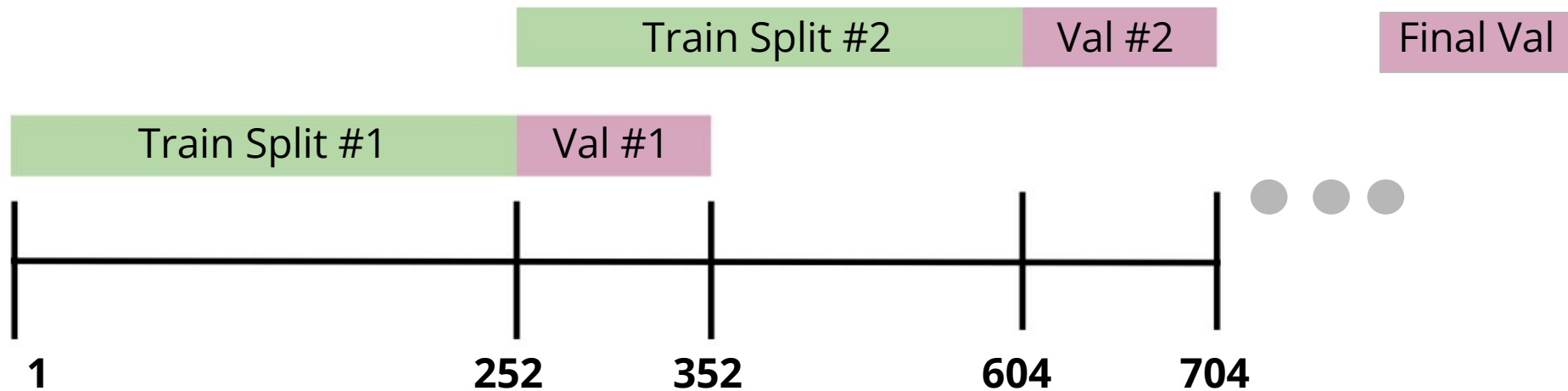
→ Scheduled depart and arrival hour (local time)

→ Class + Recency

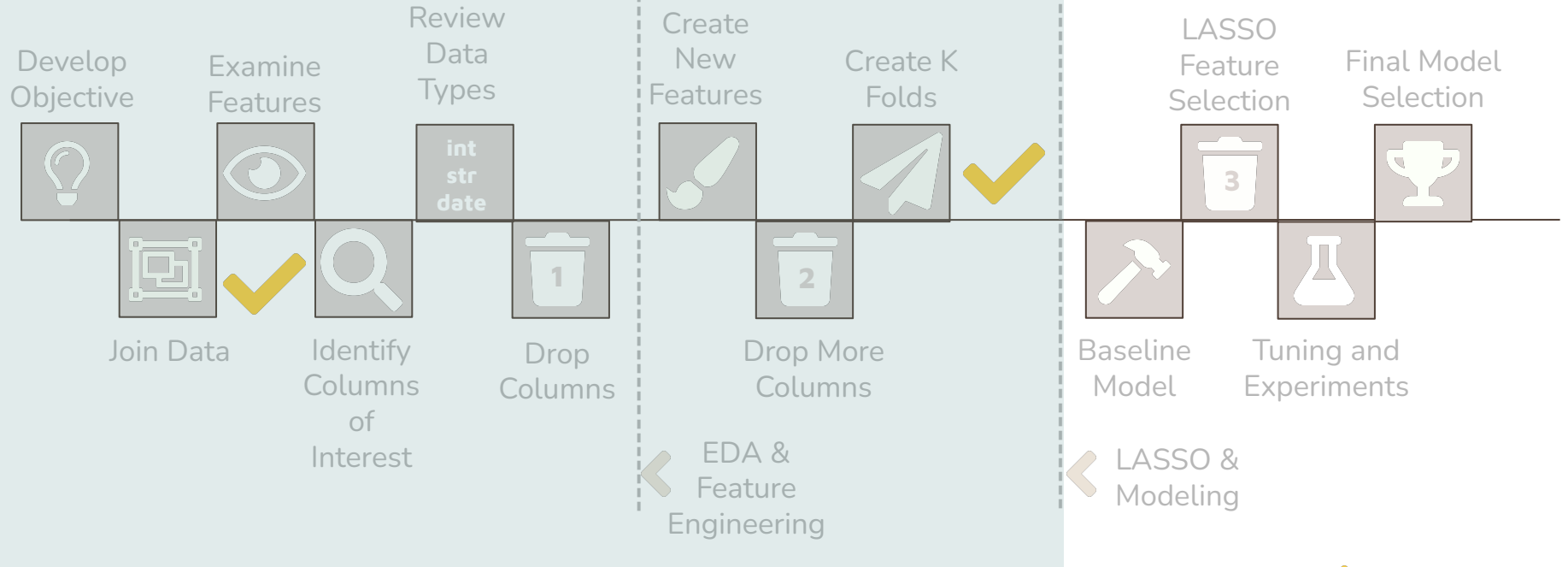
# EDA



# Cross Folds



# EDA & Feature Engineering



# Modeling Workflow

## Logistic regression

- 70 features
- Includes derived features on folds
- Validation data included
- Used Vector representation of features

## Lasso

- 64 features
- Excluded derived
- On folds
- Retained features kept in at least one fold
- Got 38 features

## Logistic regression

- 38 Lasso features
- Additionally 5 Derived features
- Ran on both folds as well as final validation set

## Random Forest

- One fold using lasso dataset
- Included derived columns
- No significant impact with Hyper params
- Used multiple seeds, avg for eval
- Ran on all folds
- Numtree=10, maxdepth=5

## GBDT

- One fold using lasso dataset
- Included derived columns
- No significant impact with Hyper params
- Used multiple seeds, avg for eval
- Ran on all folds
- iter=10, stepsize=1.0

# Neural Network Experiments

## MLP

- Experimented with different seeds, iterations, number of nodes in hidden layers with a 10% dataset
- Results showed maxiter=100 provided good result
- Both single hidden layer and two hidden layers come back with close results
- Lasso data used

## LSTM

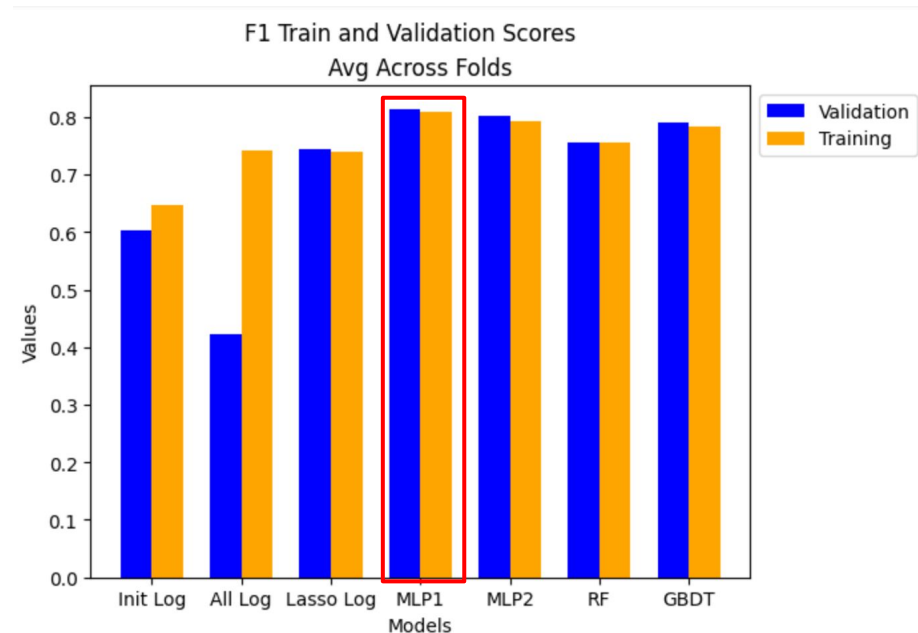
- Ran for 20 epochs, on all folds. Loss converged at 10 epochs
- Used BCELoss, Adam optimizer
- Captured scores on fold validations and final validation set
- Used only 10% of the data
- Taking too long
- Future

## Final MLP

- Single hidden layer



# Intermediary Model Results



	Train F1	Validate F1
Initial Logistic	0.648	0.604
Secondary Logistic	0.741	0.422
LASSO Logistic	0.740	0.743
MLP 2 Hidden Layers	0.792	0.801
Random Forest	0.755	0.756
Gradient Boosted Decision Trees	0.783	0.791
MLP 1 Hidden Layer	0.808	0.814

# Final Model Results

	Train F1	Test F1
Final MLP	0.805	0.803

# Conclusion

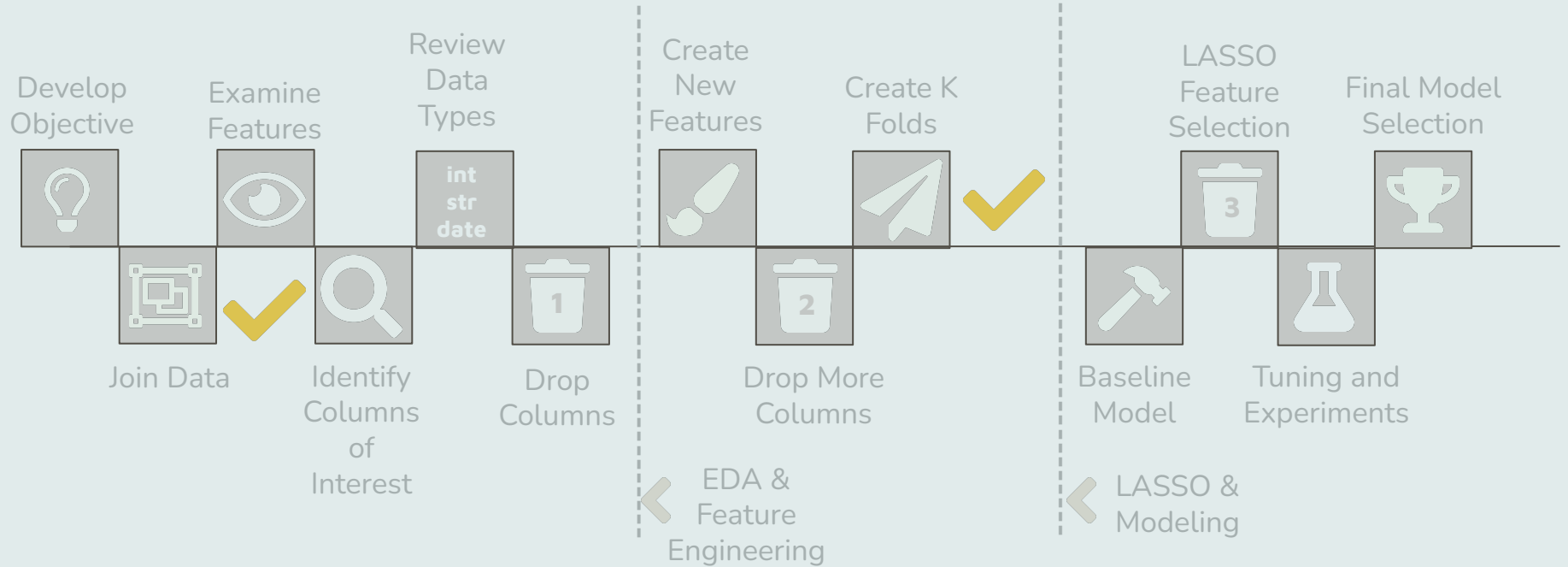
## Final Model

- 38 LASSO Features + 5 Derived Features
- 35% increase in F1 score from first to last model

## Next Steps

- More Derived Features
- Additional Hyperparameter Tuning
- Ensemble Methods
- Deep Learning and More Complex Models

# Thank You!



Data Checkpoint