

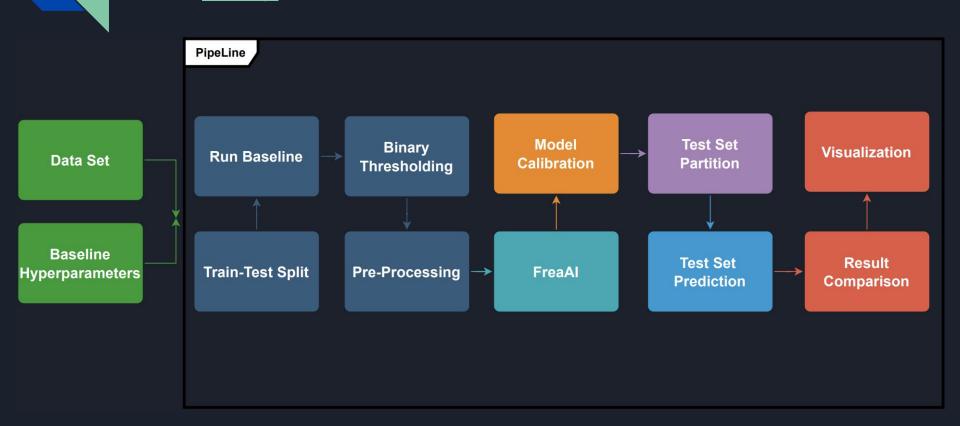
Boaz Zimbler, Adam Vinestock, Itay Shapira, Omri Newman Reichman University January 2, 2023

# Table of Contents

- Pipeline Architecture
- Pre-Processing
- FreaAl
- Model Calibration
- Test Set Partition and Prediction
- Additional Optimization Approaches
- Code Demo
- Results and Comparisons
- Next Steps
- Q&A

# Pipeline Architecture

GitHub Repo



# Pre-Processing

#### **Categorical**

VehGas	Count
Regular	345877
Diesel	332136

Area	Count
С	191880
D	151596
Е	137167
Α	103957
В	75459
F	17954

#### **Numerical**

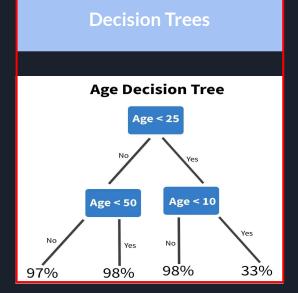


#### FreaAl

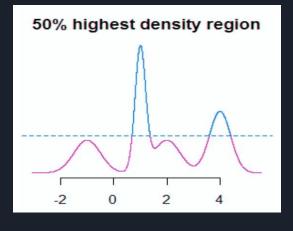
What

**Automated Extraction** of Data Slices to Test ML Models Identifying **Weak Data** which **underperforms in the model** 

How



HDR



#### FreaAl

Business **Decision Making** 

"Focusing" on specific market

**Rely on model** only in **specific** market

Why

Improving model by dealing with weak data

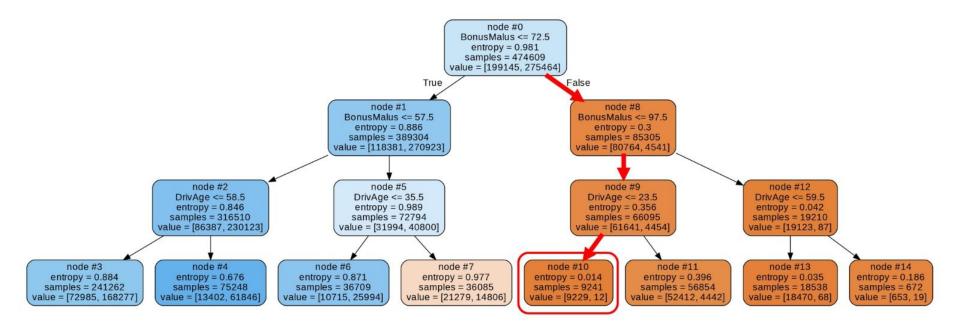
**Getting more Weak Data** to train on

**Splitting models** to reduce generalization error





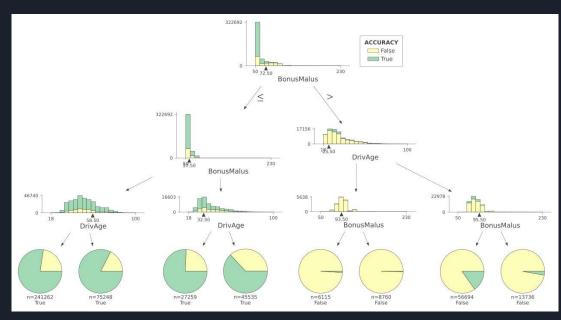
#### FreaAl



### Model Calibration



- Based on FreaAl we identify 'Low quality' subsets of data
- Extract Low quality data from the training set splitting it into two - Low and High quality data
- Train two new XGBoost models calibrated to data sets



### Test Set Partition and Prediction



Partition test set into two
High quality data | Low quality data
as determined by FreaAl decision trees



Predict on partitioned test sets using corresponding models from Calibration step

# Additional Optimization Approaches

'High quality' vs 'Low quality' data performance

features = [features]

Creating more robust data partition - low quality data detection using ensemble of trees

```
[31] # Extracting more robust poor quality data
    def mask_repeated_elemer [108] # Preproccessing the data
       unique elements, cou
                            features flagged = nn.concatenate((french df.columns.["FLAGGED"]))
       repeated elements =
                                               ----- Test data summary using threshold of 4 tree's
       return repeated elem
                            mask = frenc
                            flags = np.w
                                               Baseline
    def flagged data(df, out
                            X flagged =
                                                predictions within threshold: 64.02%, MSE: 31.215869537406565
       poor data = np.empty
       n trees = len(output
                            df flagged =
       for i in range(n tre
        features = output.
                                               Low quality data
                            df flagged['
        if isinstance(feat
                                                predictions within threshold: 19.01%, MSE: 67.26537746807695
          pass
                            df flagged['
        else:
                            df flagged.h
           features = [fe
                                               High quality data
        tree model = outpu
        leaf = output.iloc
                                                predictions within threshold: 76.07%, MSE: 9.933800932185752
                                Area Ve
        leaf indices = tre
        bad indices = np.a
                             0.0
        poor data = np.cor
                                               Weighted sum
       flagged data = mask
                                                predictions within threshold: 65.91%, MSE: 31.422290499466996
       print(f"Num of min a
                             1 0.0
       print(f"Total size o
       return flagged data
                                                           52.0
                                                                                                          22.0 0.000000
                                                                                                                             0.0
    def classify samples(X,
                                           7.0
                                                   0.0
                                                                       50.0
                                                                                  0.0
                                                                                          1.0
                                                                                                  76.0
                                                                                                          72.0 0.000000
                                                                                                                             0.0
       low quality indices
       for i in range(trees
                             4 1.0
                                           7.0
                                                   0.0
                                                           46.0
                                                                       50.0
                                                                                  0.0
                                                                                          1.0
                                                                                                  76.0
                                                                                                          72.0 1.190476
                                                                                                                              0.0
        features = trees.i
        if isinstance(feat
          pass
        else:
```

# Low-quality Data Generation

```
Test data summary using CWGAN-GP

Baseline
predictions within threshold: 64.02%, MSE: 31.215869537406565

IMPLOW quality data
predictions within threshold: 8.02%, MSE: 46.486560008961625

High quality data
predictions within threshold: 76.07%, MSE: 27.923652375874138

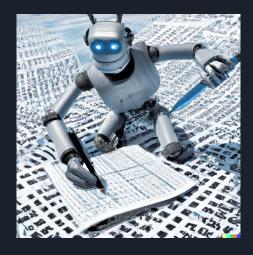
Weighted sum
predictions within threshold: 64.04%, MSE: 31.228953100811516
```

Oversampling using Gaussian Copulas

CWGAN-GP data generation using YData

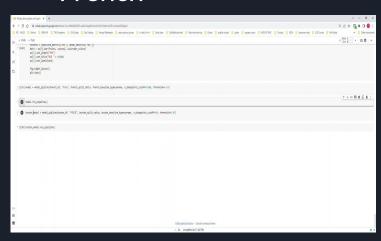
# **Splitting models** to reduce generalization error



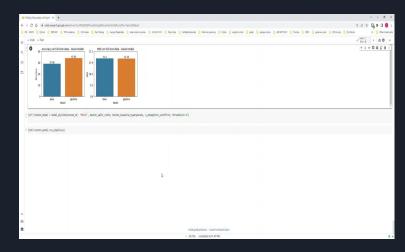


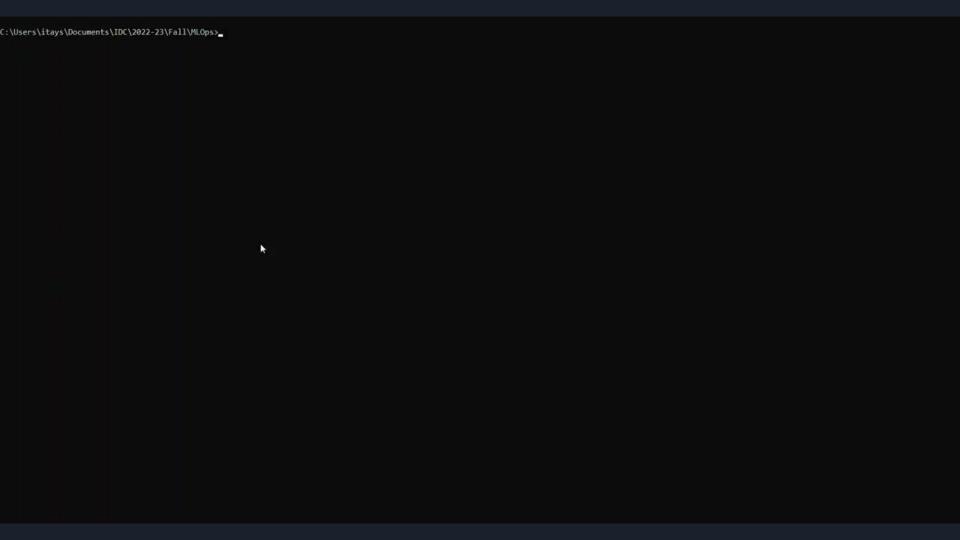
## Code Demo

## French

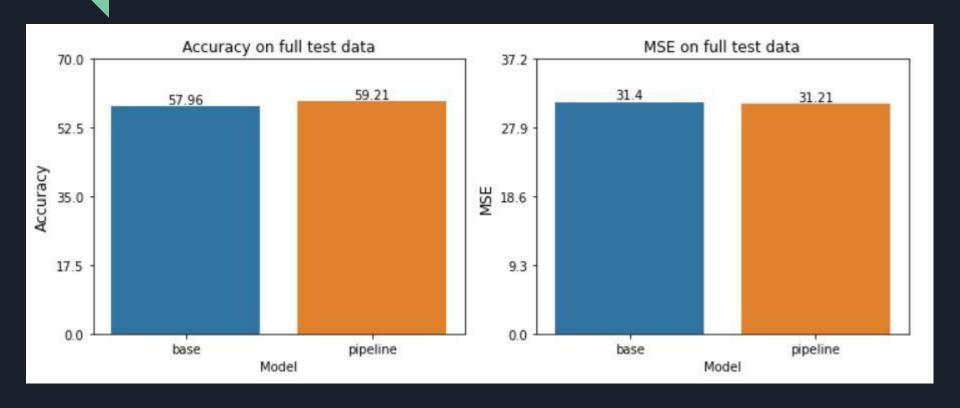


#### Boston

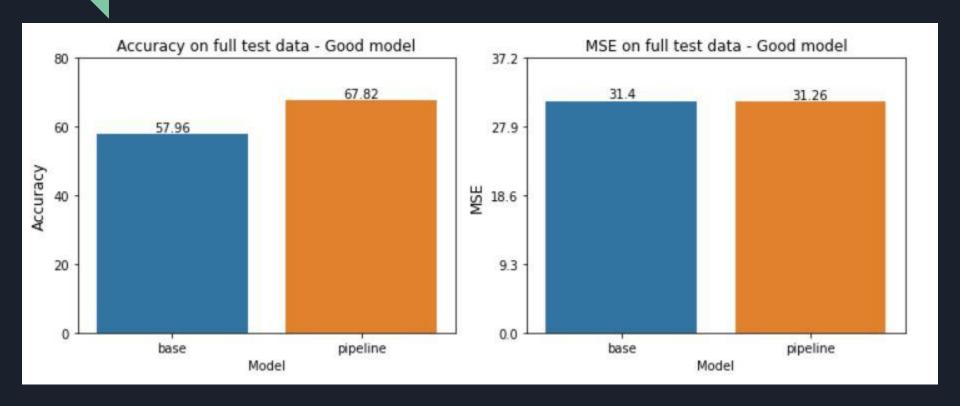




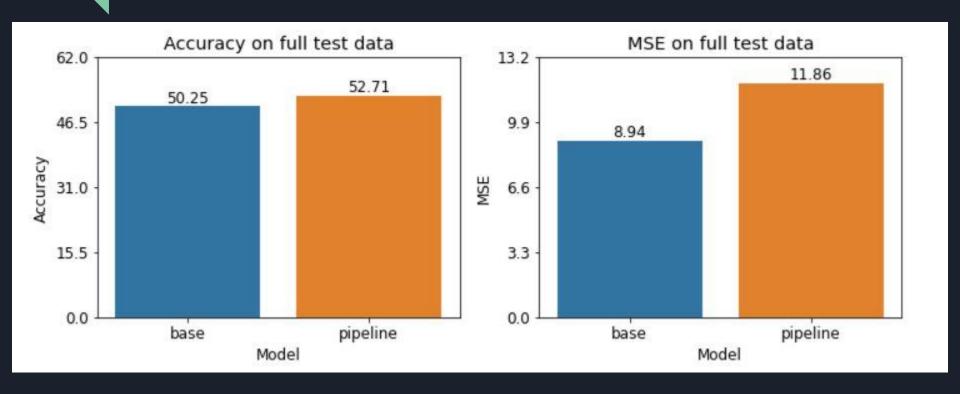
# Results and Comparison (French)



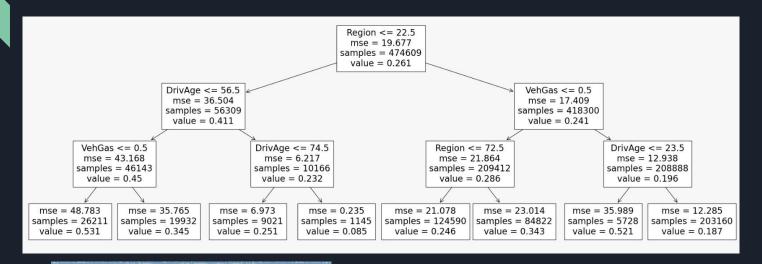
# Results and Comparison (French)



## Results and Comparison (Boston)



## Next Steps



We Want More Data

- More Data Collection
- Implementing HPD on top of Decision Trees
- Further explore FreaAl using Decision Tree Regressors

Q&A?

Thanks!

