



# Vnomics-1 Final Presentation



Daniel He



Matthew Taruno



James Sastrawan



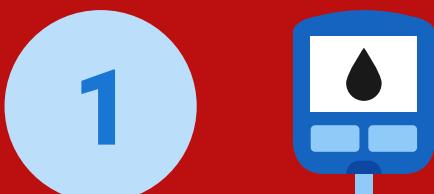
Linzan Ye



Vatsal Agarwal

# BUSINESS CONTEXT

## Core Issue



DPF failure is costly.



Sudden Breakdown

**GOAL**



Build a model, based on data science techniques, to predict DPF failures in trucks at least two weeks in advance

## DRIVING BEHAVIOR

1

2

3

4

Idling

Shifting

Accelerating  
Speeding

## TRIP CONTEXT

1

2

3

4

Load

Truck Config

Environment  
Speeding

# REQUIREMENTS

I

Produce a model to accurately predict DPF failure in trucks by establishing early indicators of such failure

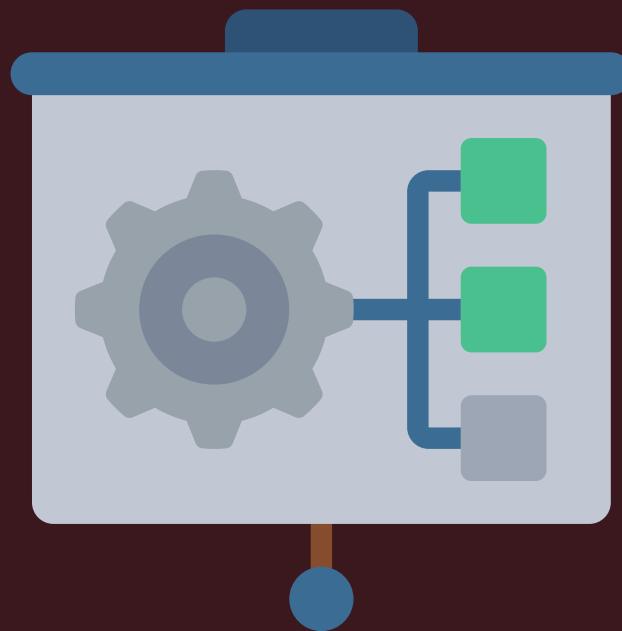
II

Alert customers of potential need for regeneration ahead of time to avoid costly roadside breakdowns

III

Ideally, a truck with any failure indicators should not be missed, so the model should have a high recall score.

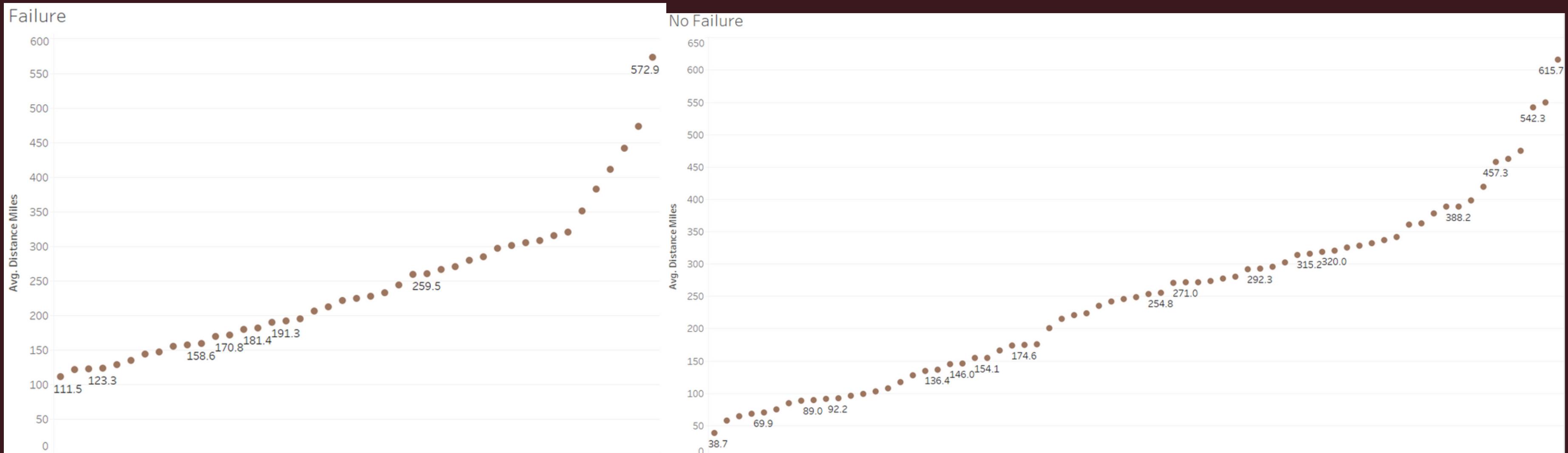
# *Project Deliverables*



- 1 Baseline model and iteration over Models to improve performance
- 2 Weekly Collaboration and Reports Tracking Progress
- 3 Final Model with High Recall Score

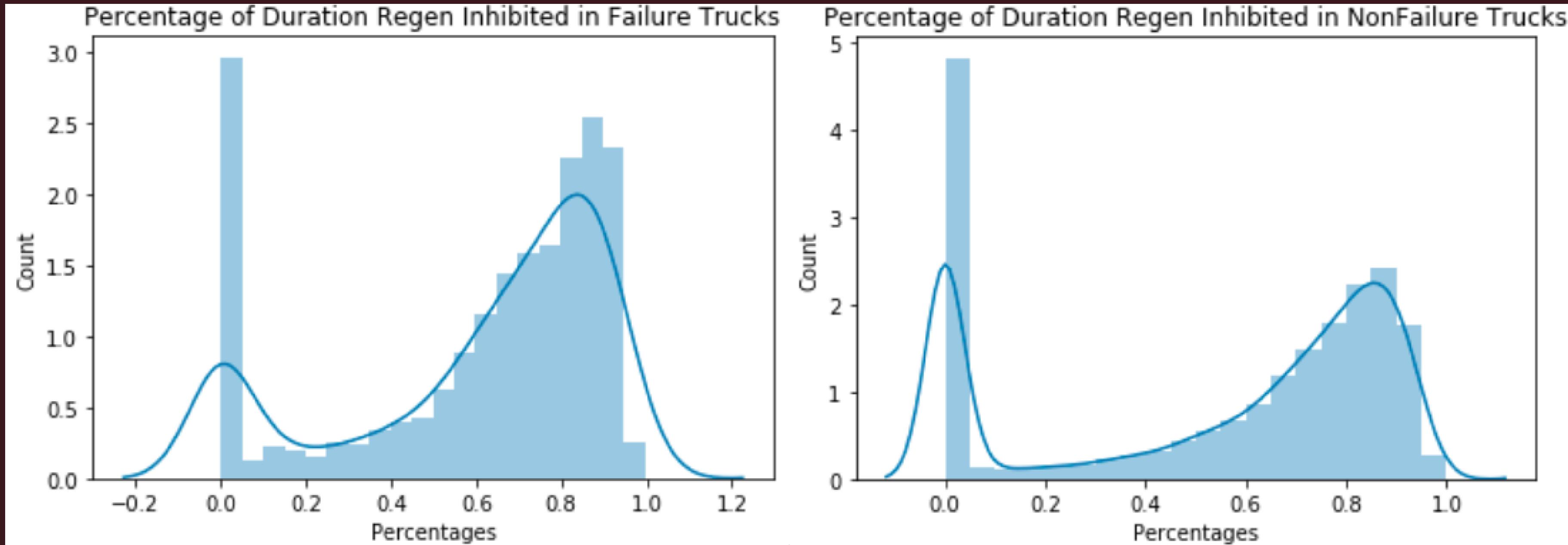


# DISTRIBUTION OF TRUCK DISTANCE



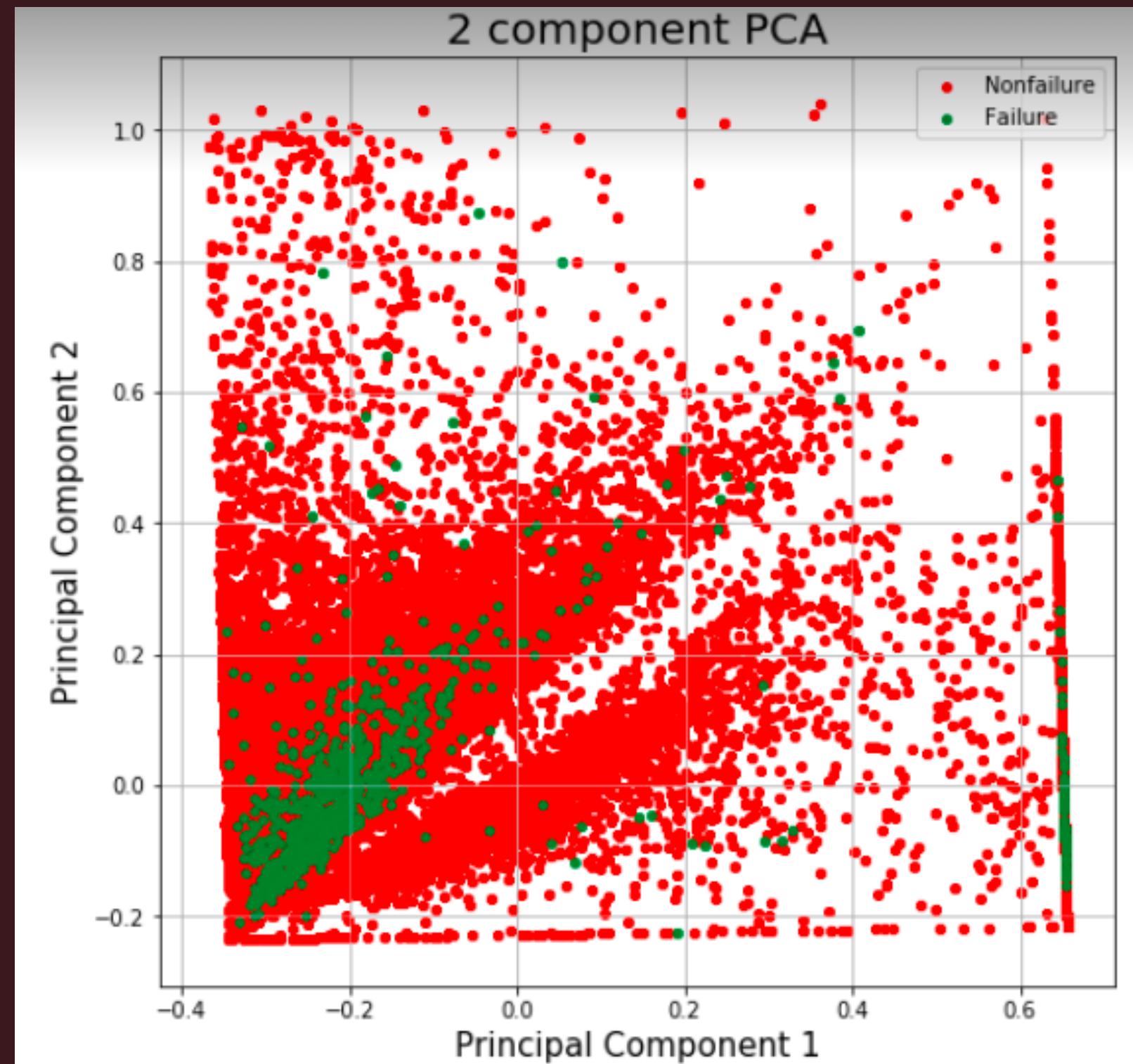
Each x-axis label is a single platform ID - we see that the average distance trucks travel range from around 90 to 600 - which shows that the variance in the distances of trucks are significant

# DURATION REGEN COMPARISON



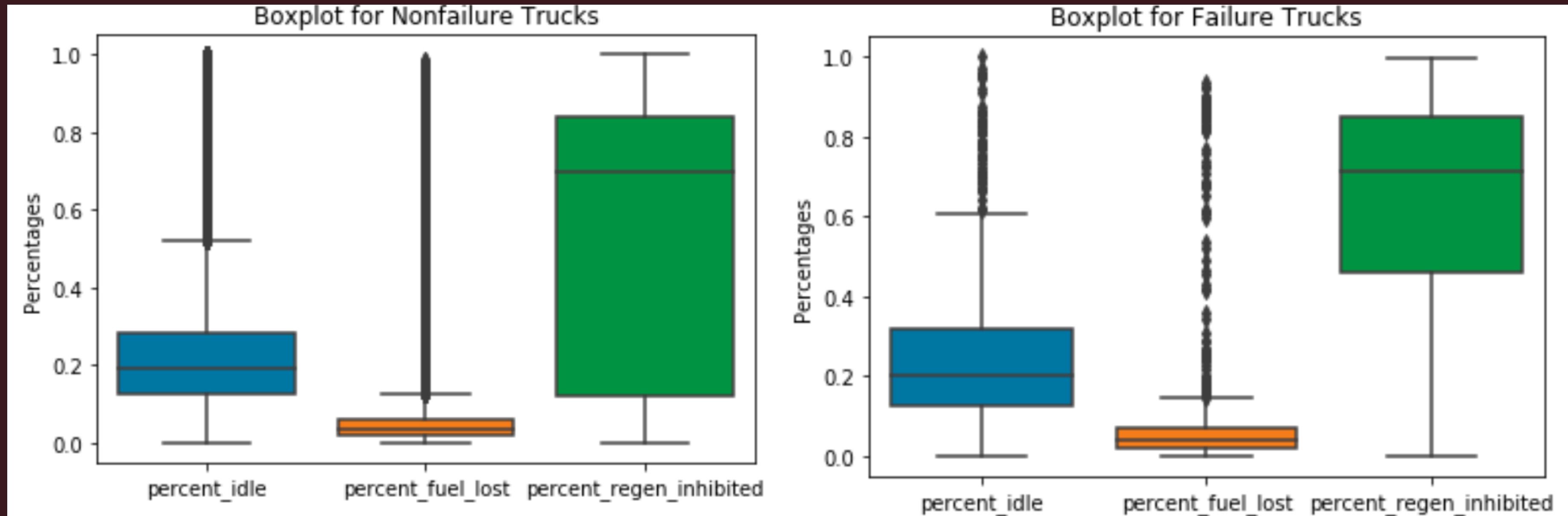
The nonfailure trucks has a stronger left skew suggesting the feature as a proxy for failure. We use this as one of our features we select later.

# 2 Component PCA



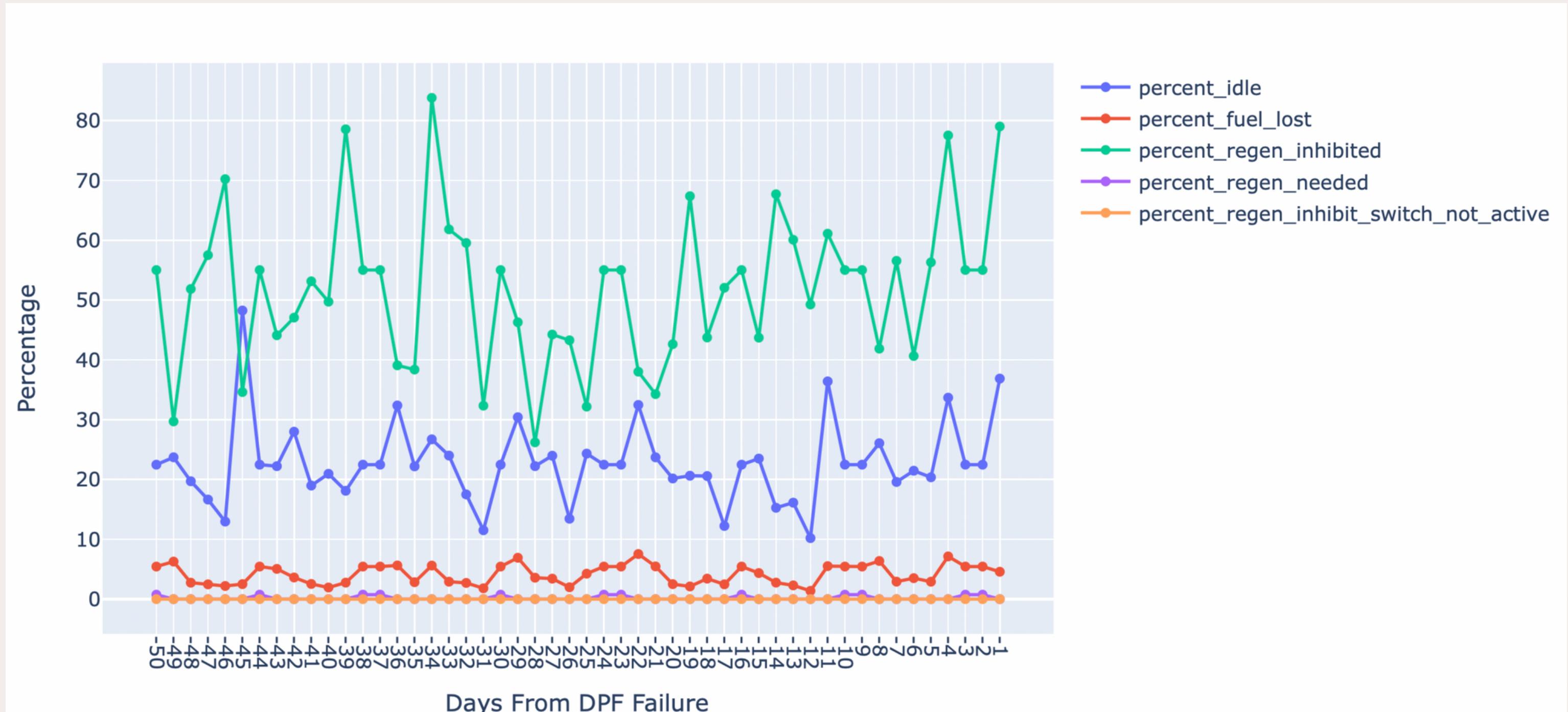
We attempted to run PCA on the original dataset for feature explainability - the results were not optimal

# BOXPLOT FOR FAILURE AND NONFAILURE



We see that failure trucks have a narrower inter quartile range for percent\_regen\_inhibited, but they mostly have similar distributions

# Key Characteristics 50 Days Before a Failure



Visualizing the data for Platform ID 300490 50 days before a failure

# EDA's Main Insights



Notable differences between  
failure and nonfailure  
trucks

Things can be explained from  
a relatively small feature set

We found clear candidates for  
features that we shouldn't use



A few features had no  
difference between  
nonfailure vs failure

Obvious class imbalance  
problem at hand

Some clear outliers for  
certain features



# Feature Creation

1



Before: Duration\_mins, mins\_idle

After: percent\_of\_time\_idle



Makes all entries on the same scale

2



Fuels Used, Miles Traveled

Fuels used per mile



Different way to aggregate the data



# Model Methods

1

Linear:

- Random Forest for Feature Importances
- PCA for Anomaly Detection

2

Nonlinear:

- Autoencoders
- TS Fresh Time Series Algorithms and Feature Transformations

# Narrowing Down to Two Features

Created a random forest model to determine the most relevant features to distinguish between the failure and nonfailure classes

Class imbalance problem can be fixed

15 days before service date = Failure

15-30 days after service date = Nonfailure



Leads to a balanced dataset



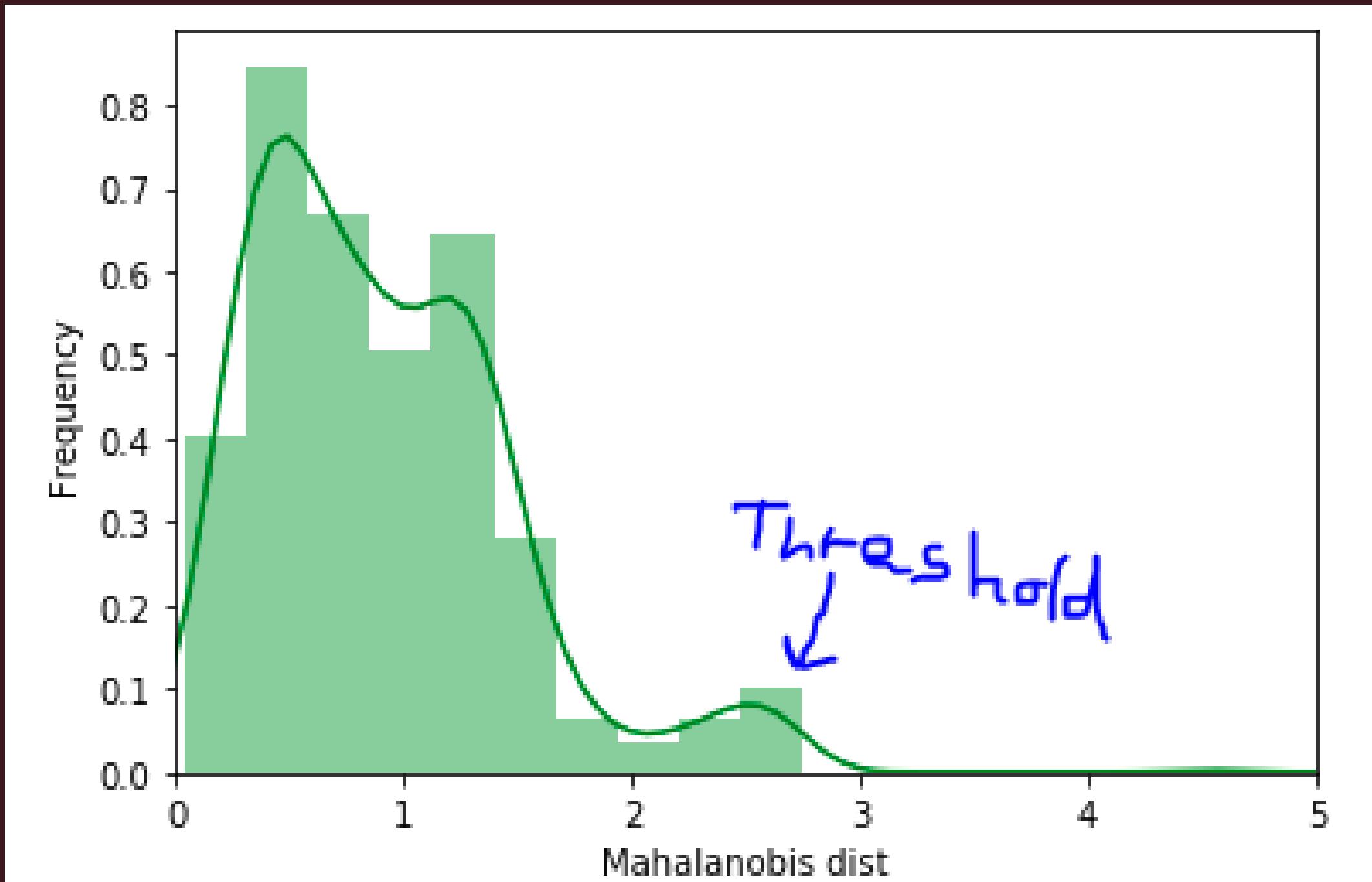
## Top 5 Important Features

- percent\_fuel\_lost
- fuels\_used\_per\_mile
- miles\_per\_minutes
- mins\_idle\_per\_mile
- percent\_regen\_inhibited

# PCA for Anomaly Detection

- 1** Reduce the original feature set to top 2 Principal Components
  - 2** Use only nonfailure trucks to calculate the covariance matrix of the normal class
  - 3** Calculate Mahalanobis distance for each of the failure trucks to the normal class
  - 4** If distance above the threshold, it would be marked as an anomaly
  - 5** Repeat the process with varying starting features
- 
- percent\_fuel\_lost, fuels\_used\_per\_mile, miles\_per\_minutes, mins\_idle\_per\_mile,  
percent\_regen\_inhibited

# Distribution of Mahalanobis Distances



The threshold is 3 standard deviations from the mean, so from this histogram, it is evident the threshold would be around 2.8

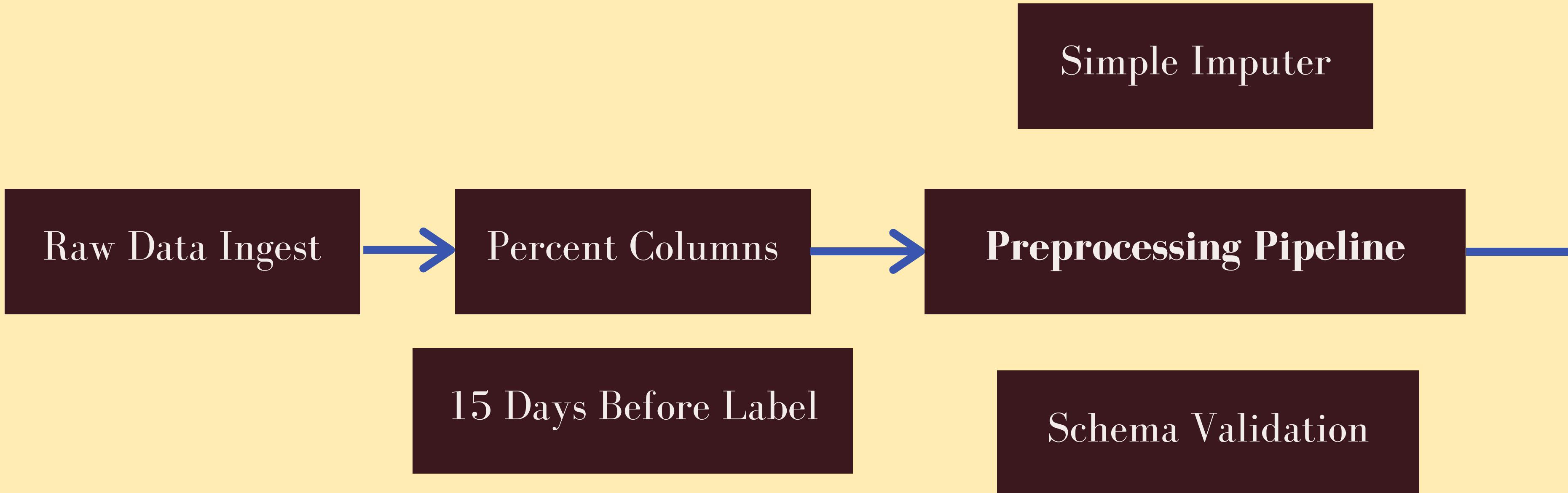
# Big Picture Pipeline



EDA\_wrangling



utilities.py



# Big Picture Pipeline



PCA

Visualizations

Random Forest



ts\_fresh



ts\_fresh



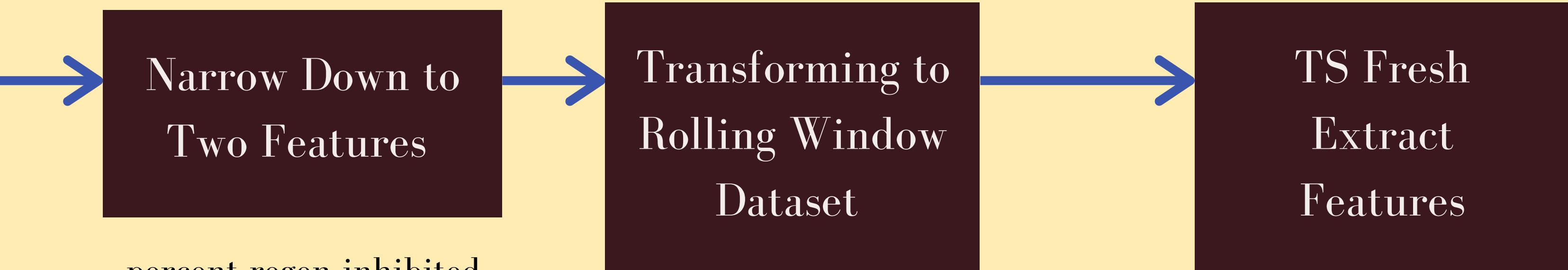
minimal FC parameters



efficient FC parameters



comprehensive FC parameters



percent\_regen\_inhibited

percent\_fuel\_lost



shift by 1  
window by 20



shift by 25  
window by 50



shift by 1  
window by 50

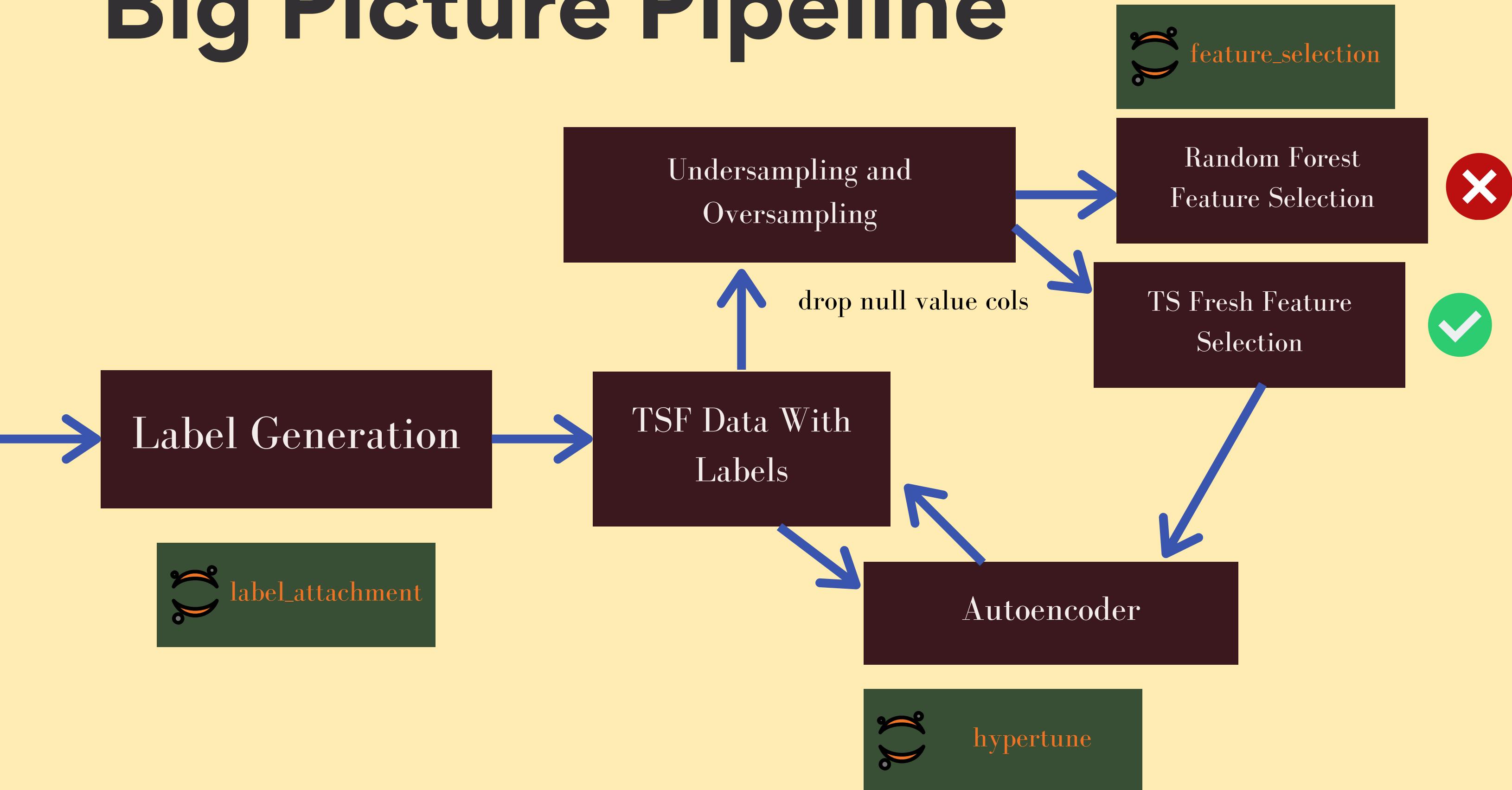


extract\_features on  
rolled window dataset



Hurdles: Runtime, parallelization,  
data formatting, data wrangling,  
seamless integration, etc.

# Big Picture Pipeline



# TS Fresh Transformations

9

Minimal

74

Comprehensive

72

Efficient

52

Medium

sum\_values  
median mean  
length  
standard\_deviation  
variance  
root\_mean\_sq  
uare  
maximum  
minimum

variance\_larger\_than\_standard\_deviation  
has\_duplicate\_max  
has\_duplicate\_min  
has\_duplicate  
sum\_values  
abs\_energy  
mean\_abs\_change  
mean\_change  
mean\_second\_derivative\_central  
absolute\_sum\_of\_changes

last\_location\_of\_maximum  
percentage\_of\_reoccurring\_values\_to\_all\_values  
percentage\_of\_reoccurring\_datapoints\_to\_all\_data\_points  
sum\_of\_reoccurring\_values  
sum\_of\_reoccurring\_data\_points  
ratio\_value\_number\_to\_time\_series\_length  
maximum minimum benford\_correlation  
time\_reversal\_asymmetry\_statistic c3 cid\_ce  
symmetry\_looking large\_standard\_deviation  
quantile autocorrelation agg\_autocorrelation  
partial\_autocorrelation number\_cwt\_peaks

# TS Fresh Transformations

9

# Minimal

74

# Comprehensive

72

# Efficient



## Too simplistic



# Runtime too high

# Memory timeout



# Runtime around 2h

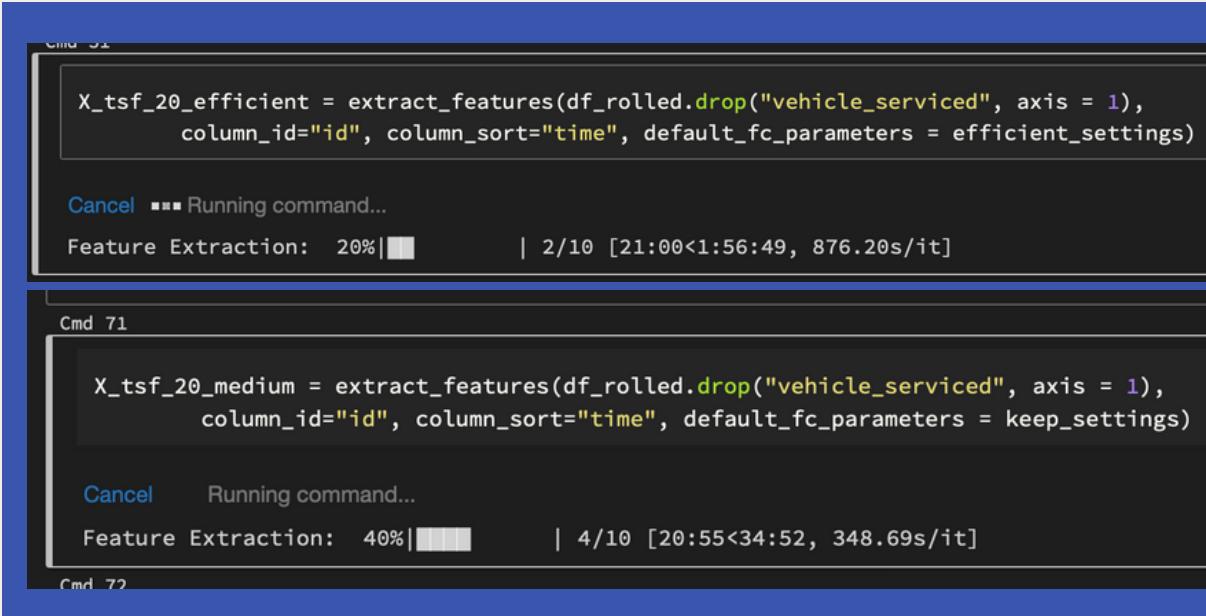
## Memory timeout



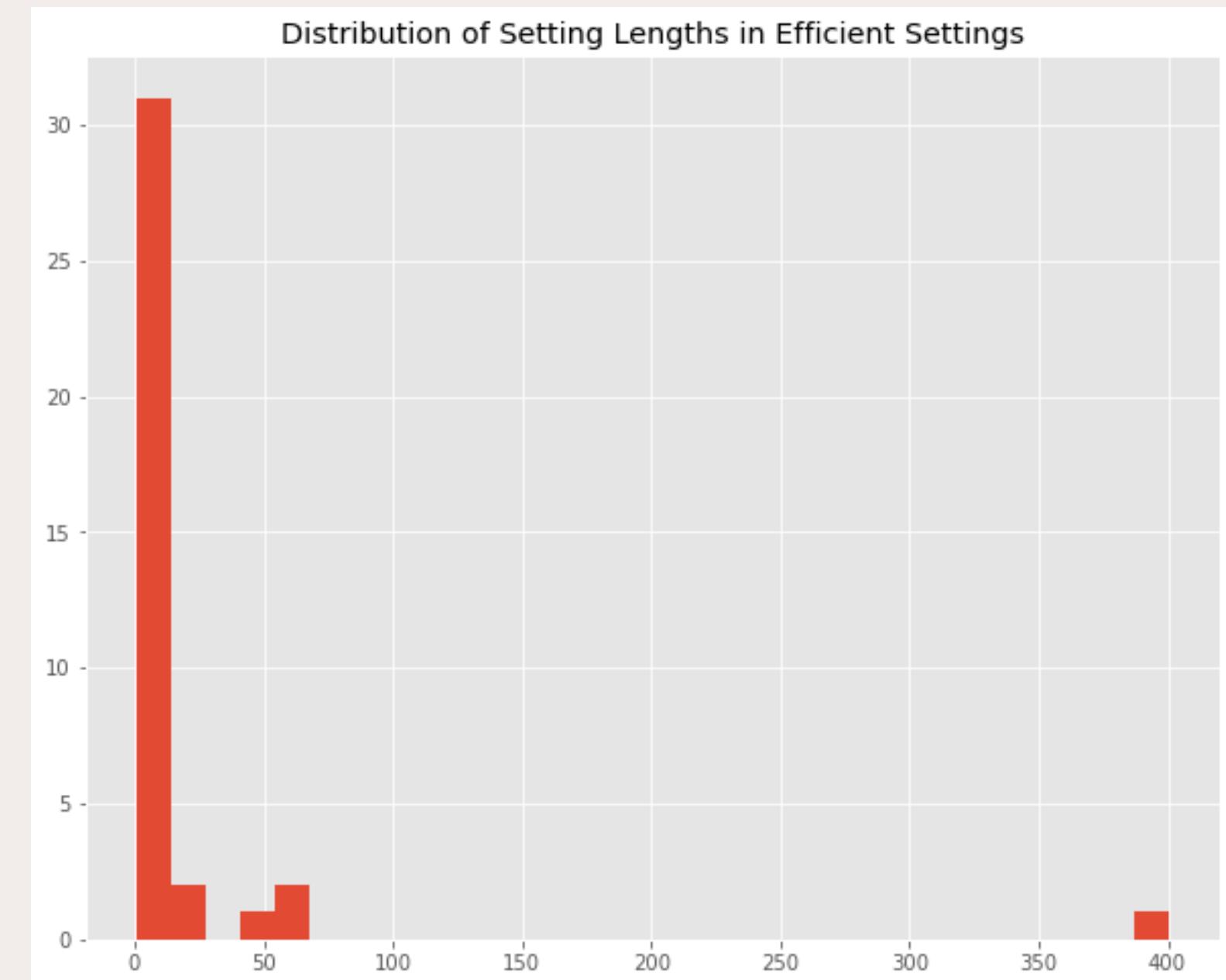
Note: Cluster is capped  
at 30GB memory

52

# Medium



Removed all instances where the transformation parameter setting dictionary value (sub transformations) exceeds 5



Tried to parallelize with Dask or  
make code Spark aware

# Class Imbalance

**Initial Dataset Distribution:**  
100,000 non failures  
1380 failures

After Transformation



**0.0909%**

of examples are a DPF failure

**97887:89**

Undersampling

Synthetic Minority  
Oversampling  
Technique (SMOTE)

**33%**

of examples are a  
DPF failure

**50%**

of examples are a  
DPF failure



imblearn

# Runs and Parameter Configurations

## Class Balance Data

Balanced

Unbalanced

## TSF Window Size and Overlap

w: 50  
s: 25

w: 20  
s: 20

w: 20  
s: 1

## RF

n\_estimators

## TS Fresh Transformations

Minimal

Comprehensive

Efficient

Medium

## Autoencoder Parameters

Autoencoder Window Size

Early Stopping

Learning Rate

Dropout

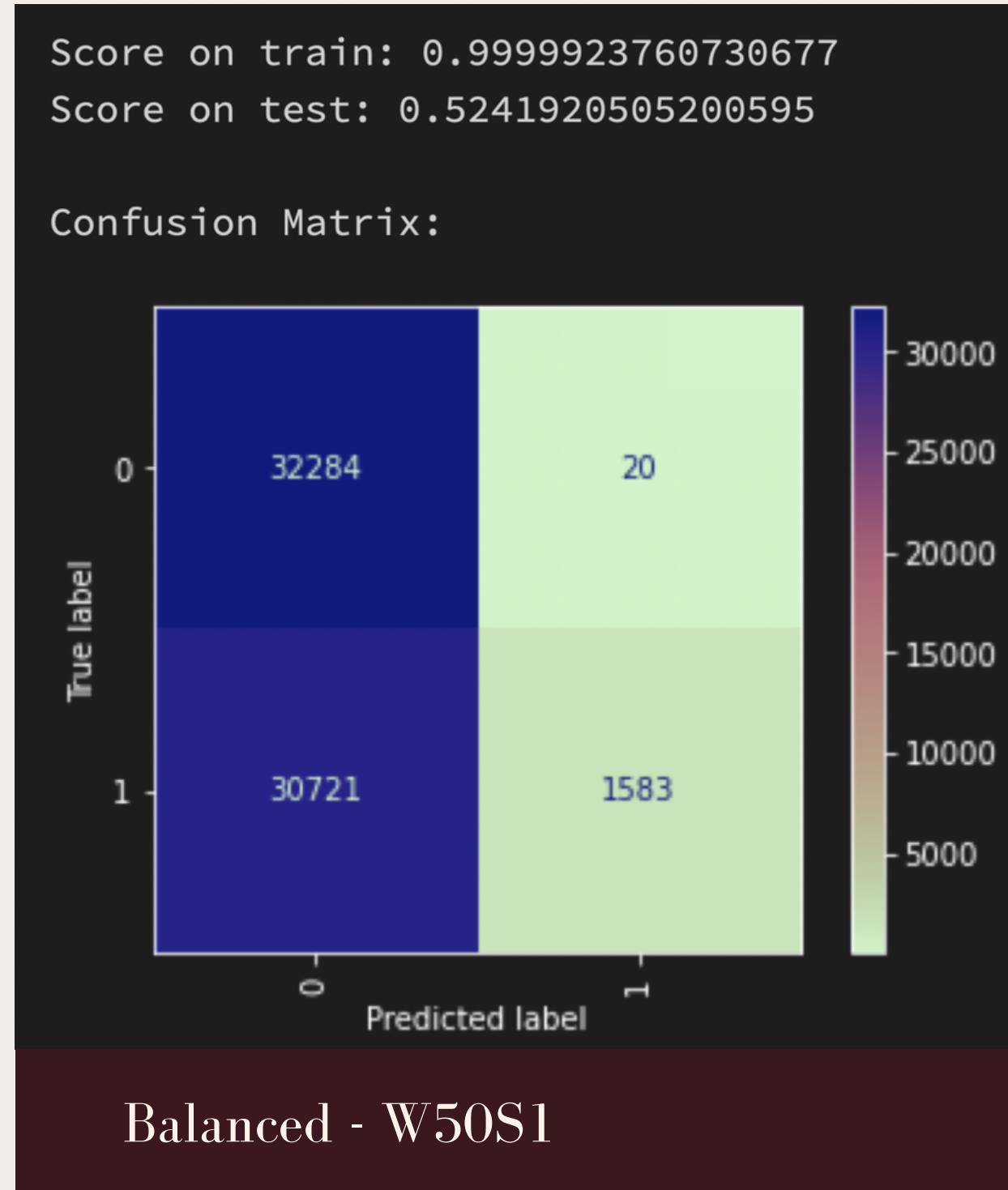
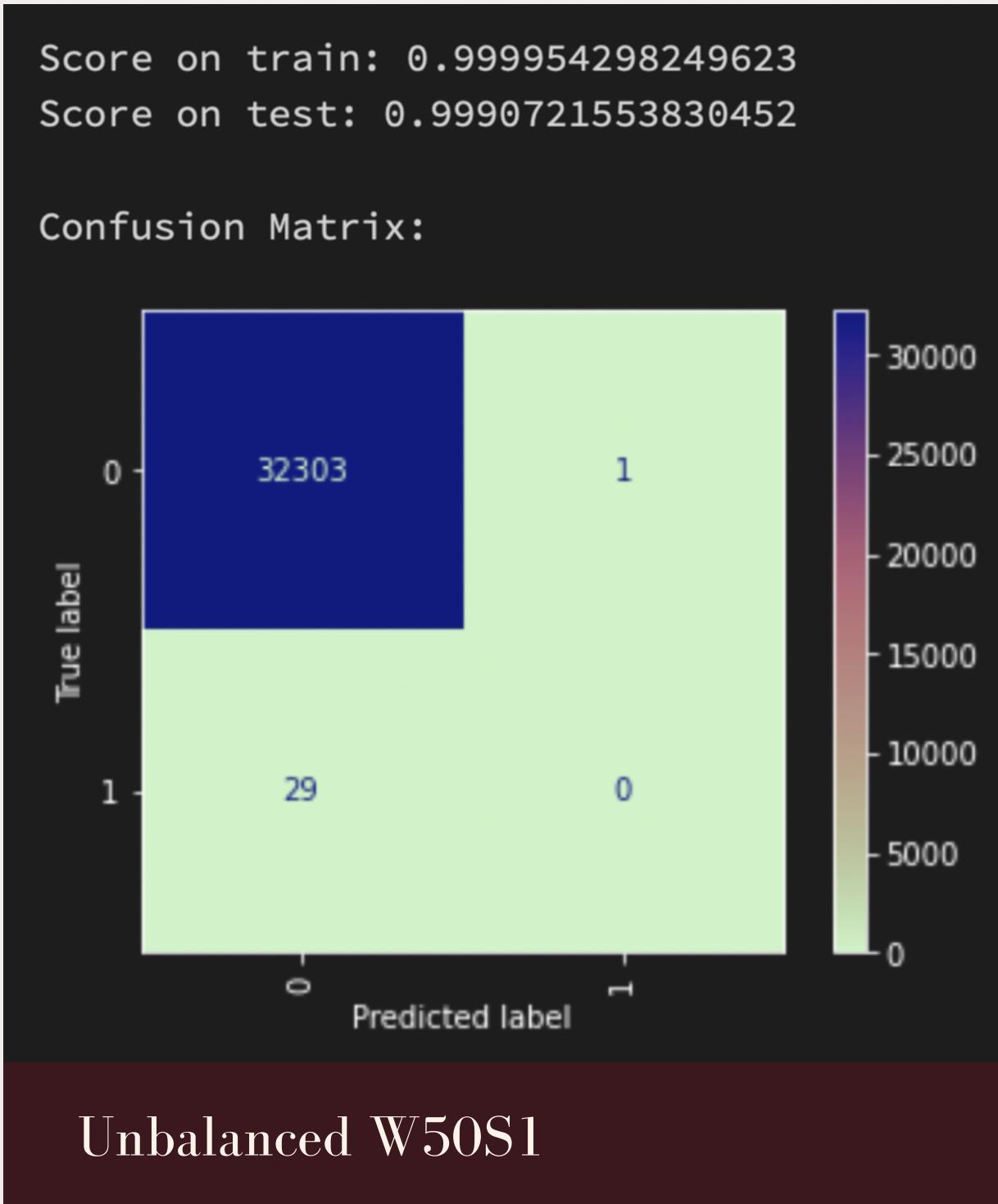
## Original Data

Features to Keep

Normalization

Imputation

# Random Forest Feature Selection



	Feature	Feature Importance
8	F1_minimum	0.083907
0	F1_sum_values	0.081044
16	F2_maximum	0.077756
2	F1_mean	0.076788
1	F1_median	0.075476
6	F1_root_mean_square	0.073237
17	F2_minimum	0.070889
10	F2_median	0.069379
11	F2_mean	0.057002
9	F2_sum_values	0.055843
7	F1_maximum	0.054174
4	F1_standard_deviation	0.053665
5	F1_variance	0.053024
14	F2_variance	0.041555
15	F2_root_mean_square	0.039783
13	F2_standard_deviation	0.036481
12	F2_length	0.000000
3	F1_length	0.000000

Balanced - W50S1

# TS Fresh Feature Selection Results

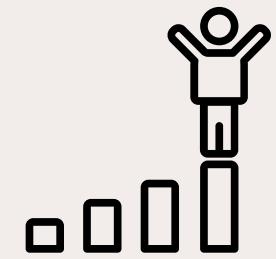
SMOTE  
Oversampled Data

tsfresh.feature\_selection  
.calculate\_relevance\_table  
  
drop features with NA

Undersampled Data



feature	type	p_value	relevant	rank
F1_has_duplicate_max	real	0.000000	True	1
F1_last_location_of_maximum	real	0.000000	True	2
F1_median	real	0.000000	True	3
F1_longest_strike_above_mean	real	0.000000	True	4
F1_root_mean_square	real	0.000000	True	5
...	...	...	...	...
F1_value_count_value_1	binary	0.000122	True	101
F1_range_count_max_1_min_-1	binary	0.000122	True	102
F2_value_count_value_1	binary	0.000244	True	103
F2_range_count_max_1_min_-1	binary	0.000244	True	104
F2_time_reversal_asymmetry_statistic_lag_3	real	0.000375	True	105
105 rows x 5 columns				



top 33  
features had  
a p\_value of  
0!

feed top 20 to  
autoencoder

## Benjamini Hochberg procedure

H0: The feature is not relevant and should not be added  
H1: The feature is relevant and should be kept



feature	type	p_value	relevant
F1_spkt_welch_density_coeff_2	real	0.057590	False
F1_ratio_value_number_to_time_series_length	real	0.058556	False
F1_variation_coefficient	real	0.060909	False
F1_percentage_of_reoccurring_datapoints_to_all_datapoints	real	0.062433	False
F1_percentage_of_reoccurring_values_to_all_values	real	0.064571	False
...	...	...	...
F2_range_count_max_0_min_1000000000000.0	constant	NaN	False
F2_range_count_max_1000000000000.0_min_0	constant	NaN	False
F2_number_crossing_m_m_-1	constant	NaN	False
F2_number_crossing_m_m_1	constant	NaN	False
F2_count_above_t_0	constant	NaN	False
132 rows x 4 columns			

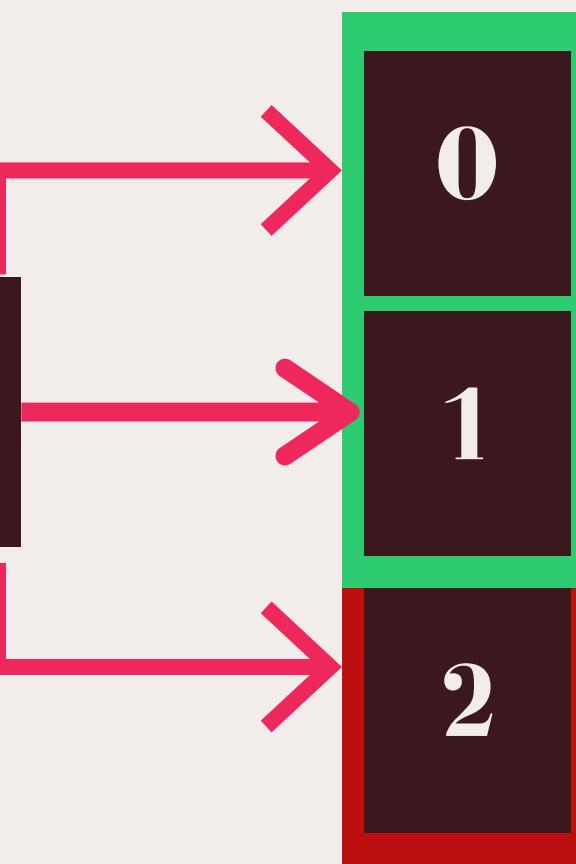
# Label Generation

has window\_end column

TSF Data (with  
extracted  
features)

shape: (111707,9)

Original Data



next row : no dpf failure

window : no dpf failure

next row : has dpf failure

window : dpf failure

any part of window has dpf  
failure - we don't want this!

TSF Data With  
Labels

Window Level Ground  
Truth Labels

combine\_label

# Autoencoder

Dense Layers

Convolutional  
Autoencoder

LSTM  
Autoencoder

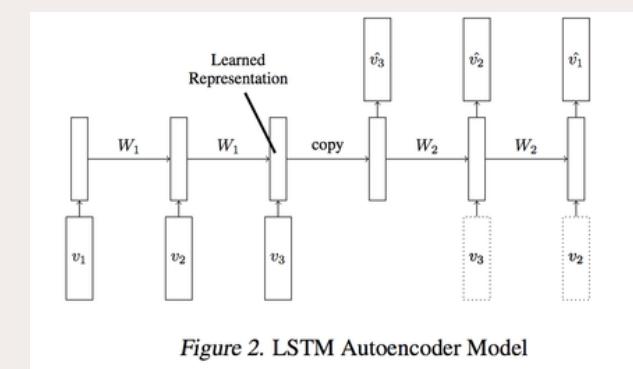
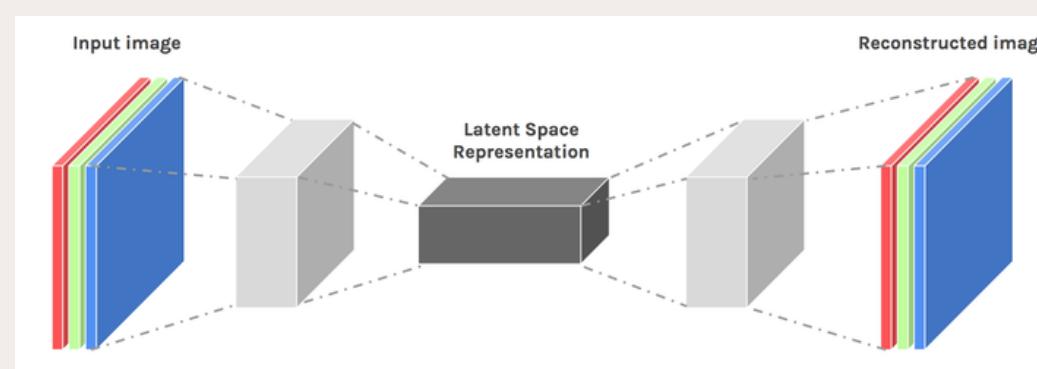
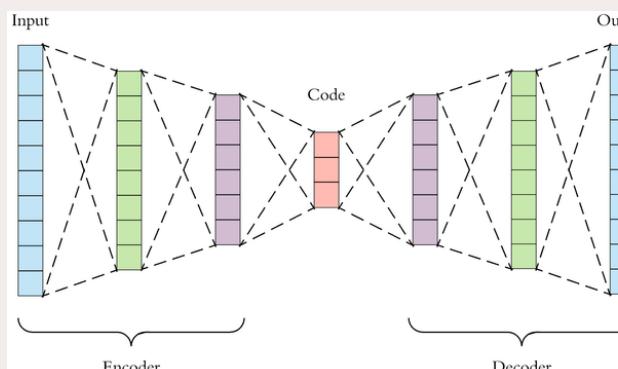
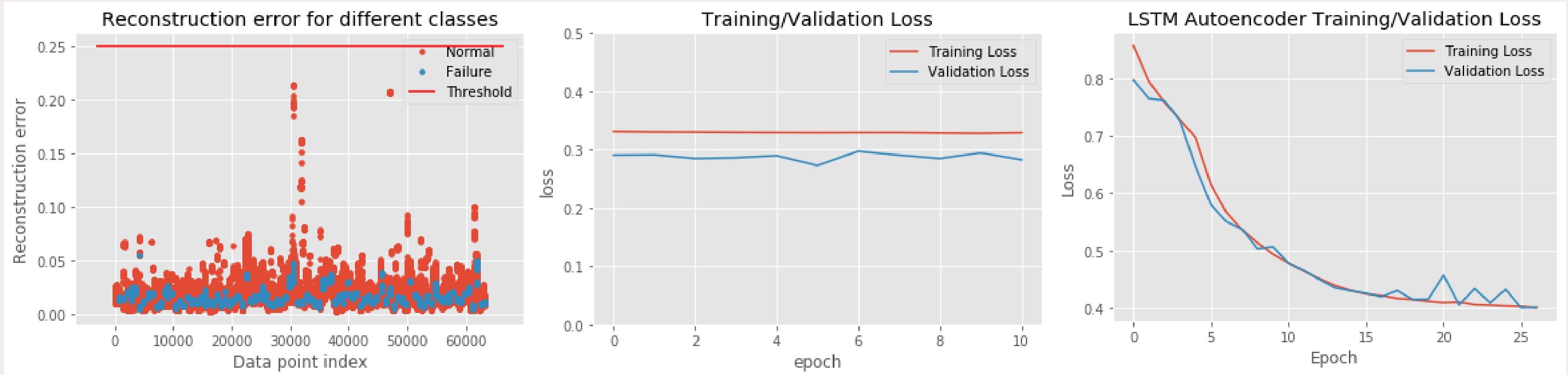
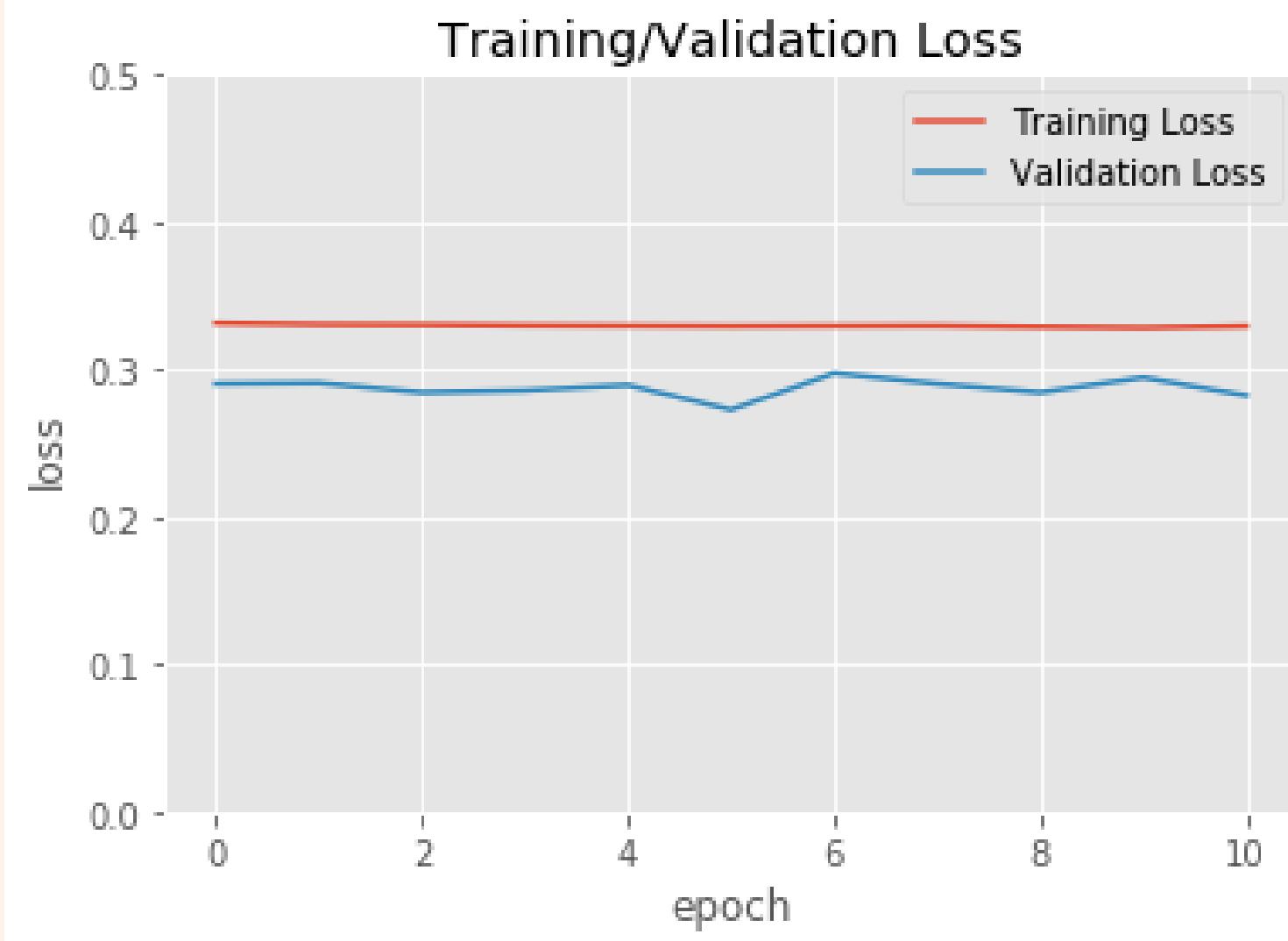


Figure 2. LSTM Autoencoder Model

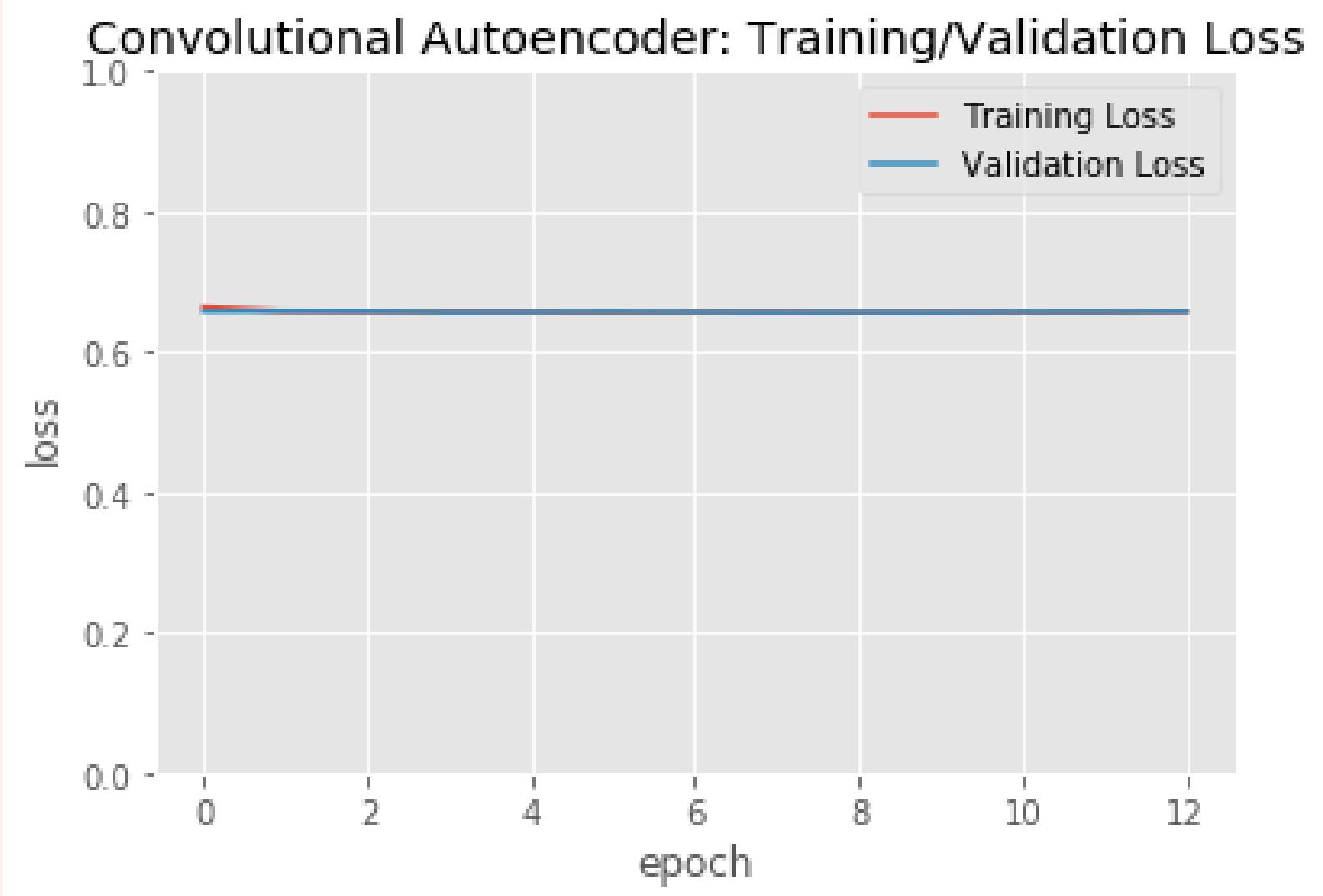
# LOSS

Convolutional Autoencoder:

Before



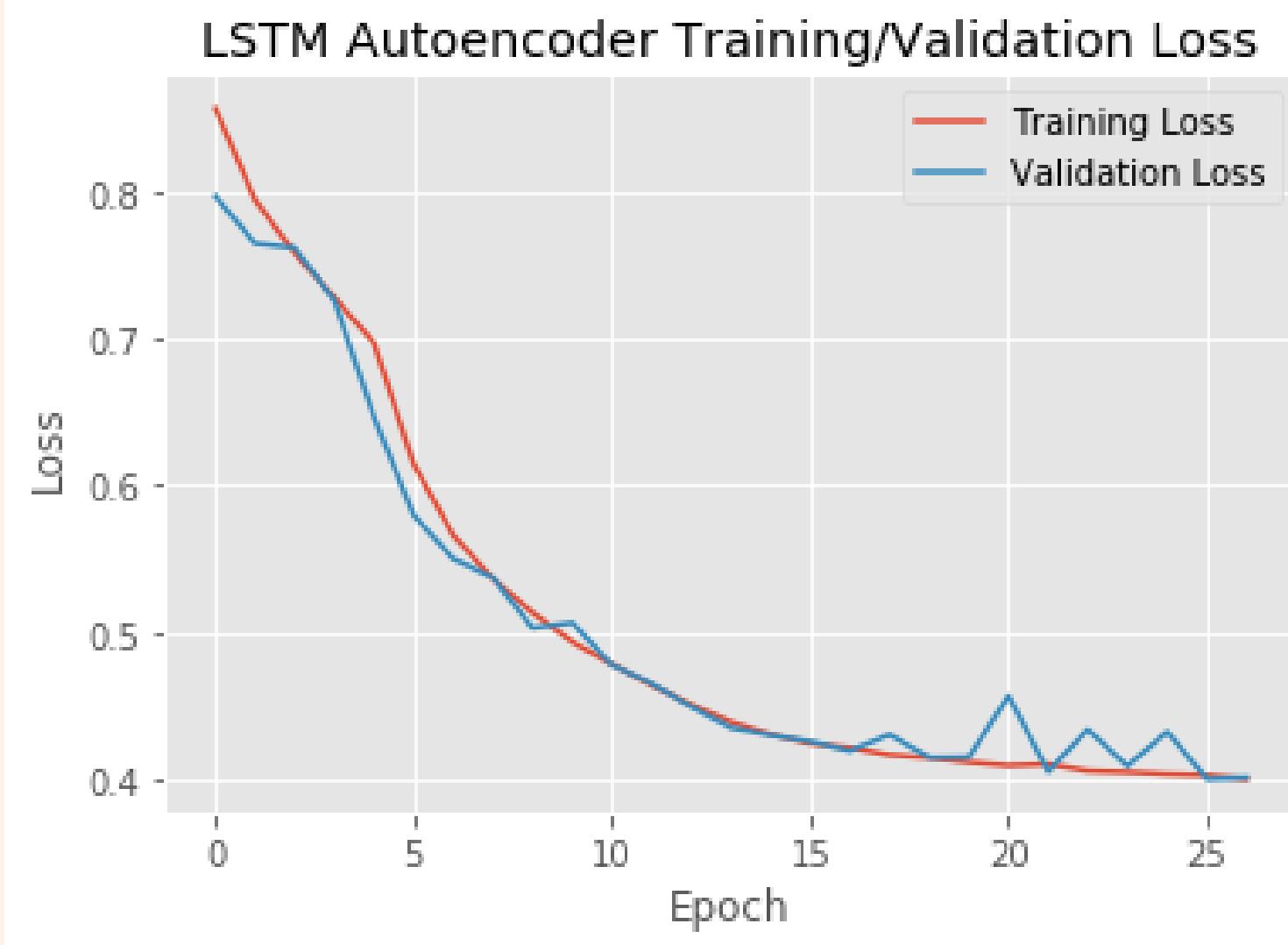
After



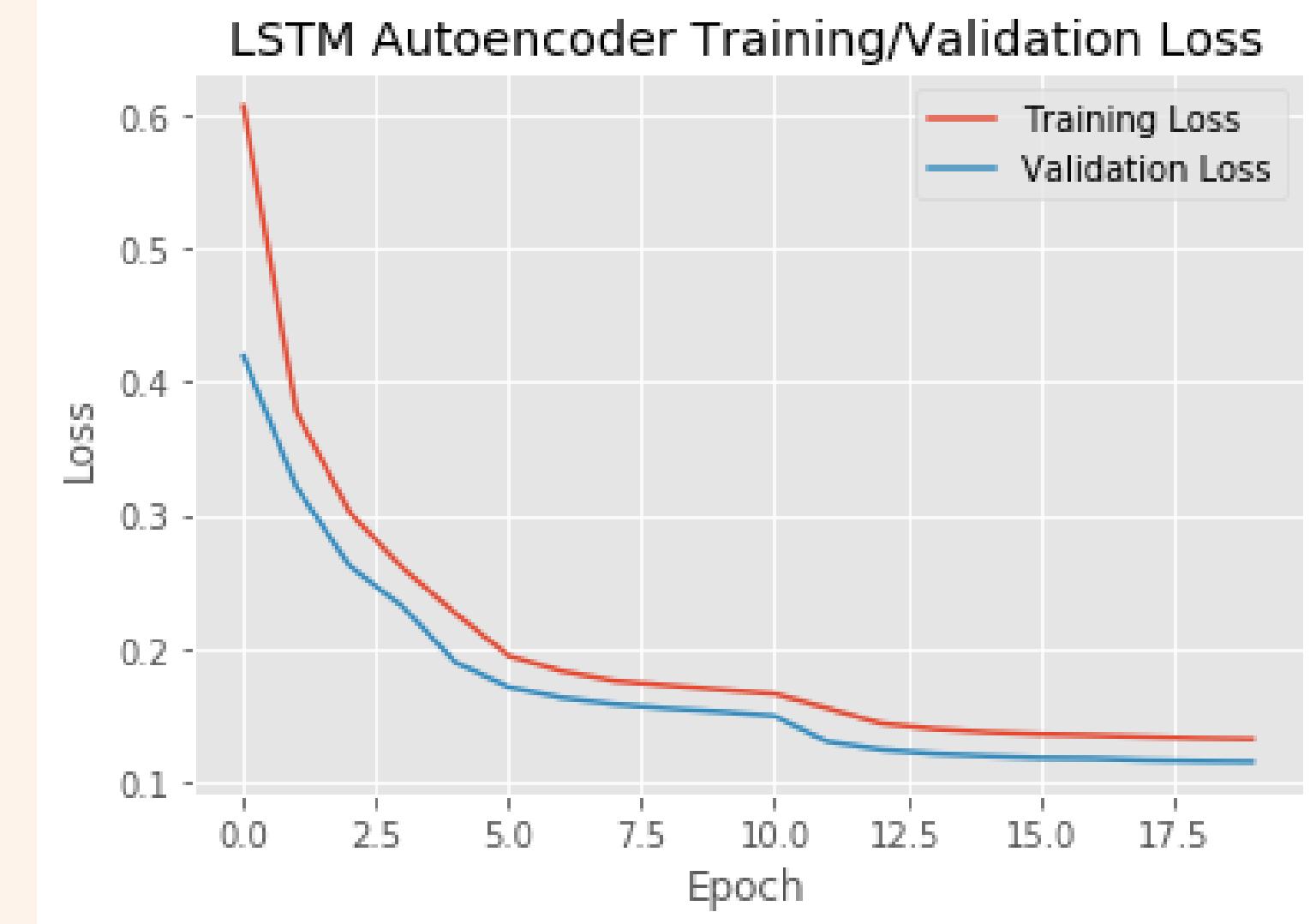
# LOSS

LSTM Autoencoder:

Before



After



# MLflow

hyperparameter to "manually" tune

Select relevant features

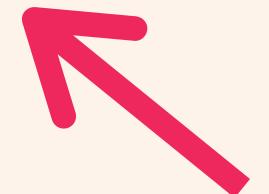
Select parameters to tune with

automated tuning with ML ops

Generate the best hyperparameters from MLflow

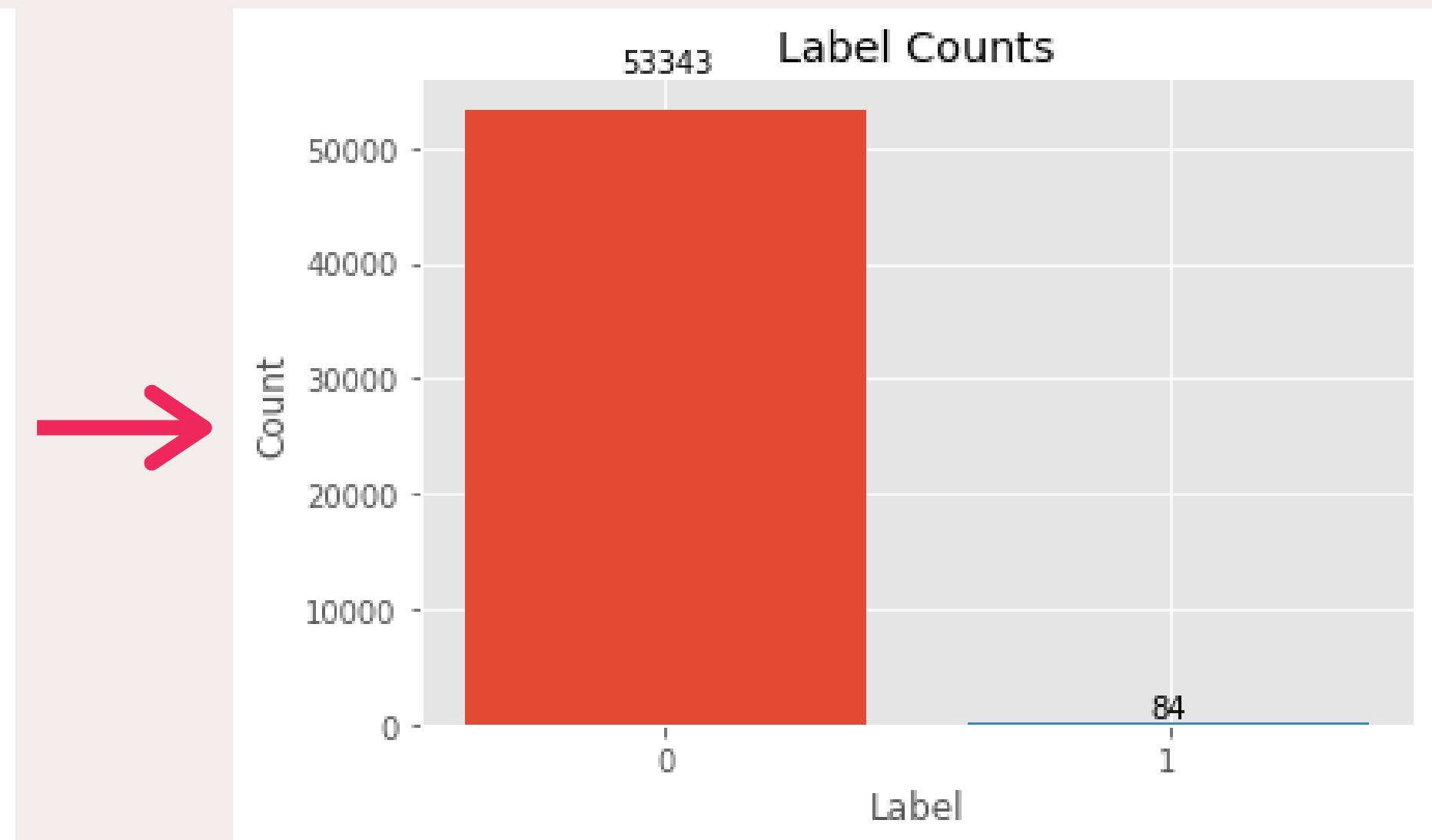
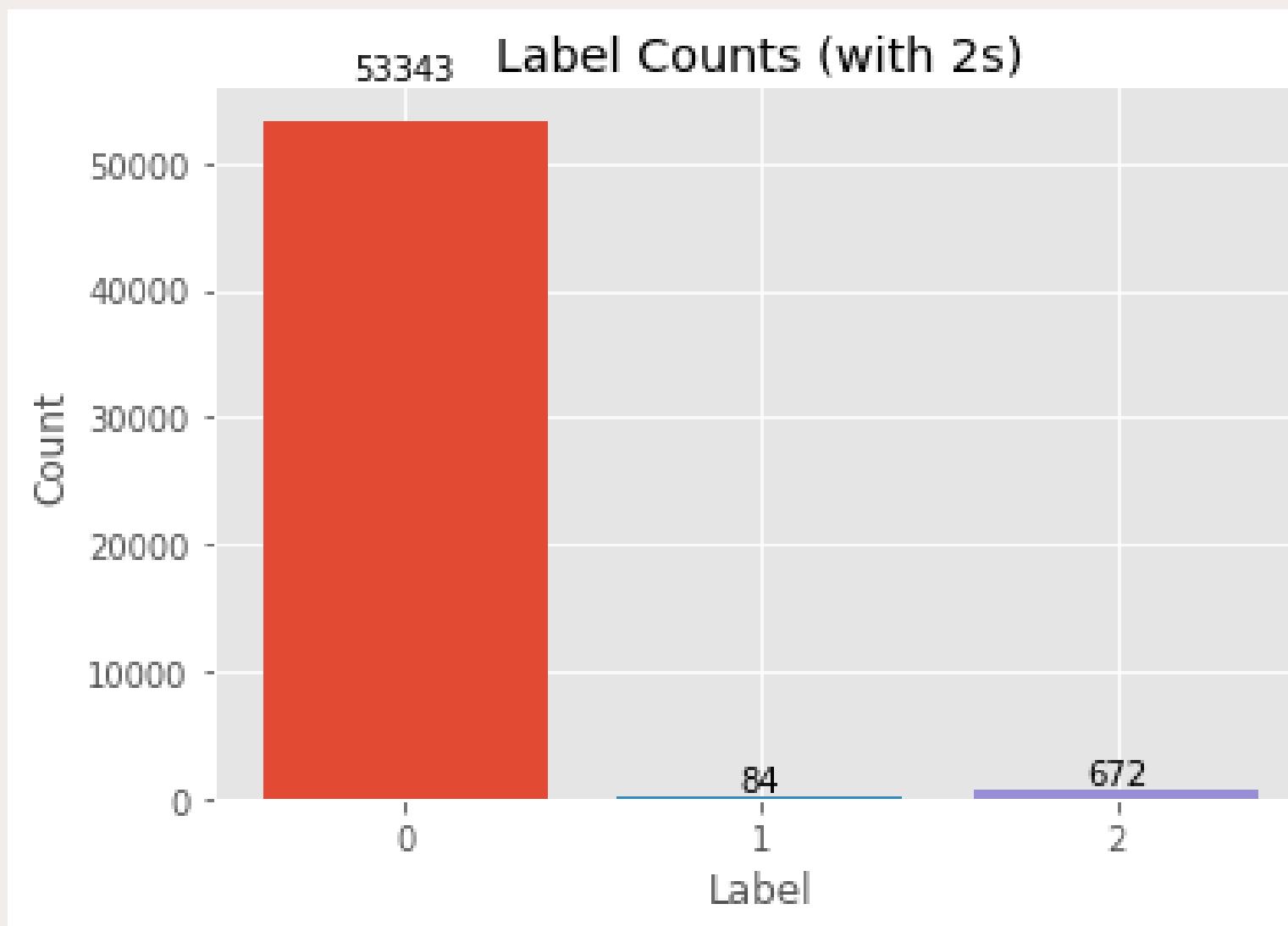
Generate labels

Examine the model on testing data



# Autoencoder Labels

with a window size of 8



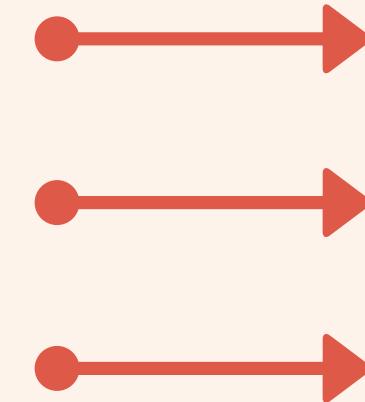
# MLflow

	RUN 1	RUN 2
<b>TS Fresh</b>	w: 50	w: 20
<b>Window Size</b>	s: 1	s: 1
<b>Autoencoder</b>	w': 8	w': 10-30
<b>Window Size</b>		
<b>Total size</b>	58	30-50

# MLflow

## Parameters to tune with:

```
LSTM layer1 units  
LSTM layer2 units  
epochs  
learning_rate  
dropout rate  
optimizer("Adadelta", "Adam")  
window_size
```



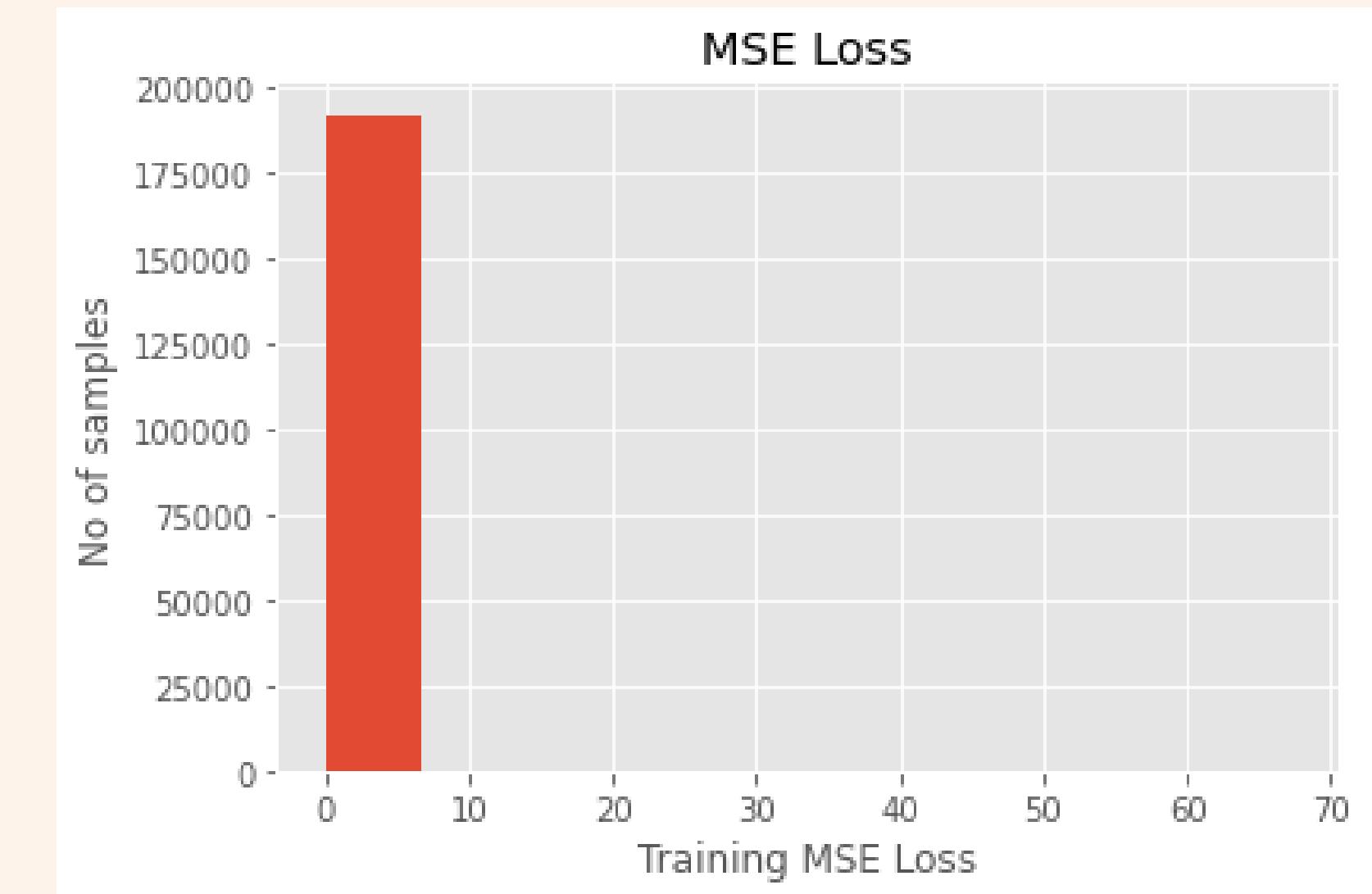
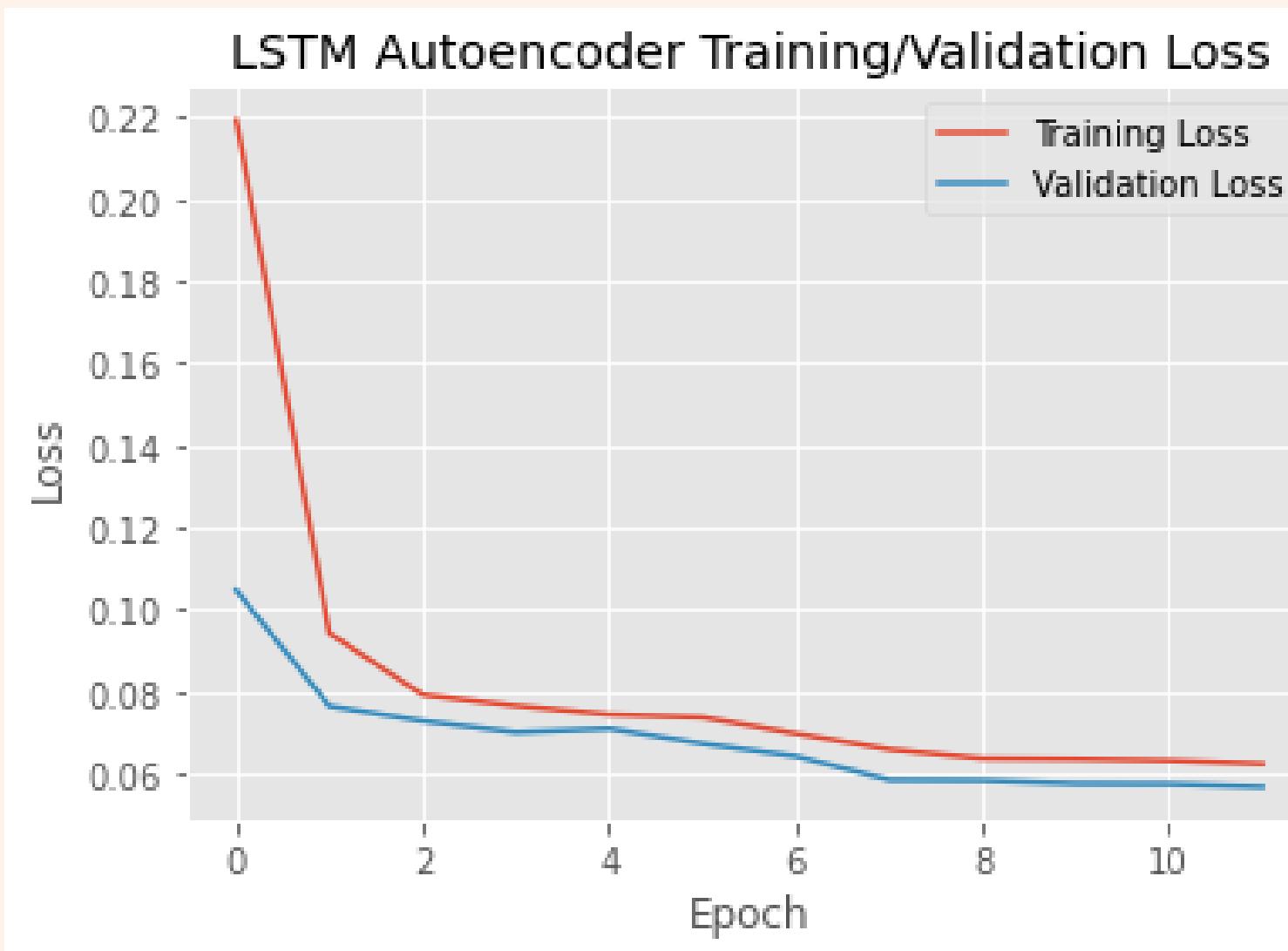
## Example output

```
'dropout': 0.02243656795453065  
'epochs': 96.0  
'learning_rate': 0.00681020946819995  
'lstm_l1': 14.0  
'lstm_l2': 8.0  
'optimizer': 'Adam'  
'window_size': 10
```

# MLflow *results*

With w=50+8

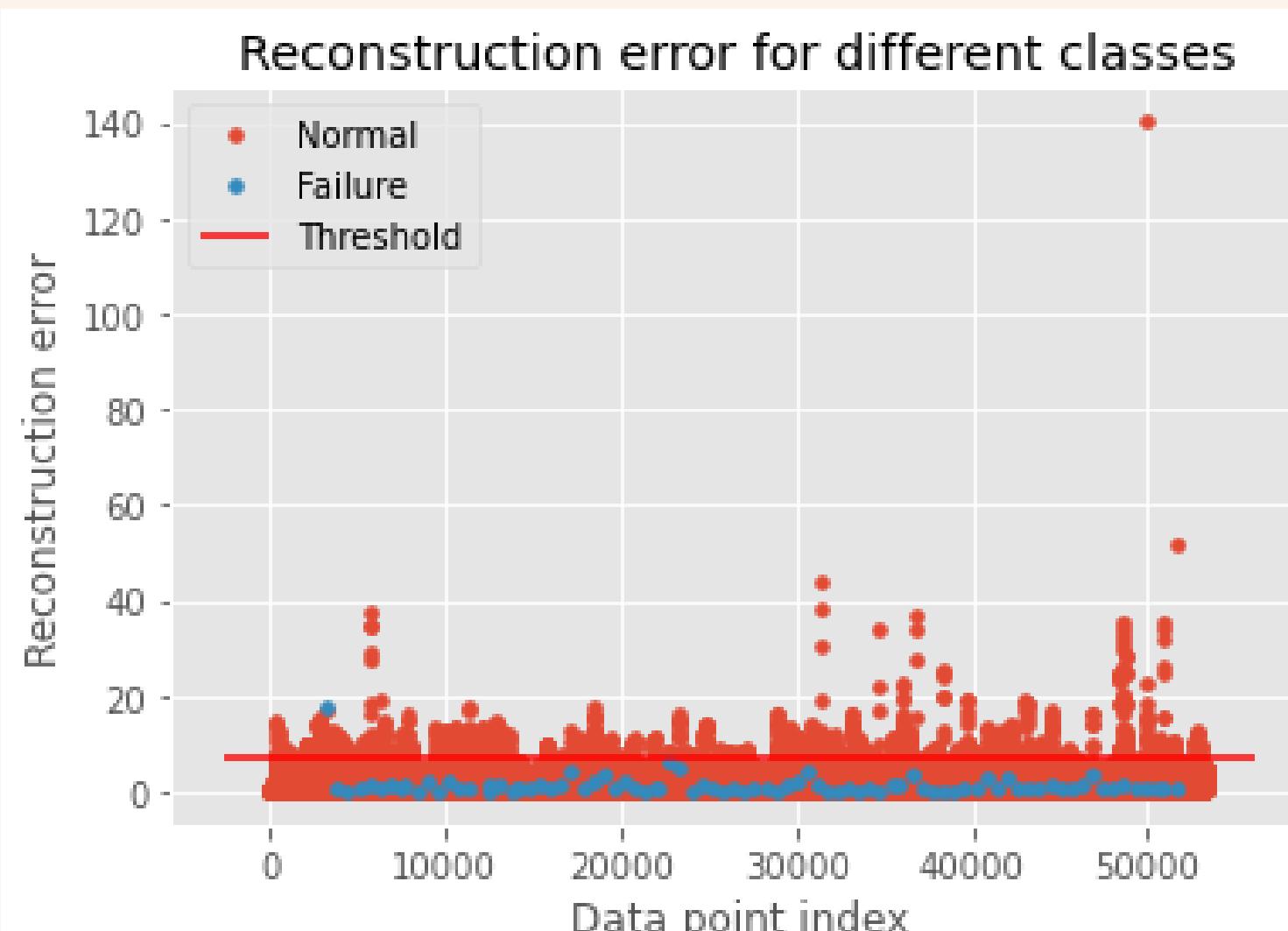
```
'dropout': 0.02243656795453065,  
'epochs': 96.0  
'learning_rate': 0.00681020946819995,  
'lstm_11': 14.0  
'lstm_12': 8.0  
'optimizer': 'Adam'
```



RUN 1

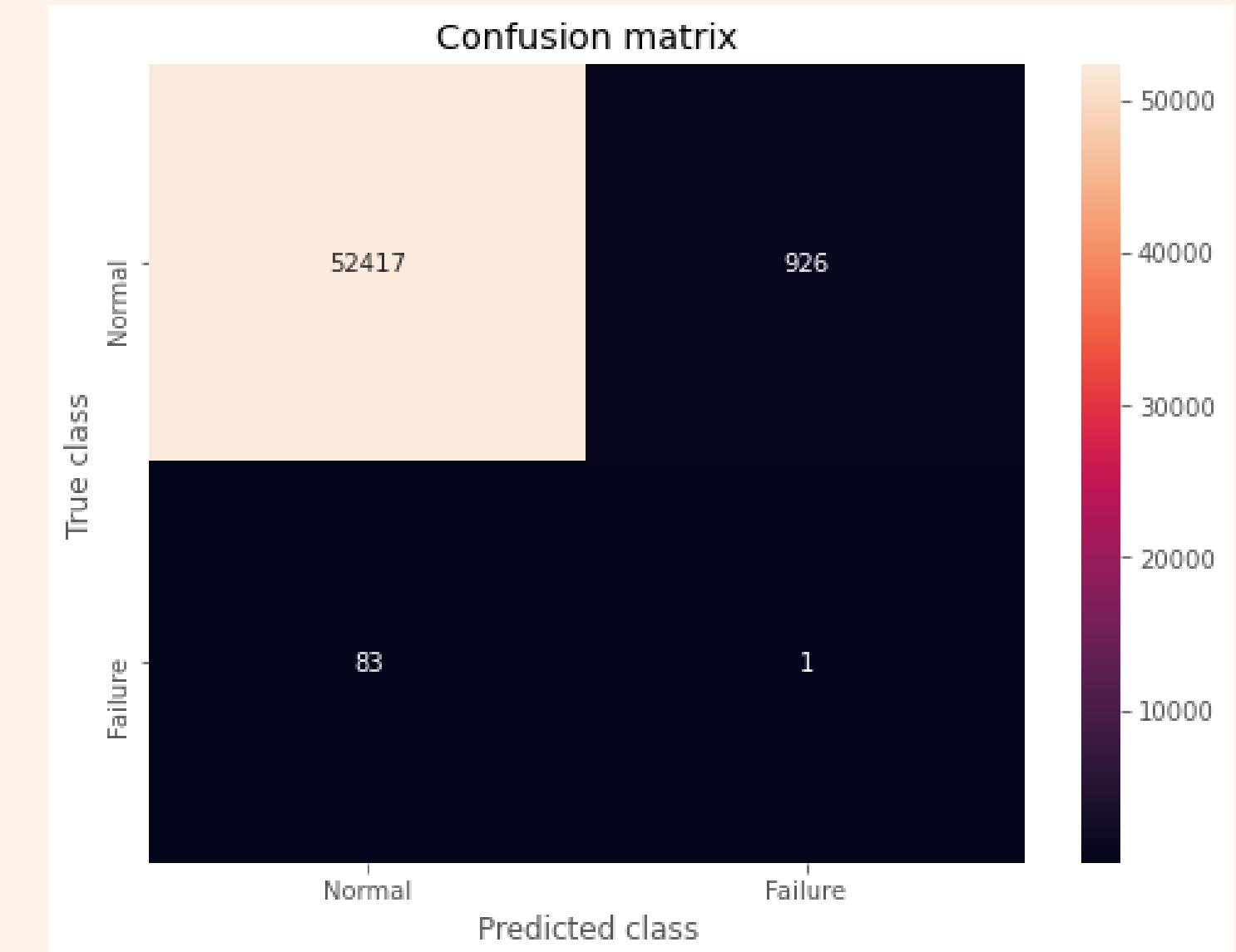
# MLflow *results*

With  $w=50+8$



The data is clustered together - bad sign

'dropout': 0.02243656795453065  
'epochs': 96.0  
'learning\_rate': 0.00681020946819995  
'lstm\_11': 14.0  
'lstm\_12': 8.0  
'optimizer': 'Adam'

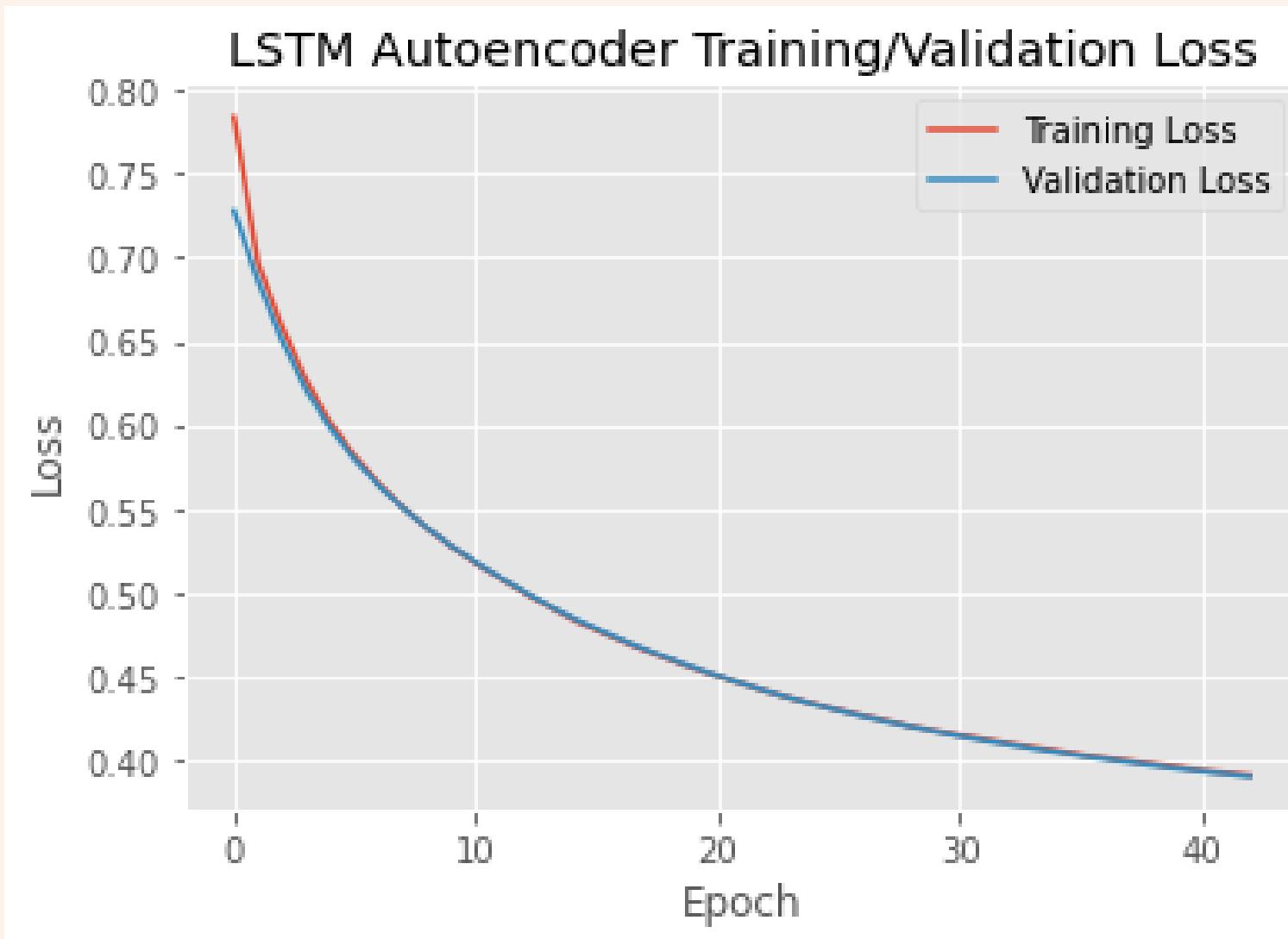


RUN 1

# MLflow *results*

With  $w=20+10 \sim 30$

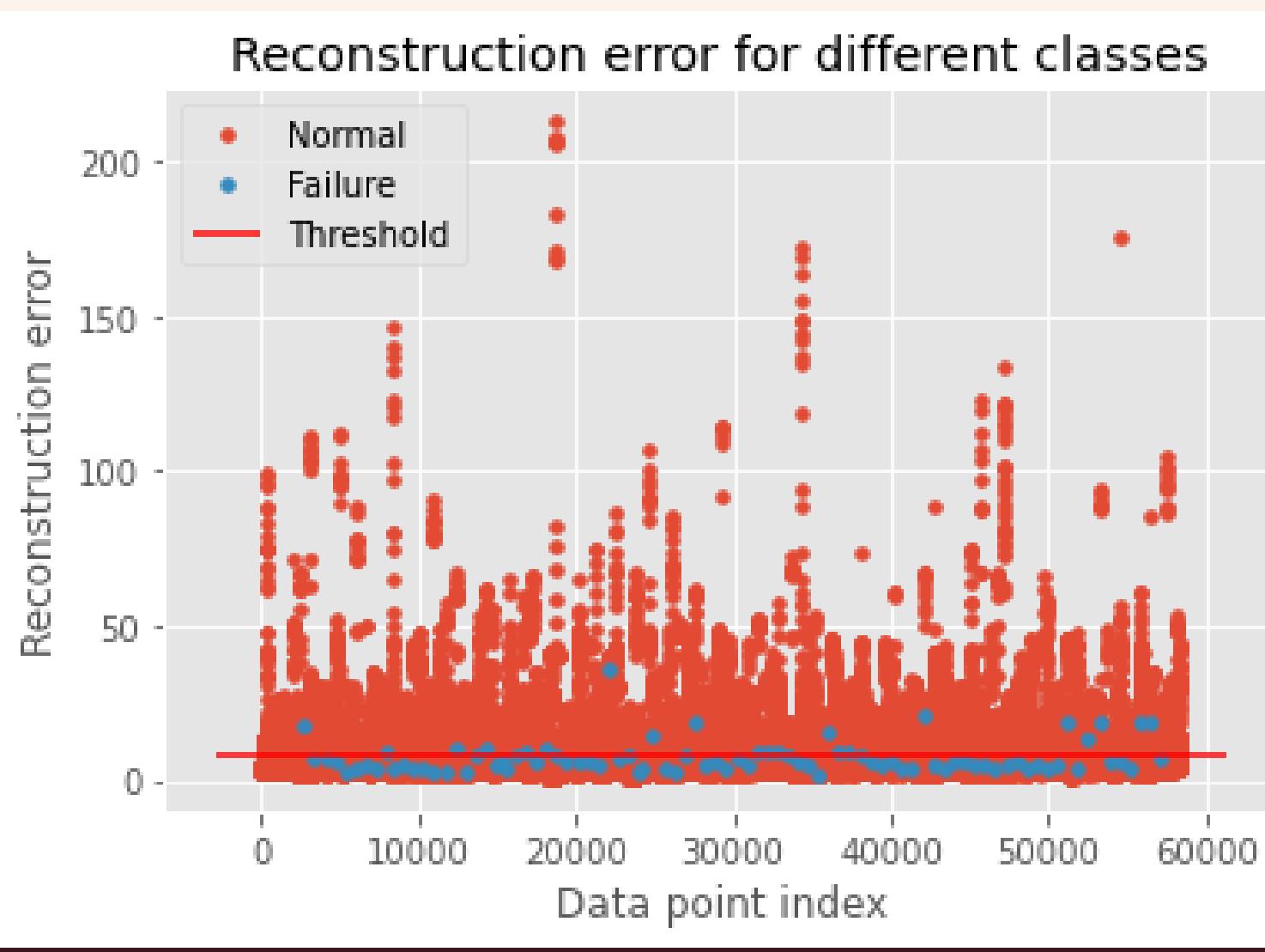
```
'dropout': 0.07031539323259835
'epochs': 55.0
'learning_rate': 0.08271996530819141
'lstm_11': 32.0
'lstm_12': 10.0
'optimizer': 'Adadelta'
>window_size': 10.0
```



RUN 2

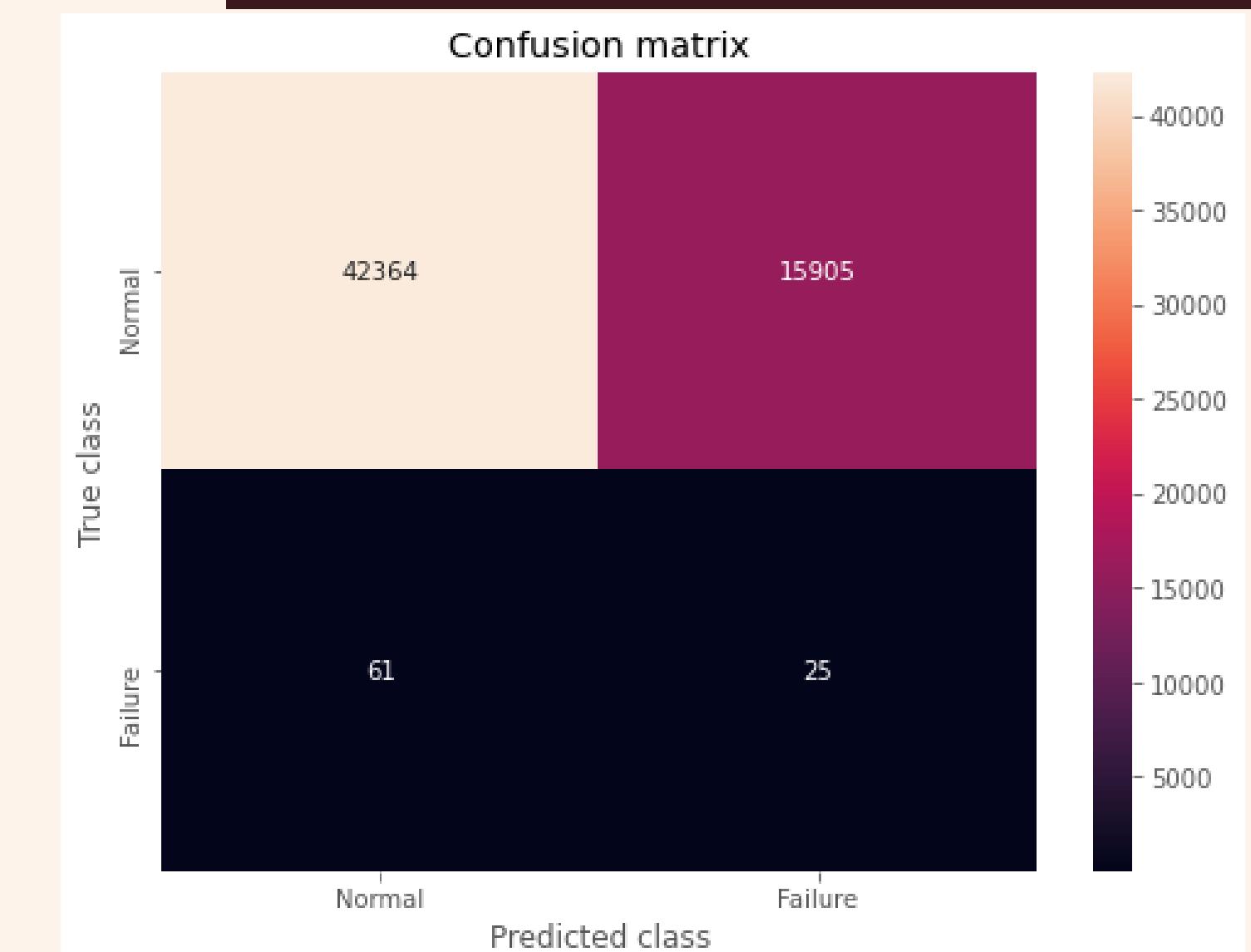
# MLflow results

With  $w=20+10 \sim 30$



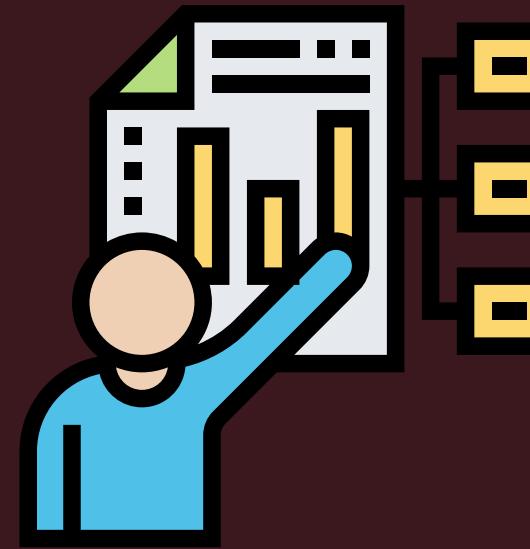
The data is clustered together - bad sign

```
'dropout': 0.07031539323259835
'epochs': 55.0
'learning_rate': 0.08271996530819141
'lstm_11': 32.0
'lstm_12': 10.0
'optimizer': 'Adadelta'
>window_size': 10.0
```



RUN 2

# *Results and Conclusion*



Lots of improvements with LSTM Autoencoder

Corrected TS Fresh Selection Method and Pipeline

Finished the complete pipeline (tuning and evaluation)



Non-failure and failure trucks MSE still clustered together

**Current Best  
Recall Scores:**

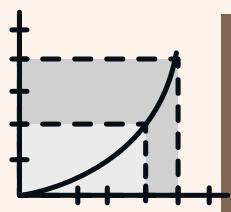
**50% and 56%**

*for DPF nonfailure      for DPF failure*

# Summary of Challenges



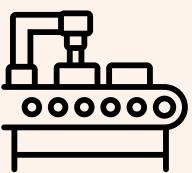
## General



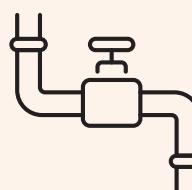
Learning Curve



Containerization



Code Productionalization



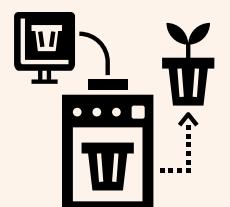
Pipeline has many moving parts



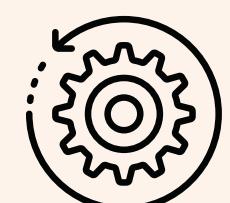
New Approaches Required



Local and cluster memory -  
runtime issues



Reproducibility



No Leaks in Pipeline



tsfresh

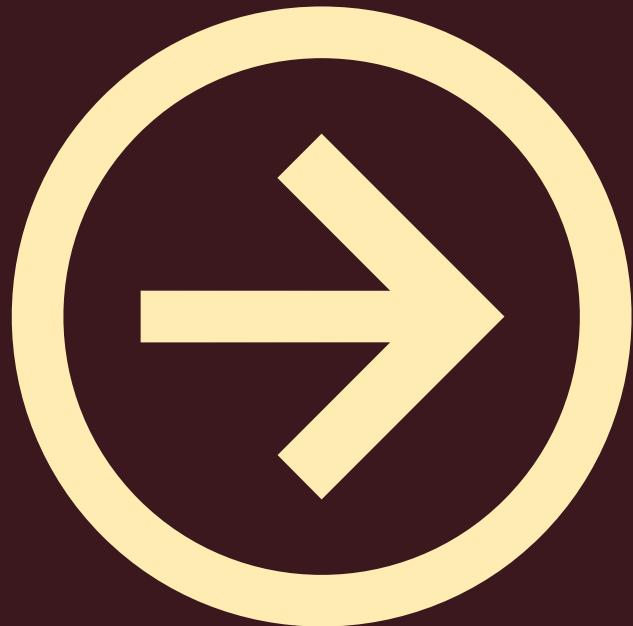


databricks®



HYPEROPT

# *Next Steps*



- 1** Continue hyperparameter tuning to optimize for recall with Hyper Opt and ML Flow
- 2** Deal with clustered reconstruction error through more feature engineering
- 3** Consider using the diagnostic data as a confounding factor
- 4** Further productionalize pipelines and code to ensure no leakage and have better deployment



# Thank You for Listening!

Questions?