

# Rod Pump Failure Analysis: Final Presentation

Preston Hart, Rohan Sura, John Mason, Ali Modak, and Tarini Thiagarajan  
NSC 325, Fall 2020



### **What is a Rod Pump?**

Sucker Rod Pumps are a type of artificial lift system used to drill for oil.

There is a large, international market for Sucker Rod Pumps consisting of several manufacturers and drilling companies.

### **The Rod Pump Challenge:**

The lost revenue from a single broken rod pump could be as much as \$5700 per day.

Some Sucker Rod Pumps are in remote or hard to access locations leaving fewer opportunities for maintenance and repairs.

### Our Solution:

We aimed to **identify** and **understand** factors which contribute to Rod Pump failure, and create a product which helps **maximize pump run time** and **minimize cost**.

Throughout this process, we have **developed a framework for Rod Pump telemetry analysis** which can be applied to a pumps in additional regions.

### What have we learned?

Understanding the nature of our raw data was crucial to making **informed data cleaning decisions** and identifying **key correlations** that could reveal hidden insights.

Domain knowledge was key to transforming our statistical results into **useful insights** by **contextualizing** and **confirming** correlations present in the data.

## Project Update

### MVP 1

#### Model Building:

Identified true null values and developed a preliminary procedure for data cleaning.

Implemented sklearn SVC, KNN, and random forest classifiers to establish a baseline for model performance.

Created a plan for improvements to our data selection and model selection procedures.

#### Dashboard:

Initial classifications and Key Insights on PowerBI

### MVP 2

#### Model Building:

Improved data cleaning procedure: added imputation and class balancing functions

Further developed KNN and random forest classifier models using hyperparameter grid search and cross validation

Implemented the SHAP explainer to understand the importance of individual variables

#### Dashboard:

Revamped initial classifications and added data visualization for timeline of failure. Added SHAP visuals and explainer plots

### MVP 3

#### Model Building:

Introduced noise variable and removed collinear variables to perform better feature reduction

Created a function to score the performance of KNN Imputer

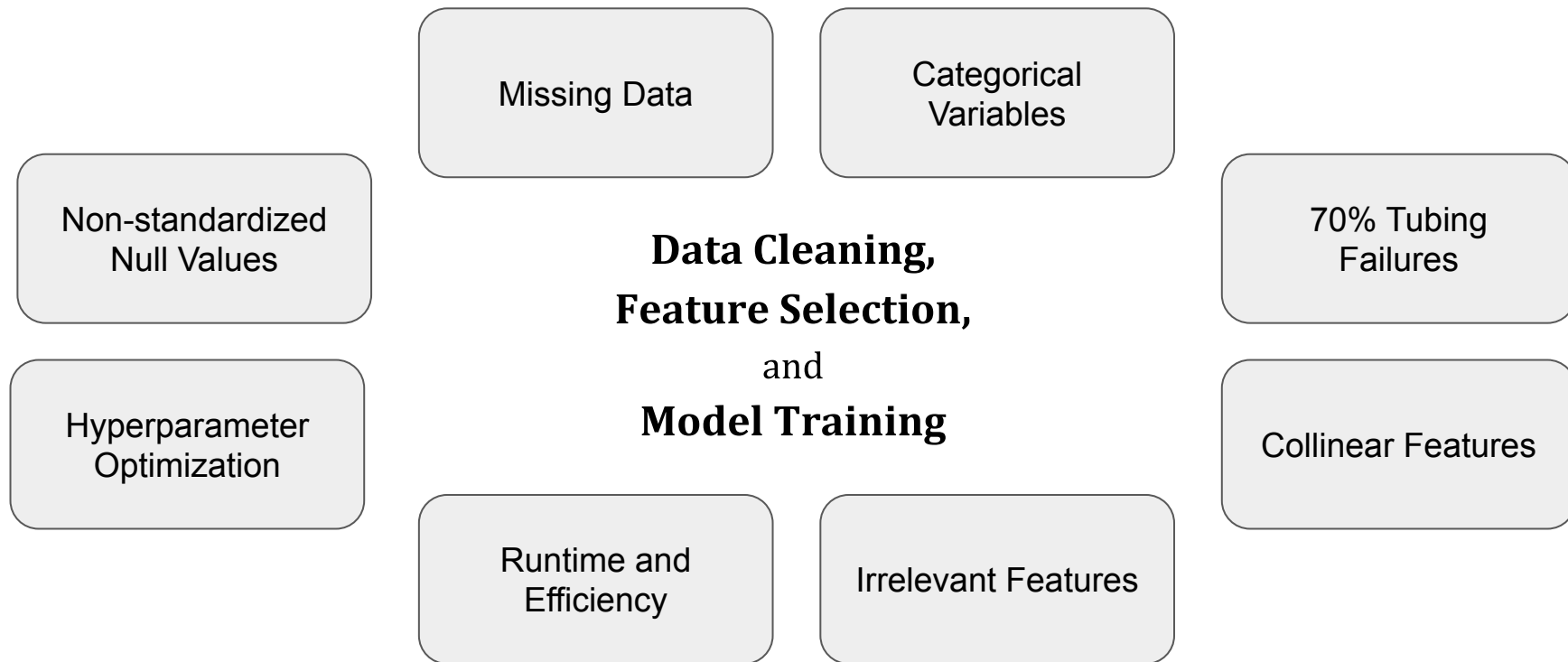
Improved SHAP explainer interpretations using affirmations from PCA analysis

#### Dashboard:

Identifying all variable to validate key influencers of shap models

Brokedown data representing factors infleuncing rod pump failure

## Our Procedure and Challenges



# Our Procedure

## Our functions:

Grid Search:  
Data Cleaning Parameters

SHAP Explainer

Feature Reduction

Grid Search:  
Hyperparameters

Standardize Null Values

Address Null Columns:  
Drop column or impute

Categorical Variables:  
One-Hot Encoder

Uneven Categories:  
Balance Dataset

Address Null Columns:  
Run KNN imputer

```
data_cleaning.py - Visual Studio Code
File Edit Selection View Go Run Terminal Help

data_cleaning.py X
1 import pandas as pd
2 from sklearn.preprocessing import LabelEncoder
3 from sklearn.preprocessing import OneHotEncoder
4 from sklearn.impute import KNNImputer
5 from sklearn.metrics import getPercentNull
6 import numpy as np
7 from numpy import random
8
9 def load(source):
10     if source == 'null transformed':
11         return pd.read_csv('floridaman/data/jmw_raw.csv', parse_dates=['lifetime_start', 'lifetime_end'])
12     else:
13         return
14
15 def preprocess(df):
16     # Calculate lifetime from lifetime_start and lifetime_end
17     df['lifetime'] = (df['lifetime_end'] - df['lifetime_start']).dt.days
18
19     # Remove procedural and redundant columns
20     del df['FAILLSTART']
21     del df['WOODID']
22     del df['IDWELL']
23     del df['REPORTID']
24     del df['tbguid']
25     del df['IDRECJOBHULL']
26
27     # Remove rows corresponding to columns
28     df = df[df['FAILURETYPE'] != 'notnull']
29
30 def encode(df, column):
31     le = LabelEncoder()
32     df[column] = le.fit_transform(df[column])
33
34     # classes = list(le.classes_)
35     # print(column + '\n')
36     # for key in range(len(classes)):
37     #     print(f'{key} : {str(classes[key])}')
38     # print('\n')
39     return df
40
41 def encode(df):
42     categorical_columns = [col for col in df.columns if df[col].dtypes == 'O']
43     categorical_columns.remove('UNIT')
44     categorical_columns.remove('model')
45     categorical_columns.remove('FAILURETYPE')
46     df = pd.get_dummies(df, columns=categorical_columns)
47     return df
48
49 def generate_candidate_dataset(data_in, COLUMN_DROP_THRESHOLD, COLUMN_IMPUTE_THRESHOLD, neighbors):
50     this_data = data_in.copy()
51     categorical_columns = [col for col in this_data.columns if this_data[col].dtypes == 'O']
52     quantitative_columns = [col for col in this_data.columns if this_data[col].dtypes == 'float']
53
54     for column in this_data:
55         this_percent_null = getPercentNull(this_data, column)
56
57         if (this_percent_null >= COLUMN_DROP_THRESHOLD):
58             del this_data[column]
59         elif (this_percent_null >= COLUMN_IMPUTE_THRESHOLD):
60             this_data.dropna(subset=[column], axis=0, inplace=True)
61         else:
62             if (column in categorical_columns):
63                 this_data.dropna(subset=[column], axis=0, inplace=True)
64
65     categorical_columns = [col for col in this_data.columns if this_data[col].dtypes == 'O']
66
67 jupyter Python 3.6.4 64bit @ 9/2/8
```



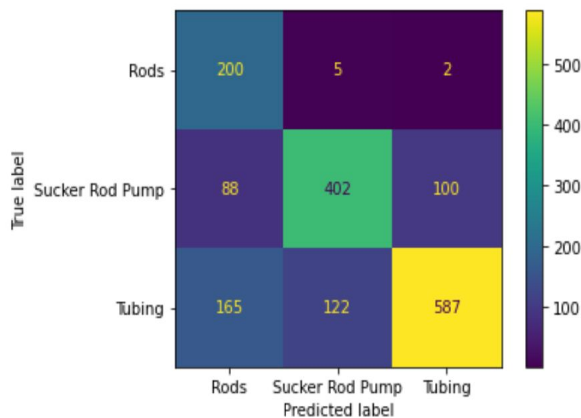
TEXAS



## Random Forest Model Results:

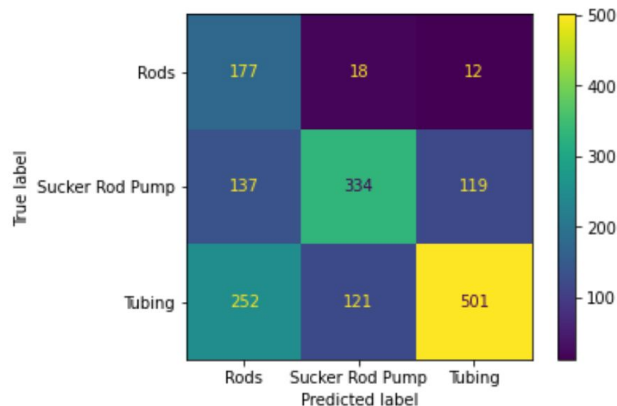
### Balancing After Imputation

*Accuracy: .71*



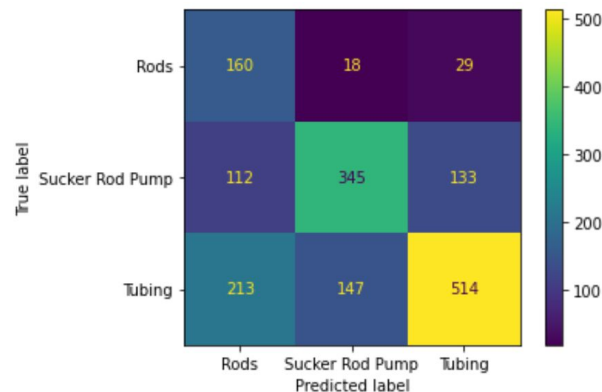
### Balancing Before Imputation

*Accuracy: .61*



### Reduced Features

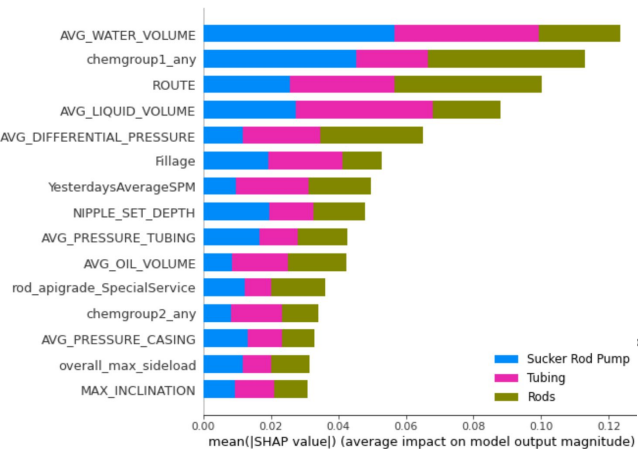
*Accuracy: .62*



## Feature Importance:

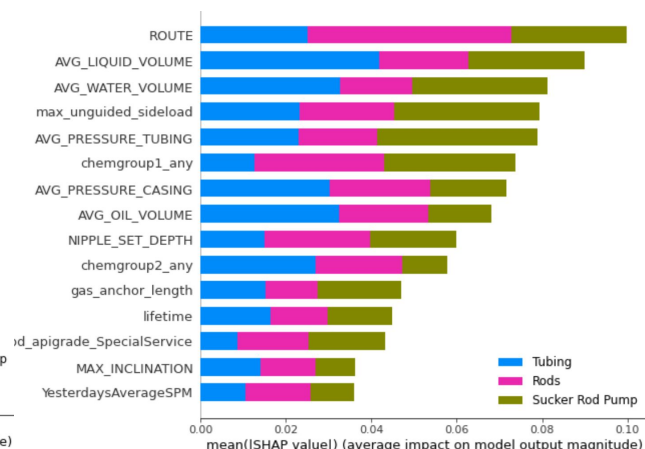
### Balancing After Imputation

*Accuracy: .71*



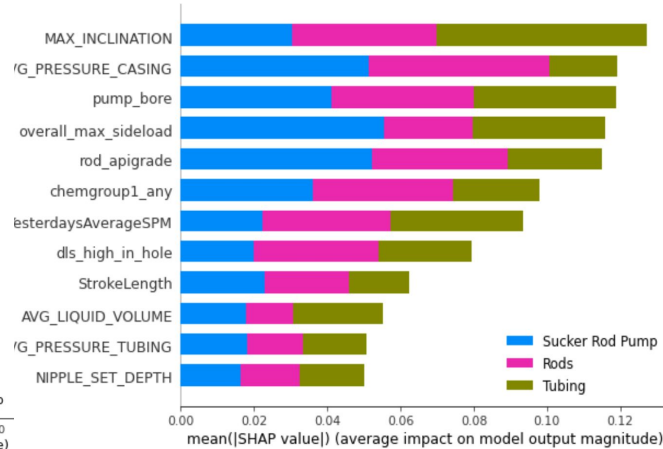
### Balancing Before Imputation

*Accuracy: .61*



### Reduced Features

*Accuracy: .63*



TEXAS



## Dashboard and SHAP Update

Live Demo

## Business Impact and Recommendations:

### Recommendations:

1. When **max inclination** is ...
  - a.  $> \sim 2.15$ , rods were **8.37%** more likely to classify as a **Tubing** failure.
  - b.  $> \sim 2.678$ , rods were **13.75%** less likely to classify as a **SRP** failure.
  - c.  $> \sim 3.22$ , rods were **8.06%** more likely to classify as a **Rod** failure.
2. When **avg\_pressure\_casing** is ...
  - a.  $> \sim 112.5$  rods were **23.18%** more likely to classify as a **Tubing** failure.
  - b.  $< \sim 78.95$ , rods were **13.13%** more likely to classify as a **SRP** failure.
  - c.  $> \sim 123.8$ , rods were **15.01%** less likely to classify as a **Rod** failure.

### Business Impact:

1. In light of the above insights regarding failure type classification likelihood, ConocoPhillips can work towards preemptively addressing a rod's failure
2. Help capture majority of lost profit of rod pump failure (\$5700/day)
3. Reduce time spent by ConocoPhillips on understanding leading factors contributing to rod pump failure (10 hrs/week)
4. A prepared dashboard and series of shap plots providing data breakdown and likeliness of each failure

## Conclusion & Next Steps:

- Improve balancing function and decrease bias introduced by random sampling by implementing a method which selectively chooses samples which well represents the distribution as a whole
- Explore feature decomposition, splitting features into subcategories based on thresholds observed in SHAP values
- Test a one-versus-all methodology which will allow us to train a unique model for each failure type
- New Model Exploration: Gradient Boosted Trees
- Implement Survival Regression Analysis to determine timeline of failures and predict when they will occur