

PREDICTIONS IN A MINING PROCESS

TEAM MEMBERS



Geovane Marquez



Ronak Parikh



Anthony Castillo



Sam Sheridan

INTRODUCTION



THE TARGET IS TO PREDICT THE % OF SILICA IN THE END OF THE PROCESS, WHICH IS THE CONCENTRATE OF IRON ORE AND ITS IMPURITY (WHICH IS THE % OF SILICA).

THE DATASET IN QUESTION SPECIFICALLY HAS THE PERCENTAGE OF SILICA IMPURITY AS THE RESPONSE VARIABLE. THUS, THE AUTHOR OF THE DATASET INTENDED FOR USERS TO RESEARCH HOW MUCH SILICA IMPURITY IS IN THE IRON ORE SO MANUFACTURING PROCESSES CAN BE IMPROVED.

WE PERFORM RANDOM FOREST, DECISION TREES, AND GRADIENT BOOSTING TREE ALGORITHMS TO MAKE THE MOST ACCURATE PREDICTIONS.

DATA I

WITH OVER 737,453 RECORDS AND 24 FEATURES, THE ANALYSIS WILL CONSIST OF USING ONLY 22 OF THE FEATURES. THE FEATURES 'DATE' AND '% IRON FEED' WILL BE EXCLUDED FROM THE ANALYSIS. THE FEATURE OF INTEREST IS THE % SILICA FEED IN APROPOS TO THE OTHERS.

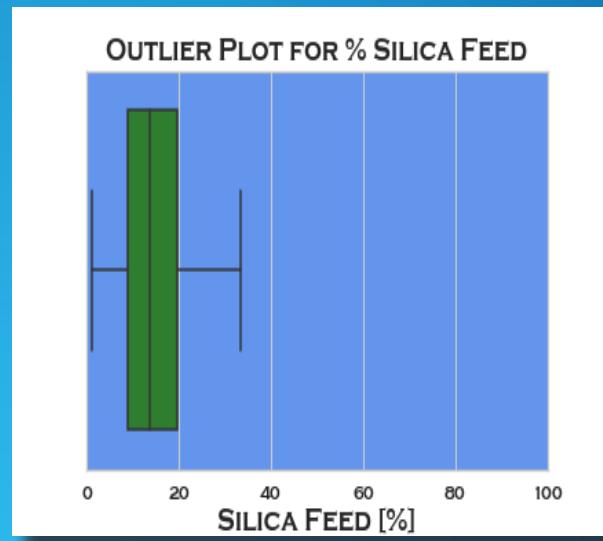
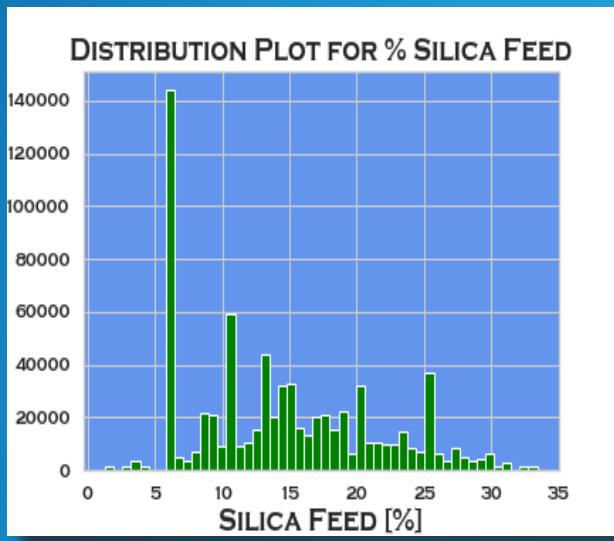
DATA SOURCE: [HTTPS://WWW.KAGGLE.COM/EDUMAGALHAES/QUALITY-PREDICTION-IN-A-MINING-PROCESS](https://www.kaggle.com/edumagalhaes/quality-prediction-in-a-mining-process)

FEATURES	DATA TYPE
DATE	STRING
% IRON FEED	DOUBLE
% SILICA FEED	DOUBLE
STARCH FLOW	DOUBLE
AMINA FLOW	DOUBLE
ORE PULP FLOW	DOUBLE
ORE PULP PH	DOUBLE
ORE PULP DENSITY	DOUBLE
FLOTATION COLUMN 01 AIR FLOW	DOUBLE
FLOTATION COLUMN 02 AIR FLOW	DOUBLE
FLOTATION COLUMN 03 AIR FLOW	DOUBLE
FLOTATION COLUMN 04 AIR FLOW	DOUBLE
FLOTATION COLUMN 05 AIR FLOW	DOUBLE
FLOTATION COLUMN 06 AIR FLOW	DOUBLE
FLOTATION COLUMN 07 AIR FLOW	DOUBLE
FLOTATION COLUMN 01 LEVEL	DOUBLE
FLOTATION COLUMN 02 LEVEL	DOUBLE
FLOTATION COLUMN 03 LEVEL	DOUBLE
FLOTATION COLUMN 04 LEVEL	DOUBLE
FLOTATION COLUMN 05 LEVEL	DOUBLE
FLOTATION COLUMN 06 LEVEL	DOUBLE
FLOTATION COLUMN 07 LEVEL	DOUBLE
IRON CONCENTRATE	DOUBLE

THE SUMMARY STATISTICS FOR THE % OF SILICA AT THE END OF THE PROCESS

COUNT	MIN	MAX	AVERAGE	STDDEV
737,453	1.31%	33.40%	14.65%	6.81%

DATA II



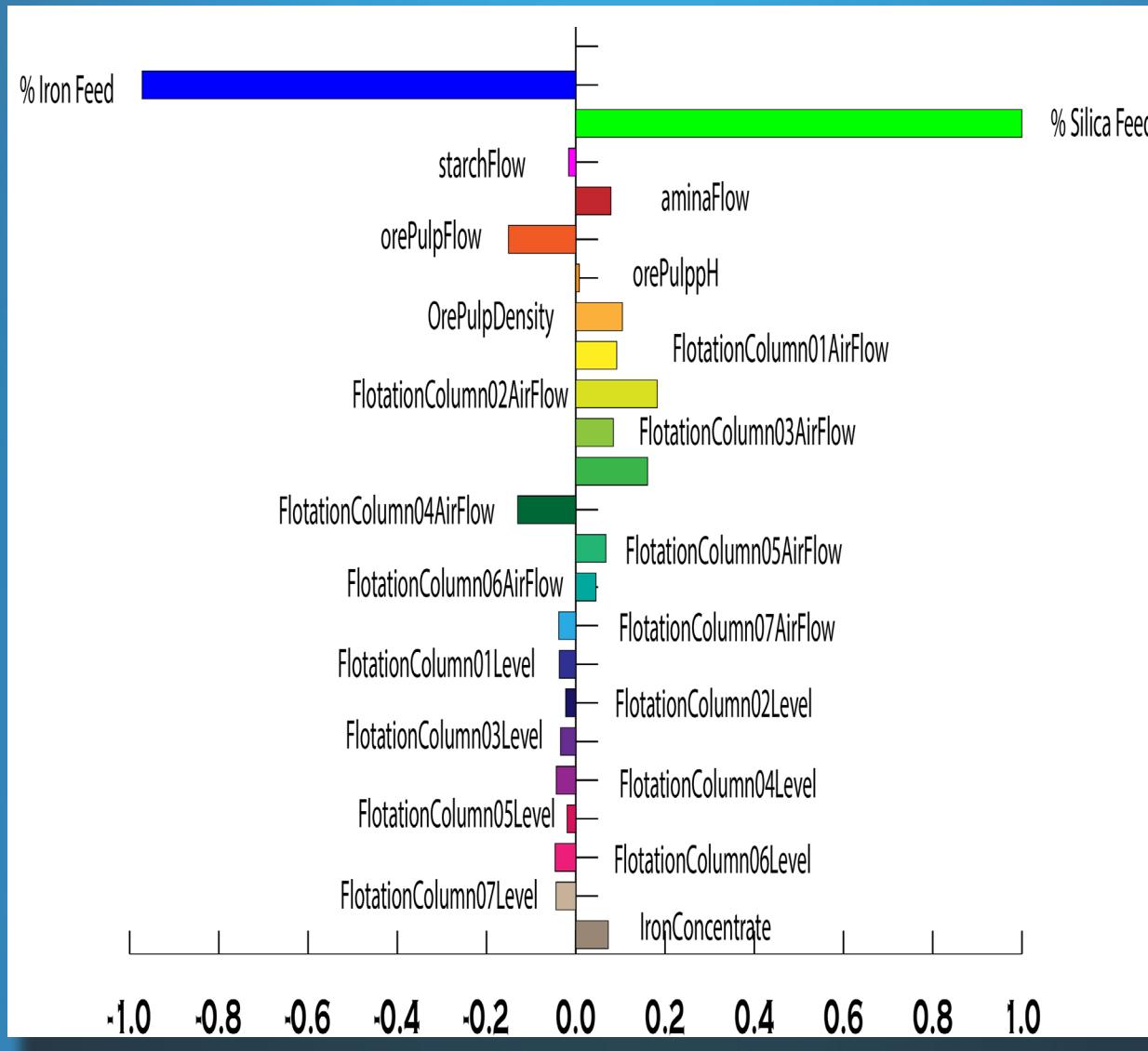
THE TARGET VARIABLE ‘% SILICA FEED’ HAS A RANGE OF PERCENTAGES BETWEEN 1.31 AND 33.40 FOR THE AMOUNT OF SILICA FOUND AT THE END OF THE MINING PROCESS.
THERE IS A HIGH COUNT AROUND 6% OF SILICA FEED.

EXPLORING FURTHER, NO OUTLIERS WERE DETECTED OUTSIDE THE MINIMUM AND MAXIMUM RANGES OF PERCENTAGES. THE DATA IS CONSISTENT, AND APPEARS TO BE MORE DENSE BETWEEN 10% AND 20% RANGE. EARLIER, THE AVERAGE STATISTIC WAS FOUND TO BE 14.65%.

IN THE FOLLOWING PAGES, WE EXPLORE THE RELATIONSHIPS THAT THE TARGET VARIABLE HAS WITH THE OTHER FEATURES. THIS BECOMES HELPFUL IN ELIMINATING UNNEEDED FEATURES.

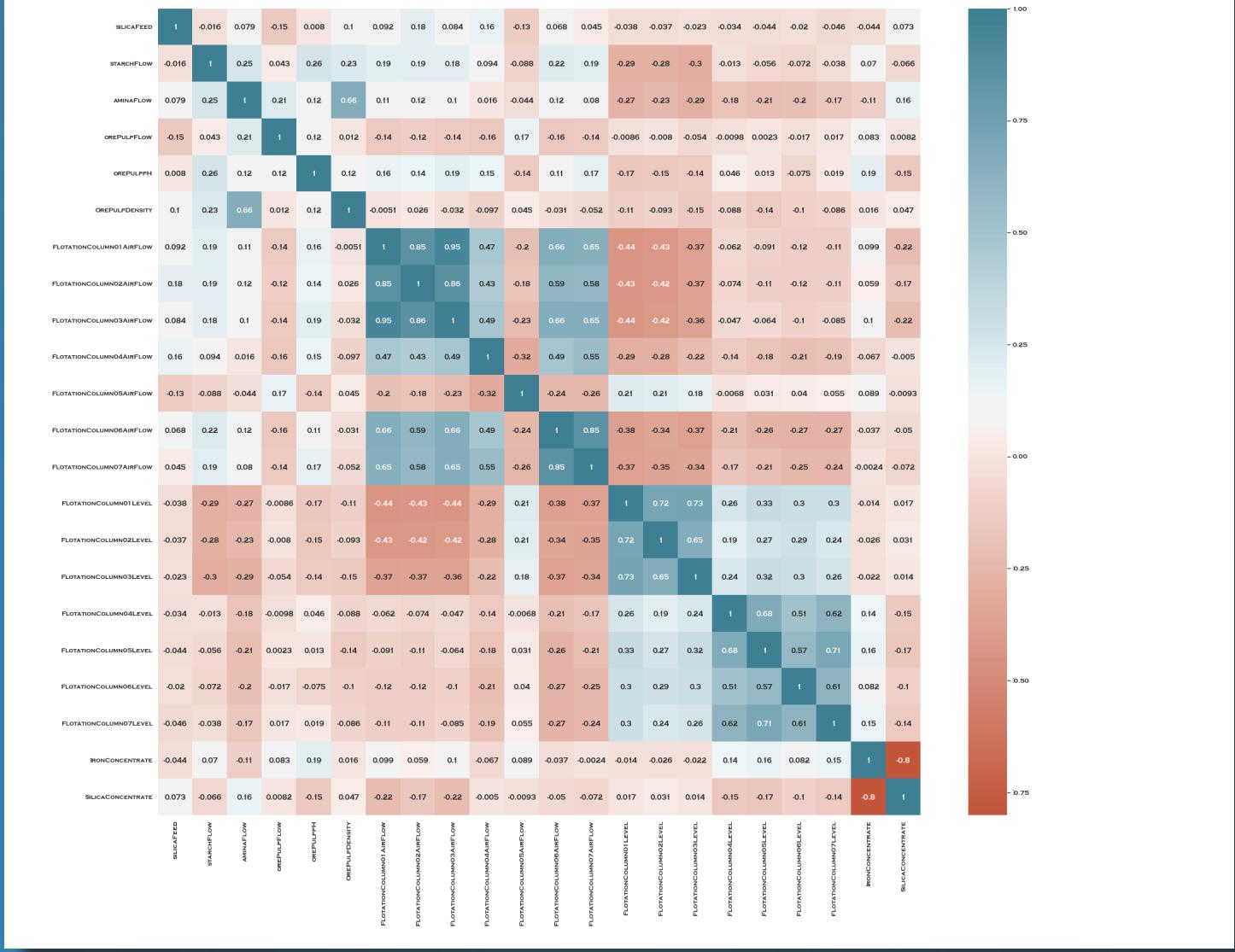
MINING DATA I

CORRELATION BUILT FROM
THE RESPONSE:
% SILICA FEED



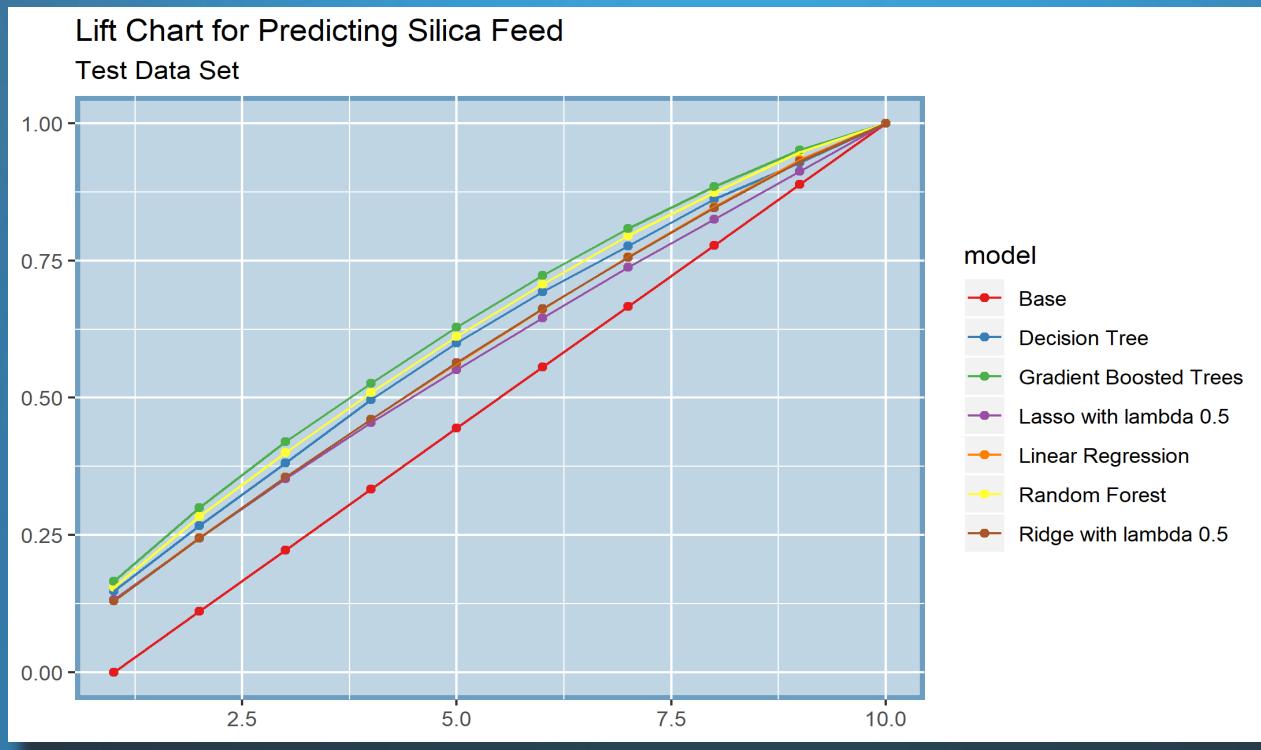
IRON FEED AND SILICA FEED ARE HIGHLY CORRELATED. IN THIS STUDY, WE ARE ONLY INTERESTED IN EXPLORING THE PREDICTIONS FOR % SILICA FEED, THUS WE REMOVED % IRON FEED OUT OF THE MODEL. IN THE NEXT PAGE, WE EXPLORE THESE RELATIONS WITH A CORRELATION HEAT MATRIX.

MINING DATA II

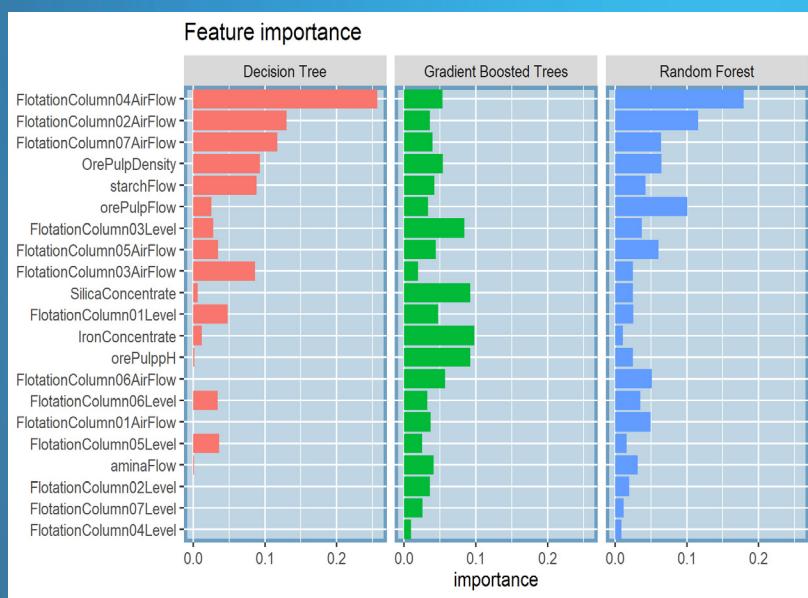


THE PLOT SHOWS THE CORRELATIONS BETWEEN THE FEATURES. WE USE THIS TO FIND OUT WHICH FEATURES AFFECT % SILICA FEED THE MOST. THE PLOT GIVES US A GREAT OVERVIEW OF THE FEATURE CORRELATION RELATIONSHIPS.

MINING DATA III

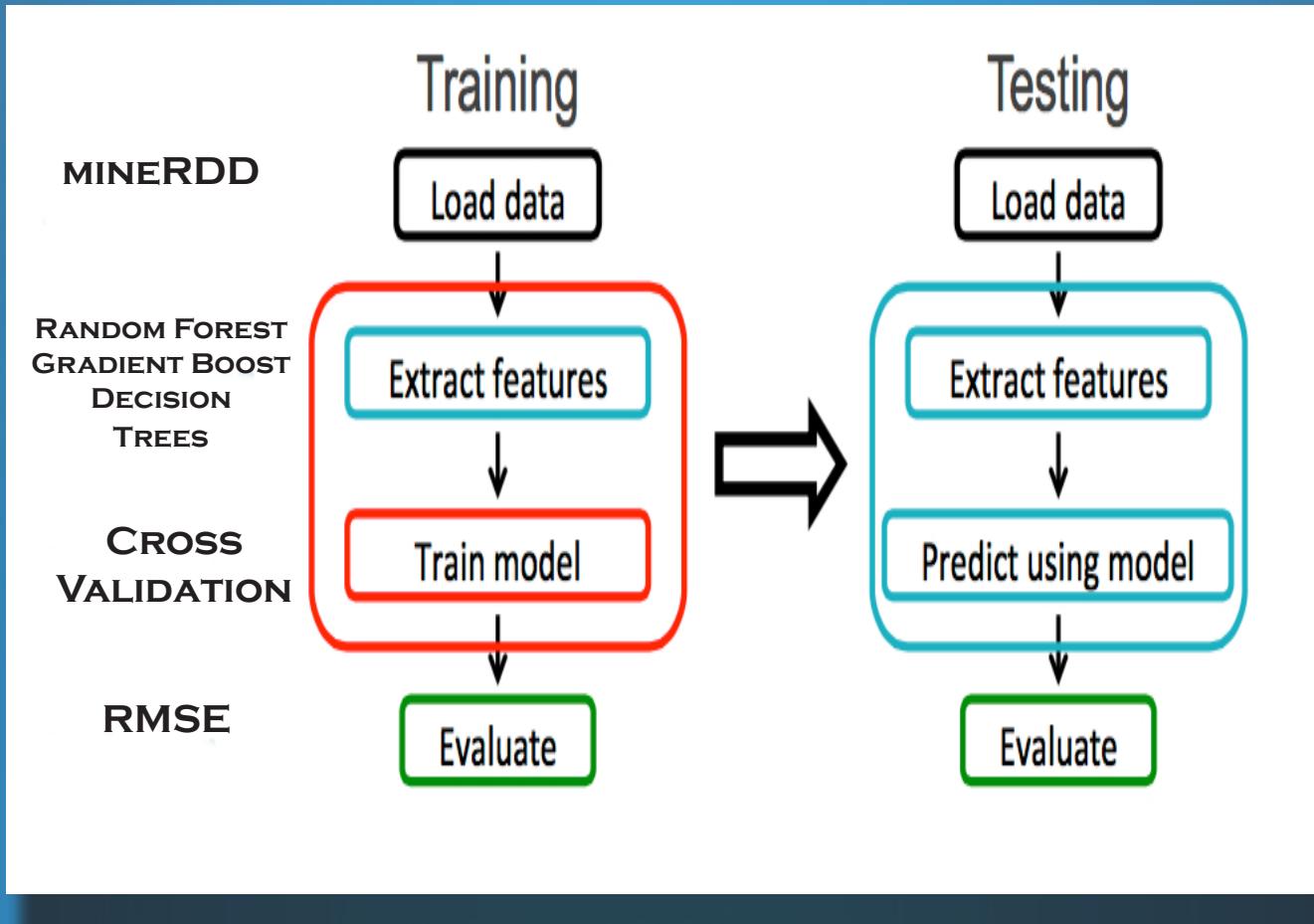


AFTER SPLITTING THE DATA: 75% TRAINING, 25% TEST. SEVERAL MACHINE LEARNING MODELS WERE CONSTRUCTED WITH THE TRAINING DATA AND TESTED AGAINST THE TEST DATA, AS SHOWN IN THE PLOT ABOVE. THE PLOT ABOVE TELLS THE STORY THAT THE BEST PREDICTIONS WILL COME FROM GRADIENT BOOSTED TREES, RANDOM FOREST, AND DECISION TREES. THUS, FINE-TUNED MODELS WILL BE FOCUSED ON THESE THREE.



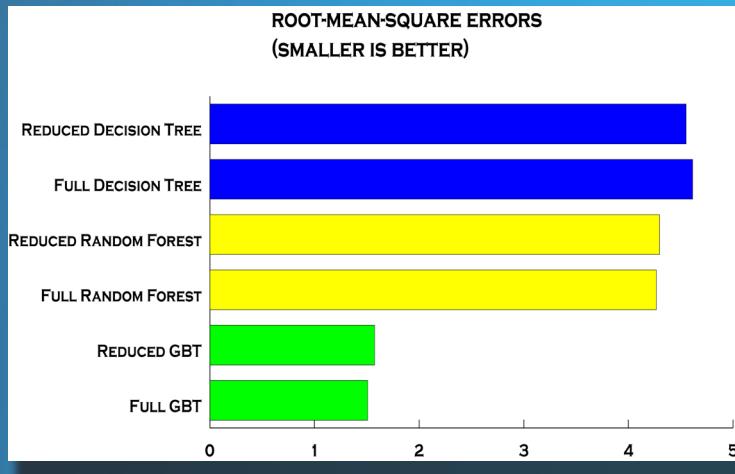
BASED ON THE PLOT TO THE LEFT, A REDUCED MODEL WILL BE CONSTRUCTED. THE REDUCED MODEL WILL THEN BE TESTED AGAINST THE FULL MODEL. BY OBSERVATION, THE LEAST IMPORTANT FEATURES ARE: AMINAFLOW, AND FLOTATION COLUMN LEVELS 02,04,07. THESE THREE WILL BE REMOVED, AND THE REST OF THE FEATURES WILL BE KEPT IN THE REDUCED MODEL.

PIPELINES



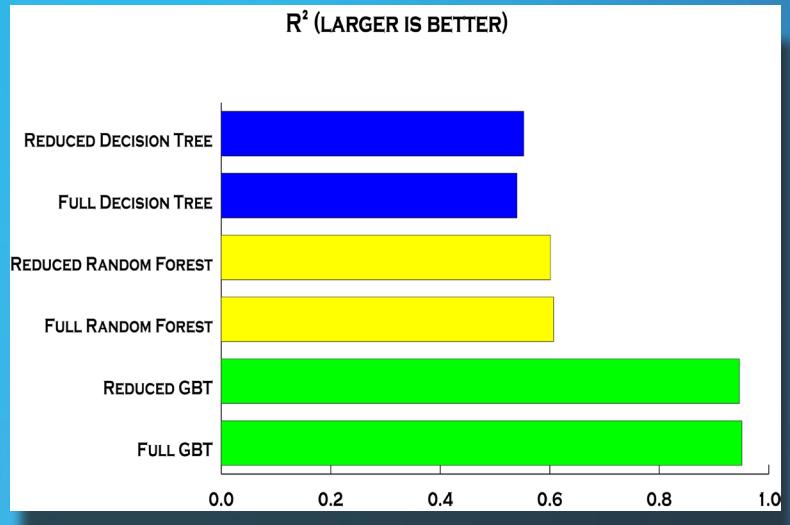
THE GRAPHIC ABOVE ILLUSTRATES THE PROCESS THAT WE TOOK TO GENERATE A PREDICTIVE FITTED MODEL. FIRST WE LOADED THE DATA INTO PYSPARK AS AN RDD. SECONDLY, WE FITTED THREE SUBMODELS USING GRADIENT BOOSTING, DECISION, AND RANDOM FOREST TREES. THIRDLY, WE PERFORMED CROSS VALIDATIONS TO GENERATE THE BEST PARAMETER METRICS. LASTLY, THE ERRORS USING RMSE WERE PLOTTED AGAINST EACH OTHER FOR COMPARISON, AND THE MODEL WITH THE LOWEST ERROR WAS CHOSEN TO BE USED TO PREDICT THE % OF SILICA FEED.

PERFORMANCE AND VALIDATION



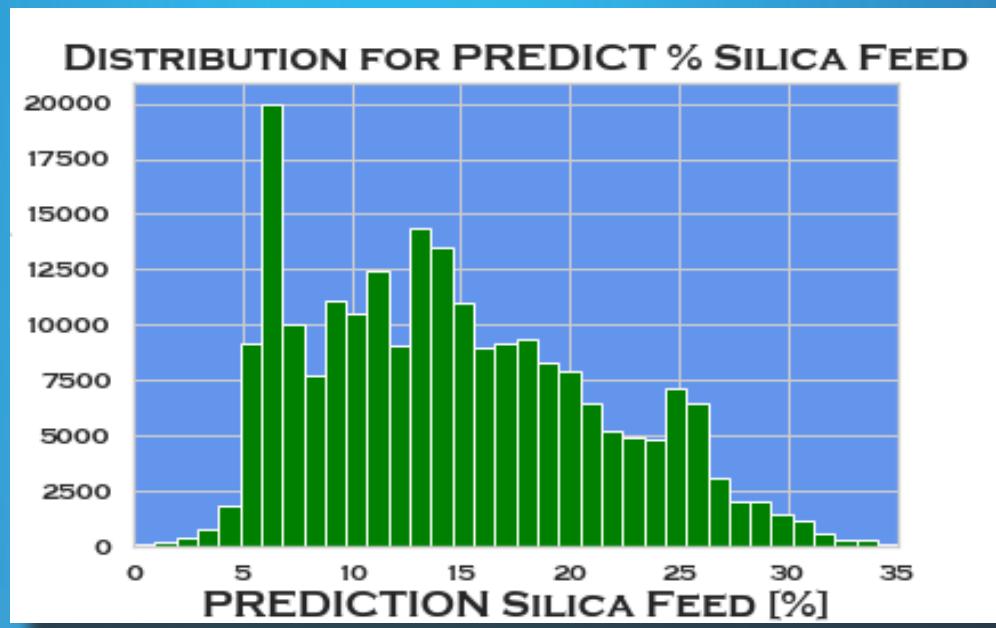
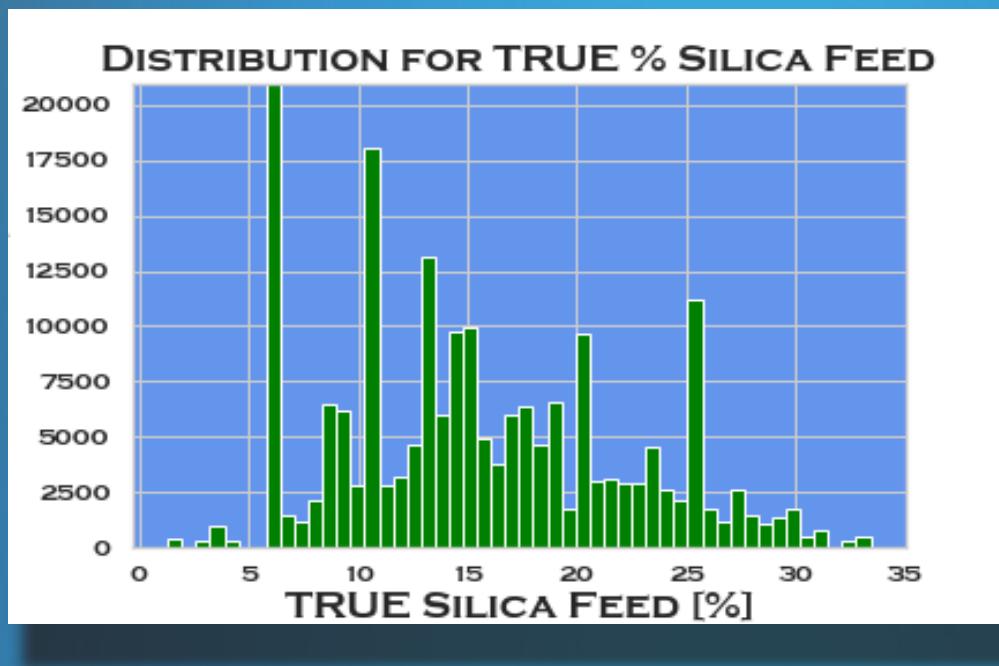
AFTER PERFORMING CROSS VALIDATIONS USING PARAGRID, AS PREDICTED BY THE EARLIER LIFT CHART, THE MODEL WITH THE SMALLEST ROOT-MEAN-SQUARE ERROR (RMSE) WAS THE GRADIENT BOOSTING TREE MODEL. THE FULL MODEL WAS INFINITESIMALLY BETTER THAN THE REDUCED MODEL WHEN COMPARING THEIR RMSES.

THE PROPORTION OF VARIANCE EXPLAINED BY THE MODEL WAS HIGHEST IN THE GRADIENT BOOSTING TREE MODELS. SIMILARLY, THE FULL GBT MODEL WAS INFINITESIMALLY BETTER THAN THE REDUCED MODEL WHEN COMPARING THEIR R-SQUARED



IN CONCLUSION, THE MODEL WE CHOSE TO CONDUCT THE PREDICTIONS FOR % SILICA FEED WAS DETERMINED TO BE THE FULL MODEL TRAINED UNDER A GRADIENT BOOSTING TREES ALGORITHM. THE HIGHLY COMPUTATIONAL EXPENSIVE MODEL CONSISTED OF AN ENSEMBLE OF MAX DEPTH OF 10 OF EACH DECISION TREE. AS WELL AS, OVER 200 ITERATIONS FOR THE NUMBER OF TREES IN EACH BOOSTING TREE ENSEMBLE.

PREDICTIONS



ABOVE WE SEE THE ORIGINAL TRUE DATA FOR % SILICA FEED. RIGHT BELOW THIS GRAPH WE HAVE THE NEWLY PREDICTED DATA FOR % SILICA FEED. THE SHAPE OF THE PREDICTION PLOT HAS SLIGHT SKEWNESS TO IT, HOWEVER, THE DISTRIBUTION LOOKS NICELY CONSISTENT.

LASTLY, THE PREDICTIONS WERE CONSTRUCTED FROM A GRADIENT BOOSTING TREE ALGORITHM THAT UTILIZED ALL THE FEATURES (EXCEPT DATE AND % IRON FEED), AND USED % SILICA FEED AS THE TARGET. THIS TRAINED MODEL CAN HELP RESEARCH HOW MUCH SILICA IMPURITY IS IN THE IRON ORE SO MANUFACTURING PROCESSES CAN BE IMPROVED.

SparklyR Appendix

Mining Project - PSTAT 135

```
library(sparklyr)
library(dplyr)
library(tidyverse)

spark_install(version='2.4')
sc <- spark_connect(master='local',version='2.4')

d <- spark_read_csv(sc, path='editedMine.csv', header=TRUE, infer_schema = TRUE)

mine_rdd = d
```

Feature Selection by CV

```
# source: https://www.eddjberry.com/post/2018-12-12-sparklyr-feature-selection/
get_feature_cols <- function(tbl, response, exclude = NULL) {

  # column names of the data
  columns <- colnames(tbl)

  # exclude the response/outcome variable and
  # exclude from the column names
  columns[!(columns %in% c(response, exclude))]
}

feature_cols <- get_feature_cols(mine_rdd, response ='silicaFeed')
#feature_cols
```

Split data

```
ds <- sdf_random_split(mine_rdd, train=0.75, test=0.25, seed=123)
names(ds)

train_data <- ds$train
test_data <- ds$test
```

Full Model

```
ml_formula <- formula(silicaFeed ~ .)
```

Reduced Model

```
red.mod <- formula(silicaFeed ~ starchFlow+orePulpFlow+orePulppH+
OrePulpDensity+FlotationColumn01AirFlow +
FlotationColumn02AirFlow+FlotationColumn03AirFlow+
FlotationColumn04AirFlow+FlotationColumn05AirFlow+
FlotationColumn06AirFlow+FlotationColumn07AirFlow+
FlotationColumn01Level+FlotationColumn03Level+
FlotationColumn05Level+FlotationColumn06Level+
IronConcentrate+SilicaConcentrate)

reducedRDD <- select(mine_rdd, -c(aminaFlow, FlotationColumn02Level,
FlotationColumn04Level, FlotationColumn07Level))
```

Pre-Tuning ML Models

```
# Decision Tree
ml_dt <- ml_decision_tree(train_data, ml_formula)

# Random Forest
ml_rf <- ml_random_forest(train_data, ml_formula)

# Gradient Boosted Tree
ml_gbt <- ml_gradient_boosted_trees(train_data, ml_formula)

# linear regression
ml_linreg <- ml_linear_regression(train_data, ml_formula)

# lasso
ml_lasso5 <- ml_linear_regression(train_data, formula = ml_formula,
elastic_net_param=1, reg_param = 0.5)

# Ridge
ml_ridge5 = ml_linear_regression(train_data, ml_formula,
elastic_net_param =0, reg_param = 0.5)
```

Validation

```
ml_models <- list(
"Decision Tree" = ml_dt,
"Random Forest" = ml_rf,
"Gradient Boosted Trees" = ml_gbt,
"Linear Regression" = ml_linreg,
"Lasso with lambda 0.5" = ml_lasso5,
"Ridge with lambda 0.5" = ml_ridge5
)

# Create a function for scoring
```

```

score_test_data <- function(model, data = test_data){
  pred <- sdf_predict(data, model)
  select(pred, silicaFeed, prediction)
}

# Score all the models
ml_score <- map(ml_models, score_test_data)

```

LIFT CHART

```

# Lift function
calculate_lift <- function(scored_data) {
  scored_data %>%
    mutate(bin = ntile(desc(prediction), 10)) %>%
    group_by(bin) %>%
    summarize(count = sum(silicaFeed)) %>%
    mutate(prop = count / sum(count)) %>%
    arrange(bin) %>%
    mutate(prop = cumsum(prop)) %>%
    select(-count) %>%
    collect() %>%
    as.data.frame()
}

# Initialize results
ml_gains <- tibble(
  bin = seq(from = 1, to = 10),
  prop = seq(0, 1, len = 10),
  model = "Base"
)

# Calculate lift
for(i in names(ml_score)){
  ml_gains <- ml_score[[i]] %>%
    calculate_lift %>%
    mutate(model = i) %>%
    bind_rows(ml_gains, .)
}

# Plot results
myplot = ggplot(ml_gains, aes(x = bin, y = prop, color = model)) +
  geom_point() +
  geom_line() +
  scale_color_brewer(palette = 'Set1') +
  labs(title = "Lift Chart for Predicting Silica Feed",
       subtitle = "Test Data Set",
       x = NULL,
       y = NULL) +
  theme(
    panel.background = element_rect(fill = "#BFD5E3", colour = "#6D9EC1",
                                    size = 2, linetype = "solid"),

```

```

    panel.grid.major = element_line(size = 0.5, linetype = 'solid',
                                    colour = "white"),
    panel.grid.minor = element_line(size = 0.25, linetype = 'solid',
                                    colour = "white")
)
myplot

```

FEATURE IMPORTANCE CHART

```

# Initialize results
feature_importance <- tibble()

# Calculate feature importance
for(i in c("Decision Tree", "Random Forest", "Gradient Boosted Trees")){
  feature_importance <- ml_tree_feature_importance(ml_models[[i]]) %>%
    mutate(Model = i) %>%
    rbind(feature_importance, .)
}

# Plot results
feature_importance %>%
  ggplot(aes(reorder(feature, importance), importance, fill = Model)) +
  facet_wrap(~Model) +
  geom_bar(stat = "identity") +
  coord_flip() +
  labs(title = "Feature importance",
       x = NULL) +
  theme( legend.position ="none",
        panel.background = element_rect(fill = "#BFD5E3", colour = "#6D9EC1",
                                         size = 2, linetype = "solid"),
        panel.grid.major = element_line(size = 0.5, linetype = 'solid',
                                         colour = "white"),
        panel.grid.minor = element_line(size = 0.25, linetype = 'solid',
                                         colour = "white"))
)

```

Spark Appendix

In [1]:

```
from pyspark import SparkConf
from pyspark import SparkContext
from pyspark.sql import SparkSession, SQLContext
import os
import pandas as pd
import os.path
#sc = SparkContext.getOrCreate()
conf = SparkConf().set("spark.cores.max", "8") \
    .set("spark.driver.memory", "16g") \
    .set("spark.executor.memory", "16g") \
    .set("spark.executor.memory_overhead", "16g") \
    .set("spark.driver.maxResultsSize", "0")
sc = SparkContext(appName="mineProject", conf=conf)
spark = SparkSession(sc)
```

Read Data

In [2]:

```
filename = 'MiningProcess_Flotation_Plant_Database.csv'
mine_pd = pd.read_csv(filename, decimal=',')
```

In [3]:

```
newColumns = ["date", "ironFeed", "silicaFeed", "starchFlow", "aminaFlow", "orePulpFlow", \
    "orePulppH", "OrePulpDensity", "FlotationColumn01AirFlow", \
    "FlotationColumn02AirFlow", "FlotationColumn03AirFlow", \
    "FlotationColumn04AirFlow", "FlotationColumn05AirFlow", "FlotationColumn06AirFlow", "Flot \
    ationColumn07AirFlow", "FlotationColumn01Level", "FlotationColumn02Level", "FlotationColumn03Level", \
    "FlotationColumn04Level", "FlotationColumn05Level", "FlotationColumn06Level", "FlotationColumn07Level" \
    , 'IronConcentrate', 'SilicaConcentrate']
```

PANDAS AND SNS STUFF

In [4]:

```
#rename columns
mine_pd.columns = newColumns
```

In [5]:

```
mine_pd.drop(['ironFeed', 'date'], axis=1, inplace=True)
```

In [7]:

```
# Export csv for sparklyr
mine_pd.to_csv('/Users/Geo/Desktop/project_photos/editedMine.csv', index = False)
```

In [35]:

```
import matplotlib.pyplot as plt
import seaborn as sns
import numpy as np
```

Histogram of SilicaFeed

In []:

```
sns.set(style="whitegrid", font='Copperplate Gothic Bold', rc={'axes.facecolor':'cornflowerblue', 'figure.facecolor':'cornflowerblue'})
ax=sns.distplot(mine_pd['silicaFeed'], color='green', hist_kws={ "alpha": 1}, kde=False)
plt.title('Distribution Plot for % Silica Feed', fontsize=18)
plt.ylabel('Relative Frequency', fontsize=18)
plt.xlabel('Silica Feed [%]', fontsize=18)
plt.show()
```

Log Transform Silica Feed

In []:

```
ax = sns.boxplot(x=mine_pd["silicaFeed"], color="forestgreen")
ax.set_title('Outlier Plot for % Silica Feed', fontsize=18)
ax.set_xlabel('Silica Feed [%]', fontsize=18)
ax.set(xlim=(0,100))
ax.figure.savefig("plot2.png")
```

Heatmap for Silica Feed

In []:

```
plt.figure(figsize=(30, 25))
mine_heatmap = sns.heatmap(mine_pd.corr(), annot=True, cmap=sns.diverging_palette(20, 220, n=200), )
mine_heatmap.figure.savefig("plot3.png")
```

In [6]:

```
spark.conf.set("spark.sql.execution.arrow.enabled", "true")
os.environ['ARROW_PRE_0_15_IPC_FORMAT'] = "1"
```

In [8]:

```
# Drop unneed columns
mine_pd.drop(['aminaFlow', 'FlotationColumn02Level', 'FlotationColumn04Level', 'FlotationColumn07Level'],
            axis=1, inplace=True)
```

In [7]:

```
mineRDD = spark.createDataFrame(mine_pd)
```

Data Dimensions

In [8]:

```
print((mineRDD.count(), len(mineRDD.columns)))
```

```
(737453, 22)
```

Data Types

In []:

```
mineRDD.dtypes
```

The goal is to predict the % of Silica impurity in iron ore concentrate

Summary Statistics

In [15]:

```

from pyspark.sql.functions import count, mean, stddev, min, max

In [ ]:
mineRDD.select(count('silicaFeed'),mean('silicaFeed'),stddev('silicaFeed'), \
               min('silicaFeed'), max('silicaFeed')).show()

```

Building a Correlation Matrix

We want to omit any highly correlated or unnecessary features

```

In [14]:
from pyspark.ml.stat import Correlation
from pyspark.ml.feature import VectorAssembler

In [10]:
def correlation_matrix(df, corr_columns, method='pearson'):
    vector_col = "corr_features"
    assembler = VectorAssembler(inputCols=corr_columns, outputCol=vector_col)
    df_vector = assembler.transform(df).select(vector_col)
    matrix = Correlation.corr(df_vector, vector_col, method)

    result = matrix.collect()[0][f"pearson({})".format(vector_col)].values
    return pd.DataFrame(result.reshape(-1, len(corr_columns)), columns=corr_columns, index=corr_col
umns)

In [ ]:
correlation_matrix(mineRDD, mineRDD.columns[1:]).iloc[1:2].T

```

Keep all the features.

Pipelines

```

In [10]:
from pyspark.ml import Pipeline
from pyspark.ml.regression import RandomForestRegressor
from pyspark.ml.tuning import ParamGridBuilder
from pyspark.ml.evaluation import RegressionEvaluator
from pyspark.ml.tuning import CrossValidator

```

Split training and test sets

```

In [11]:
train, test = mineRDD.randomSplit([0.7, 0.3], seed = 123)

```

```

In [12]:
from pyspark.ml.feature import VectorAssembler, VectorIndexer

```

Vector Assembler

```

In [13]:
featuresCols = mineRDD.columns
featuresCols.remove('silicaFeed')

vectorAssembler = VectorAssembler(inputCols=featuresCols, outputCol="rawFeatures")

```

```
# This identifies categorical features and indexes them.  
vectorIndexer = VectorIndexer(inputCol="rawFeatures", outputCol="features")
```

RANDOM FOREST - REDUCED MODEL

In [15]:

```
rf = RandomForestRegressor(labelCol='silicaFeed')
```

In [16]:

```
# Define a grid of hyperparameters to test:  
  
paramGrid = ParamGridBuilder() \  
.addGrid(rf.numTrees, [2,10]) \  
.addGrid(rf.maxDepth, [2,10]) \  
.build()  
  
#This tells CrossValidator how well we are doing by comparing the true labels with predictions.  
evaluator = RegressionEvaluator(metricName="rmse", labelCol=rf.getLabelCol(), predictionCol=rf.getP  
redictionCol())  
# Declare the CrossValidator, which runs model tuning for us.  
cv = CrossValidator(estimator=rf, evaluator=evaluator, estimatorParamMaps=paramGrid)
```

In [17]:

```
pipeline = Pipeline(stages=[vectorAssembler, vectorIndexer, cv])
```

In [24]:

```
pipelineRFModel = pipeline.fit(train)
```

In [26]:

```
predictions = pipelineRFModel.transform(test)
```

In [27]:

```
rmse = evaluator.evaluate(predictions)  
print("RMSE on our test set: %g" % rmse)
```

```
RMSE on our test set: 4.2974
```

\$R^2\$

In [29]:

```
y_true = predictions.select("silicaFeed").toPandas()  
y_pred = predictions.select("prediction").toPandas()
```

In [31]:

```
import sklearn.metrics  
r2 = sklearn.metrics.r2_score(y_true, y_pred)  
print('r2: {0}'.format(r2))
```

```
r2: 0.6011599040096292
```

gradient Boost Trees

In [16]:

```
from pyspark.ml.regression import GBTRegressor
```

```
In [17]:
```

```
# Takes the "features" column and learns to predict "silicaFeed"  
gbt = GBTRegressor(labelCol="silicaFeed")
```

```
In [18]:
```

```
# Define a grid of hyperparameters to test:  
paramGrid = ParamGridBuilder()\  
    .addGrid(gbt.maxDepth, [2,10])\  
    .addGrid(gbt.maxIter, [10,200])\  
    .build()  
  
#This tells CrossValidator how well we are doing by comparing the true labels with predictions.  
evaluator = RegressionEvaluator(metricName="rmse", labelCol=gbt.getLabelCol(), predictionCol=gbt.get  
tPredictionCol())  
# Declare the CrossValidator, which runs model tuning for us.  
cv = CrossValidator(estimator=gbt, evaluator=evaluator, estimatorParamMaps=paramGrid)
```

```
In [19]:
```

```
pipeline = Pipeline(stages=[vectorAssembler, vectorIndexer, cv])
```

```
In [ ]:
```

```
pipelineModel = pipeline.fit(train)
```

```
In [34]:
```

```
best_predictions = pipelineModel.transform(test)
```

RMSE

```
In [35]:
```

```
rmse = evaluator.evaluate(best_predictions)  
print("RMSE on our test set: %g" % rmse)
```

```
RMSE on our test set: 1.57483
```

\$R^2\$

```
In [36]:
```

```
y_true = best_predictions.select("silicaFeed").toPandas()  
y_pred = best_predictions.select("prediction").toPandas()
```

```
In [37]:
```

```
r2 = sklearn.metrics.r2_score(y_true, y_pred)  
print('r2: {0}'.format(r2))
```

```
r2: 0.9464383854719923
```

Decision Trees

```
In [20]:
```

```
from pyspark.mllib.tree import DecisionTree  
from pyspark.ml.regression import DecisionTreeRegressor
```

```
In [21]:
```

```
# decisionTree train model
dt = DecisionTreeRegressor(labelCol="silicaFeed")
```

```
In [22]:
```

```
# Define a grid of hyperparameters to test:
paramGrid = ParamGridBuilder() \
    .addGrid(dt.maxDepth, [2,10]) \
    .build()

#This tells CrossValidator how well we are doing by comparing the true labels with predictions.
evaluator = RegressionEvaluator(metricName="rmse", labelCol=dt.getLabelCol(), predictionCol=dt.getPredictionCol())
# Declare the CrossValidator, which runs model tuning for us.
cv = CrossValidator(estimator=dt, evaluator=evaluator, estimatorParamMaps=paramGrid)
```

```
In [23]:
```

```
pipeline = Pipeline(stages=[vectorAssembler, vectorIndexer, cv])
```

```
In [33]:
```

```
pipelineTreeModel = pipeline.fit(train)
```

```
In [35]:
```

```
predictions = pipelineTreeModel.transform(test)
```

```
In [36]:
```

```
rmse = evaluator.evaluate(predictions)
print("RMSE on our test set: %g" % rmse)
```

```
RMSE on our test set: 4.55198
```

\$R^2\$

```
In [39]:
```

```
import sklearn.metrics
y_true = predictions.select("silicaFeed").toPandas()
y_pred = predictions.select("prediction").toPandas()
r2 = sklearn.metrics.r2_score(y_true, y_pred)
print('r2: {0}'.format(r2))
```

```
r2: 0.5525064735055689
```

FULL MODEL RMSE

RandomForest

MAKE SURE TO RELOAD THE DATA MINERDD BUT WITH FULL FEATURES

Reload mineRDD, drop date and % IronFeed, convert pandasDF to SparkRDD

Random Forest

In [18]:

```
pipelineRFModel = pipeline.fit(train)
predictions = pipelineRFModel.transform(test)
rmse = evaluator.evaluate(predictions)
print("RMSE on our test set: %g" % rmse)
```

```
RMSE on our test set: 4.26634
```

In [19]:

```
y_true = predictions.select("silicaFeed").toPandas()
y_pred = predictions.select("prediction").toPandas()
import sklearn.metrics
r2 = sklearn.metrics.r2_score(y_true, y_pred)
print('r2: {0}'.format(r2))
```

```
r2: 0.6069056057099547
```

Decision Tree

In [24]:

```
pipelineTreeModel = pipeline.fit(train)
predictions = pipelineTreeModel.transform(test)
rmse = evaluator.evaluate(predictions)
print("RMSE on our test set: %g" % rmse)
```

```
RMSE on our test set: 4.6144
```

In [25]:

```
y_true = predictions.select("silicaFeed").toPandas()
y_pred = predictions.select("prediction").toPandas()
r2 = sklearn.metrics.r2_score(y_true, y_pred)
print('r2: {0}'.format(r2))
```

```
r2: 0.5401485664010969
```

GBT

In [20]:

```
pipelineModel = pipeline.fit(train)
```

In [21]:

```
predictions = pipelineModel.transform(test)
```

In [22]:

```
rmse = evaluator.evaluate(predictions)
print("RMSE on our test set: %g" % rmse)
```

```
RMSE on our test set: 1.50837
```

In [24]:

```
import sklearn.metrics
y_true = predictions.select("silicaFeed").toPandas()
y_pred = predictions.select("prediction").toPandas()
r2 = sklearn.metrics.r2_score(y_true, y_pred)
print('r2: {0}'.format(r2))
```

```
PYTHON [12]:
```

```
r2: 0.9508639723695737
```

Best Predictions - GBT Full Model

We convert RDD back to Pandas for easy graphics

In [33]:

```
# silicaFeed: the True % Silica Feed from real data
# prediction: our predicted % Silica Feed
rdd_pred = predictions.select('silicaFeed', 'prediction')
```

In [34]:

```
# Spark to Pandas
predictions_pd = rdd_pred.toPandas()
```

Bar Plot

In []:

```
sns.set(style="whitegrid", font='Copperplate Gothic Bold', rc={'axes.facecolor':'cornflowerblue', 'figure.facecolor':'cornflowerblue'})
ax=sns.distplot(predictions_pd['silicaFeed'], color='green', hist_kws={ "alpha": 1}, kde=False)
plt.title('Distribution for TRUE % Silica Feed', fontsize=18)
plt.ylabel('Relative Frequency', fontsize=18)
plt.xlabel('TRUE Silica Feed [%]', fontsize=18)
plt.ylim((0,21000))
plt.show()
ax.figure.savefig("plotf1.png")
```

In []:

```
sns.set(style="whitegrid", font='Copperplate Gothic Bold', rc={'axes.facecolor':'cornflowerblue', 'figure.facecolor':'cornflowerblue'})
ax=sns.distplot(predictions_pd['prediction'], color='green', hist_kws={ "alpha": 1}, kde=False)
plt.title('Distribution for PREDICT % Silica Feed', fontsize=18)
plt.ylabel('Relative Frequency', fontsize=18)
plt.xlabel('PREDICTION Silica Feed [%]', fontsize=18)
plt.xlim((0,35))
plt.show()
ax.figure.savefig("plotf2.png")
```