

# MIE1628 Final Project Report

Group Number: 8

Student Name: Zijian Wei

Student Number: 1002276823

Submission Data: Dec 22<sup>nd</sup>, 2020

## Contents

1. Dataset Selection and Size .....	3
2. Problem Definition.....	3
2.1 Background Information.....	3
2.2 Scoping .....	3
2.3 Objective .....	4
3. Literature Review.....	4
4. Data Cleaning and Processing .....	4
4.1 Target Variable Description .....	4
4.2 Data Cleaning and Feature Preparation.....	4
5. Model Development.....	5
5.1 Simple Moving Average (SMA) Model.....	5
5.2 ARIMA Model.....	6
5.3 Linear Regression Model .....	6
5.4 Random Forest.....	8
5.5 Gradient Boosted Tree (GBT) .....	8
6. Result Discussion .....	9
7. Future Improvement and Personal Contribution .....	9
8. References .....	10
9. Appendix .....	11
9.1 Literature Review Plot.....	11
9.2 Comparison of Predictions and Actual Values for 5 Models .....	12
9.3 Table of SMAPE for 5 Models .....	14
9.4 Code (see following pages) .....	15

## 1. Dataset Selection and Size

My group chose Gbatteries project and the dataset is offered rather than selected. The reason of choosing such project is that we want to attempt an actual problem in industry and apply what we have learned in this course. With this project, we are solving a challenge that can bring real benefits to industrial process, which is exciting. We also think that such dataset offered has enough complexity for us to work with.

The raw data set has 32266 records in total for 48 batteries, and 10 fields. The file size of the raw data in .csv format is 2.8 MB. The size of dataset is small. The details of dataset are discussed later after defining the problem.

## 2. Problem Definition

### 2.1 Background Information

Nowadays, rechargeable battery plays an important role in many different fields for its ability of storing electrical energy and recharging. One of the measurements of performance of a rechargeable battery is its durability, which can be represented by the number of time it can be recharged while still being able to release enough electrical energy.

In this Project, the metric of durability is the number of cycles that a battery is charged and then discharged before reaching its end of life (EOL). There are two types of cycling: regular cycle and RPT cycle. RPT cycles are specialized cycles to measure the retentive capacity of battery. EOL happens when the RPT discharge capacity crosses 80% of the nominal value of capacity.

Testing durability of batteries with cycling is a lengthy process and can take months. Therefore, being able to forecast the RPT discharge capacity accurately can greatly shorten the durability testing process.

### 2.2 Scoping

The retentive capacity of battery has a general decreasing trend as cycle number increases. Although RPT cycle is specialized cycle to measure the retentive capacity accurately, it can sometimes be volatile. In order to determine the EOL correctly, the trend of RPT discharge capacity is focused regardless of volatility. Notice that volatility is complex and can differ significantly depending on the set of charging parameters, which are not given in the data and are different for each unique battery id. The dataset is of small size and would not satisfy the need of investigating volatility. In addition, the trend is eventually the key of predicting EOL.

Moreover, behavior of battery changes significantly after EOL. Since the main interest is to predict the number of cycles until EOL takes place, the prediction of RPT discharge capacity before EOL taking place is focused. The data of cycles before EOL are taken for the analysis and modeling process. If there is need of predicting battery behavior after EOL, another model would be suggested.

### 2.3 Objective

Develop a model that predict RPT cycle discharge capacity at various steps into the future before EOL.

## 3. Literature Review

There are many models can work with time series data. For univariate models, there are different types of regression [1] and ARIMA [2]. For multivariable model, ensemble learning [3] can be performed including decision tree, gradient boosted tree, random forest, and proposed ensemble model. Moreover, client mentions in the project proposal that fitting a square root function has been applied.

Through cycling, degradation of the lithium-ion batteries, battery type in this project, has in fact been investigated. As shown in Figure 1 in Appendix, the capacity has a relationship with respect to cycle that can be represented by linear. Furthermore, when higher and higher elevated charge voltage is applied, the relationship becomes more and more quadratic-like which is concave down, as shown in Figure 2 in Appendix. [4]

Linearity is potentially a good local or even global approximation of the relationship between capacity and number of cycles.

## 4. Data Cleaning and Processing

### 4.1 Target Variable Description

The target variable in this project is RPT discharge capacity, which is considered as an accurate measurement of retentive capacity of battery. RPT discharge capacity has a general decreasing trend with some volatility, with respect to cycle number. In addition, RPT discharge capacity is not recorded in a continuous way. Each RPT cycle takes place after a series of regular cycles (approximately 20 cycles). Sometimes consecutive RPT cycles take place, with early RPT cycles stabilizing the battery and the last RPT cycle being the actual measurement, as mentioned in the provided project proposal.

### 4.2 Data Cleaning and Feature Preparation

The raw data has each record corresponding to a state of charge (SOC) within a cycle. A cycle has 4 SOC (thus 4 records) and each SOC has 3 measurements as features. The

set of 4 records corresponding to same cycle are combined as one, such that each record has 12 features along with discharge capacity.

Since the target variable is RPT discharge capacity, data of all batteries that only have regular cycles are eliminated. There are 22 batteries that have RPT cycles. The remaining data have no null value. Some cycles have data for less than 4 SOC, but they belong to the batteries that have no RPT cycle. As mentioned in Section 4.1, there are some sets of consecutive RPT cycles. For these consecutive sets of RPT cycles, only the last RPT cycle of each consecutive set is retained.

The dataset is then reformatted such that each record contains the discharge capacity and 12 features for both regular and RPT as if both measurements are performed at same cycle. Since a cycle can only be either regular or RPT, this reformatting results in a lot of null values. In addition, the dataset has some cycles missing. These null values are handled by performing linear interpolation and linear extrapolation (for initial cycles), with the assumption of local linearity, which is suggested by client. Then, a continuous time series dataset is obtained, with cycle number being the time stamp.

Lastly, there are some batteries have too few RPT cycles. Some of the models require first 100 or even 110 cycles for training to forecast the later cycles. Therefore, there are 2 batteries eliminated again because they do not have RPT cycles after cycle 100.

The dataset after cleaning and processing consists of 20 batteries. The data for each battery is a set of continuous time series data. Each record has discharge capacity and 12 features for both regular and RPT cycles.

## 5. Model Development

### 5.1 Simple Moving Average (SMA) Model

SMA is the simplest model. It is univariable model. The future values are forecasted by taking mean of past values. The model can be optimized by choosing size of window of taking mean. However, there is a crucial limitation of this model: it assumes the future values are close to past values. In this scenario, the RPT discharge capacity has a negative trend. Using data of early cycles that locate at upper part in terms of magnitude would confine the predictions also inside the upper part. Depending on size of window, the prediction eventually will converge, because it is simply taking mean of past values repeatedly. A small window size of 5 is chosen, such that the model converging at value close to the last RPT cycle, as shown in Figure 3 in Appendix.

Due to its limitation, SMA model is not expected to give accurate prediction of future RPT discharge capacities. However, it can be used as a baseline model. As will be

exhibited in the Model Result component, the SMAPEs of SMA are in fact not large, and growing as forecasting time horizon increases. This not only demonstrates an underlying trend of RPT discharge capacity cycle, but also implies that the rate of change is small. The later models are expected to result in smaller SMAPEs, or the learned trends of batteries differ from the actual trends by a lot.

## 5.2 ARIMA Model

While SMA is simply taking mean of past values and does not analyze the trend of data, ARIMA is a univariate model that allows analysis with orders of differencing, which is important detrending process for data with trend (nonstationary data). Therefore, ARIMA model is attempted after realizing the limitation of SMA. ARIMA is a complex autoregressive model that involves 3 key parameters:  $p$ ,  $d$ , and  $q$ , which are lag order, degree of differencing, and order of moving average. It has been a challenge that there is no embedded ARIMA model in Spark library like some other models (linear regression, random forest, and gradient boosted tree). A piece of code is found online and studied. By tuning the three parameters using first 100 cycles of the dataset, a set of parameters is obtained that can produce model with linear negative trend (the degree of differencing is 1). The first 100 cycles are used to train ARIMA model to forecast future values for each battery. As shown in Figure 4 in Appendix, the linear negative trend produced by ARIMA model is a good approximation of the actual negative trend.

Since the actual RPT cycles are sparse and most RPT data are filled by linear interpolation, window sizes for lagging and moving average have some degree of flexibility without changing the model significantly. The main challenge of ARIMA model is to determine the order of differencing parameter,  $d$ . By attempting different values of  $d$ , it is found that only first order of differencing can make the data stationary. This also implies that the underlying trend of RPT discharge capacity is best represented by linearity, among different polynomial order. The limitation of ARIMA is that it eventually converges to a linear regression model, thus has same limitation as linear regression, as will be discussed later. In addition, ARIMA model is complicated and a piece of code online is studied and used, therefore it is hard to implement additional techniques into it. For instance, the outlier removing for linear regression.

## 5.3 Linear Regression Model

As first order differencing can make the data stationary when developing ARIMA model, linear regression is potentially a good model. The linear regression model is univariate and take first 110 cycles to train for a linear relationship between RPT discharge capacity and cycle number. Actual discrete RPT discharge capacity data are used for developing linear regression model instead of the continuous data after interpolation and extrapolation. In fact, the process of interpolation and extrapolation is under the

assumption of local linearity. Since an RPT cycle takes place approximately every 20 cycles (can be less frequent), there are in fact very few points for training. In addition, as mentioned for SMA model, the rate of change of the negative trend is small. Therefore, any outlier due to volatility of RPT discharge capacity can impact the linear regression model significantly. With a regular linear regression that simply fits on the early RPT cycles with cycle number less than 110, the performance is poor, as showed in Table 1. The effect of outliers can be lessened to some degree by taking in more points, but a RPT cycle is expensive can require another 20 regular cycles.

The challenge of this model is to deal with outlier with such few data points. Several steps have been designed to eliminate significant outliers and make the line of linear regression model as appropriate as possible. The first step is setting upper and lower bounds for RPT discharge capacity cycles, these thresholds are very loose and just for outliers that are obviously problematic. The second step is to fit a linear regression model first and eliminate outliers based on deviation (threshold value is tuned by using two batteries as validation set) of the training data points from the line. Second step is repeated for multiple times if multiple outliers are found. However, since the data points are too few, the second point can eliminate significant outliers and might eliminate good data points when there are more than one outlier data points among the few training data points. The RPT discharge capacity measured from last RPT cycle of consecutive cycles are more reliable and less volatile than single RPT cycle, additional weights are given to training points corresponding to consecutive RPT cycles by replicating the data points twice. The third step not only help determining outliers in second step, but also lessen the effect of insignificant outliers, because those RPT discharge capacities are more assured. As showed in Table 1, the linear regression model obtained after these three steps are implemented has a much better performance.

Time Horizon	50	100	150	200	250	300	350
Regular Lin. Reg.	5.83	8.18	9.33	10.92	8.69	22.19	13.50
Enhanced Lin. Reg.	2.36	3.12	1.65	1.88	1.94	1.93	3.86

Table 1. Performance Comparison between Regular and Enhanced Lin. Reg. model

Linear regression model has two main limitation. First one is that it is trying to ignore the volatility of RPT discharge capacity, which is hard to accomplish with the amount of data given. Second limitation is shown in Figure 5 in Appendix, a line is indeed not a perfect fit for some batteries (in fact, it is perfect fit for some) and there is secondary term left over. Error accumulates from this secondary term as forecasting time horizon increases.

### 5.4 Random Forest

While the previous three models are all univariate, random forest is used as a multivariable approach. The challenge of random forest is that the data is nonstationary. When use early 100 cycles to train a random forest, the model would not recognize the later cycles because their values are all lower (except for some due to outliers) than the earlier training data. The predictions would converge to a flat line as SMA model. While ARIMA demonstrates that first order differencing can make the data stationary, the same technique is applied for random forest. However, detrending for ARIMA can be done by setting parameter  $d$ , while the detrending needs to be implemented for random forest. Random forest uses all 24 SOC features along with regular discharge capacity as features. First order differencing is applied to all features. The target variable also becomes the difference of RPT discharge capacity between two consecutive cycles. Feature selection has not been done for this model. Since first order of differencing is used, this model exhibits the limitation of linear regression.

### 5.5 Gradient Boosted Tree (GBT)

Challenge of GBT is also that it can only be applied to stationary time series data. Instead of detrending, GBT model turns to an approach different from previous four models. GBT model is designed to predict slope and intercept of linear model of RPT discharge capacity using only regular cycle data. Since frequency of actual regular cycle data is much higher than actual RPT cycle data, enough information can be obtained from regular cycle, with much less early cycles required.

Since the regular cycle SOC features have small variation during the cycling process of each battery, mean of each for first 50 cycles can represent the regular cycle SOC data of a battery. In addition, the regular cycle discharge capacity also has a linear-like relationship with respect to cycle number. Therefore, the slope and intercept of linear model of regular cycle discharge capacity derived from first 50 cycles are extracted as 2 features. These 14 features in total from regular cycles of data are used to train GBT for each of the two target variables: slope and intercept of linear model of RPT discharge capacity. The reformatted data has 20 records, each corresponding to a battery. Two batteries are kept as testing set and 18 batteries are split into training and validation set with ratio of 7:3. The labels (2 target variables) of validation and training set are derived by fitting linear regression to RPT discharge capacity data for each battery.

Feature selection is done with PCA. The performance on validation set is the best when using the top 5 features are output from PCA. In addition, other hyperparameters for GBT trees: `maxIter`, `maxDepth`, `stepSize`, are tuned. The tuning process accomplished by `TrainValidationSplit` with `paramGrid`.

Limitation of this model exhibits the limitation of linear regression model. In addition,



there are not enough data for different batteries for this approach. Moreover, the outlier removing process in linear regression is not adopted for extracting features (not working well with whole set of data for a battery, which might need implementation of RANSAC).

## 6. Result Discussion

The table summarizing the SMAPE for 5 models with different time horizons are placed in Appendix (section 9.3). ARIMA, linear regression, GBT provide reasonably good approximation of the actual data. In specific, linear regression has the best performance, which is expected. Since other model use linear interpolation, linear extrapolation, and first order differencing for detrending, and GBT is predicting a linear model of RPT discharge capacity with regular cycle data. The models all eventually end up with a linear model and exhibit the limitation of linear regression model, which can be shown by viewing the plots in Appendix (section 9.2) Moreover, as mentioned in linear regression model section, outliers can be a big performance degradation factor when the actual RPT data points are not many. With the outlier removing steps, linear regression model can find the more appropriate line for the batteries. SMA has relatively poor performance because, as mentioned, the model does not involve the trend of data. The reason of GBT having poor performance is that it has not yet been tuned and feature selection also has not yet been done, which causing it to have extremely bad performance for some batteries. Besides SMAPE, it is important to notice that GBT requires much less cycles (50 cycles in this scenario) to make predictions into the future because it makes use of regular cycle data. Making use of regular cycle data is also an expectation from the client.

## 7. Future Improvement and Personal Contribution

Since the models all exhibit limitation of linear regression, it is important for dealing with limitation of linear regression. A secondary term can be added to the fitting model on top of linear equation to lessen the accumulation of error over long time horizon forecasting. Meanwhile, if more data of RPT cycles are given, the volatility of RPT cycles might be modeled, which indeed also depends on how the RPT cycles are arranged (consecutive RPT cycles give more stable measurements). In addition, if data of more batteries are given, the performance of GBT is also expected to be increased (too few batteries currently). Moreover, hyperparameter tuning and feature selection need to be completed for random forest model. Furthermore, RANSAC might be implemented for the whole data set of a battery (the outlier removing has been designed for only earlier cycles), such that models can produce more accurate approximation of the trend.

Throughout this project, I have participated in team data cleaning and processing, and done individual data processing to understand the data better. I developed SMA and linear regression model, and given suggestions to other teammates constructing their models.

## 8. References

- [1] Ostrom, C. W. (1990). *Time Series Analysis: Regression Techniques* (Vol. 9).
- [2] Krome, C., Sander, V. Time series analysis with apache spark and its applications to energy informatics. *Energy Inform* **1**, 40 (2018). <https://doi.org/10.1186/s42162-018-0043-1>
- [3] Galicia, A., Talavera-Llames, R., Troncoso, A., Koprinska, I., & Martínez-Álvarez, F. (2018, October 12). Multi-step forecasting for big data time series based on ensemble learning. Retrieved December 22, 2020.
- [4] BU-808: How to Prolong Lithium-based Batteries. (n.d.). Retrieved December 22, 2020, from [https://batteryuniversity.com/learn/article/how\\_to\\_prolong\\_lithium\\_based\\_batteries](https://batteryuniversity.com/learn/article/how_to_prolong_lithium_based_batteries)

## 9. Appendix

### 9.1 Literature Review Plot

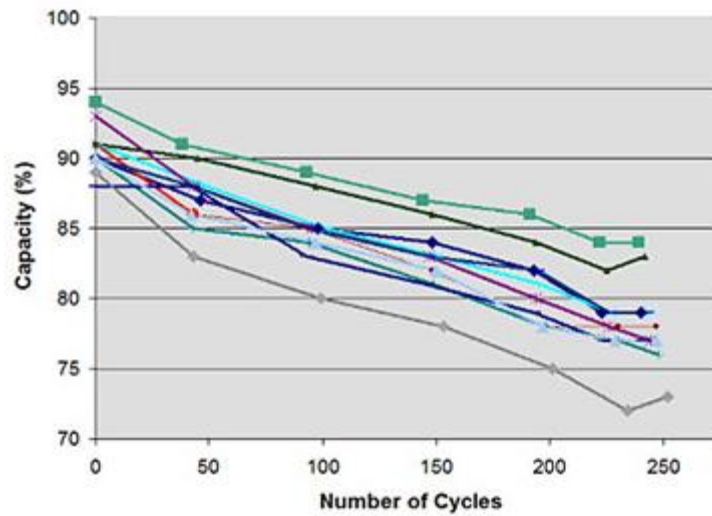


Figure 1. Ageing of Lithium-ion Battery

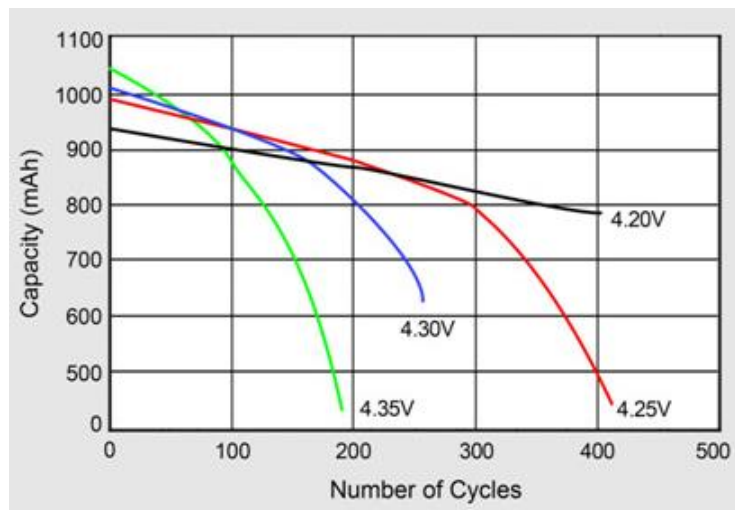


Figure 2. Ageing of Lithium-ion Battery with Elevated Charge Voltage

## 9.2 Comparison of Predictions and Actual Values for 5 Models

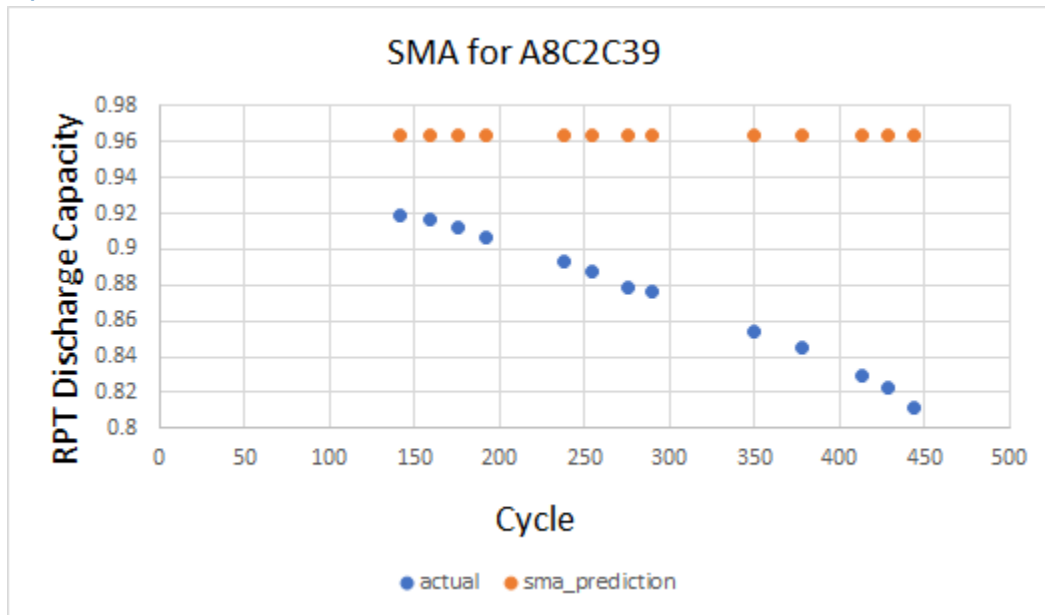


Figure 3. Predictions vs. Actual Values for SMA Model

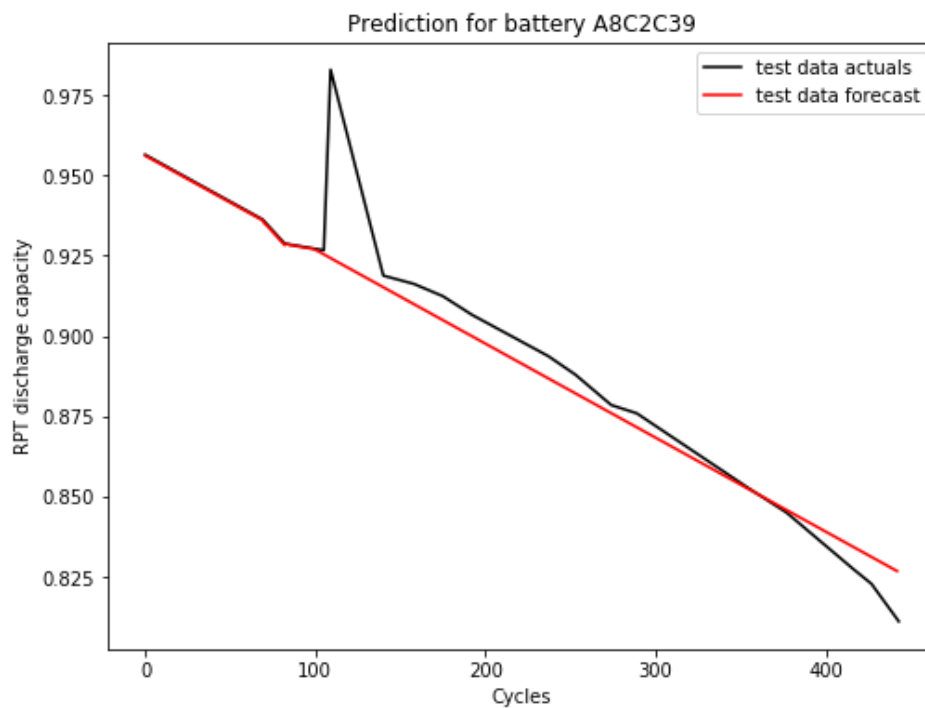


Figure 4. Predictions vs. Actual Values for ARIMA Model

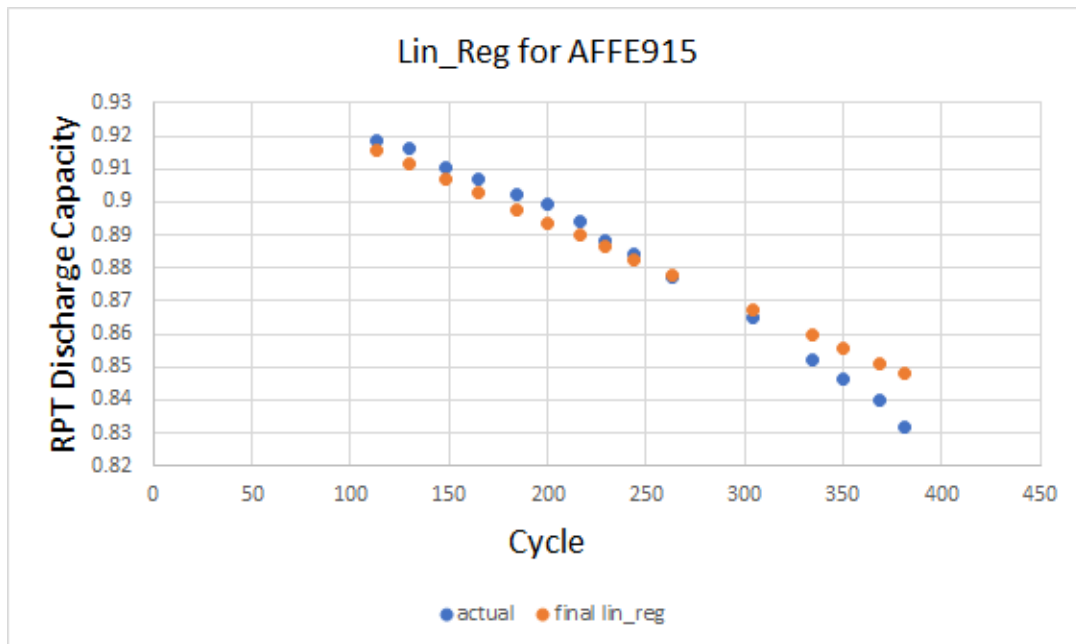


Figure 5. Predictions vs. Actual Values for Lin. Reg. Model

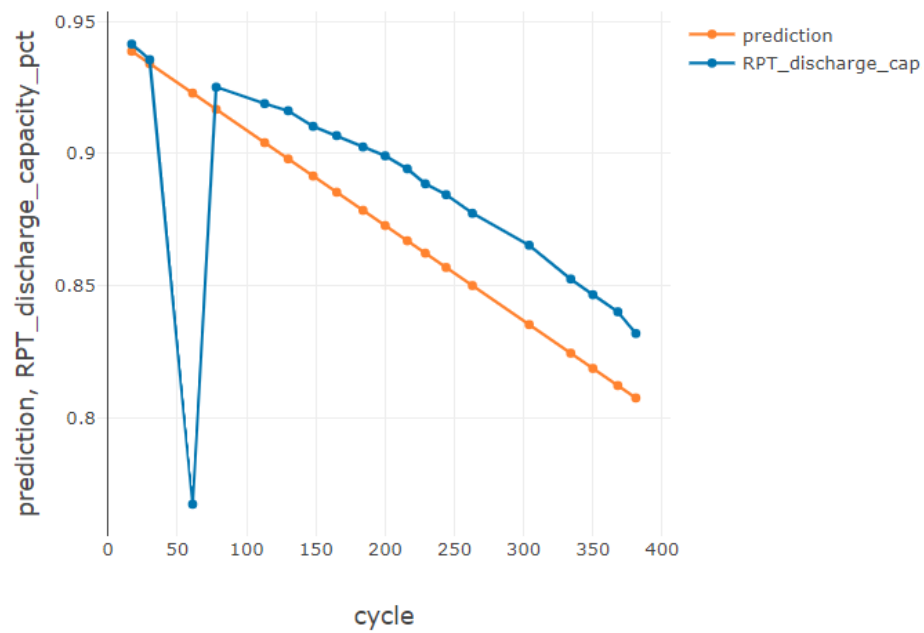


Figure 6. Predictions vs. Actual Values for Random Forest Model

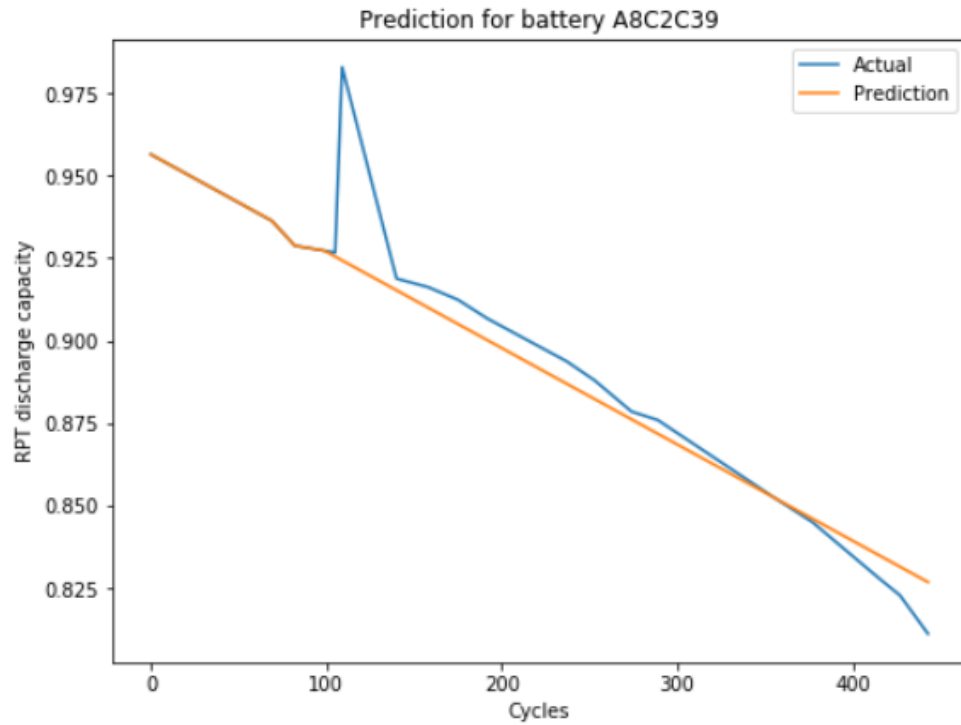


Figure 7. Predictions vs. Actual Values for GBT Model

### 9.3 Table of SMAPE for 5 Models

Time Horizon	50	100	150	200	250	300	350
SMA	3.30	4.55	6.43	7.62	8.69	10.08	12.17
ARIMA	2.43	4.51	5.55	6.92	7.21	7.39	7.39
Linear Regression	2.36	3.12	1.65	1.88	1.94	1.93	3.86
GBT	1.88	5.85	6.63	6.83	8.24	8.23	9.28
Random Forest	11.04	20.52	24.14	26.28	27.15	27.75	27.76

Table 2. SMAPE for 5 Models

9.4 Code (see following pages)

```

1  # -*- coding: utf-8 -*-
2  """group8_notebook_clone_1.ipynb
3
4  Automatically generated by Colaboratory.
5
6  Original file is located at
7  https://colab.research.google.com/drive/1fwyGSZJZ-USeQvPqzeakNz0QYouXsLdA
8  """
9
10 '''
11 MIE1628: Summative Course Project; GBatteries Team 8
12 Team members: Junzhe Liu (1002416557), Yue Niu (1002222613), David Wang (1001393351), Zijian Wei
   (1002276823), Yixin Xiao (1002115221)
13
14 GBatteries is an advanced battery technology company pioneering an alternative way to charge lithium-ion (Li
   -ion) batteries. Their charging method uses adaptive pulses instead of the conventional charging
   protocol CCCV used as the standard Li-ion charging method. Their technology enables ultra-fast charge
   without compromising the health of the battery.
15
16 The goal for this project is to develop a model to predict the RPT discharge capacity at various steps into
   the future. Labeled data is provided for thousands of cycles; this notebook will explore modeling
   solutions that will eventually play a role in a monitoring solution that will be deployed in the
   GBatteries battery testing facility. More information about the company, project requirements, and
   targets can be found in the project report (completed on a per-contributor basis).
17
18 The general steps undertaken in this notebook will be as follows:
19 * Cleaning and other transformations
20 * Time series visualization
21 * Removing non-stationarity or seasonality (not a problem in this problem domain)
22 * Partitioning and vector of past values
23 * Train and evaluate various models; model optimization
24 * Discuss results
25 '''
26
27 #####
28 # i. IMPORTS + LOADING DATASET
29 #####
30 # !pip install --upgrade pip
31 # !pip install fbprophet
32 # !pip install koalas
33 # !pip install missingno
34 # !pip install fancyimpute
35 # !pip install flint
36
37
38 from pyspark import since, SparkContext
39 from pyspark.sql.column import Column, _to_java_column, _to_seq
40 from pyspark.sql import SparkSession, Window
41 from pyspark.sql.functions import countDistinct
42 from pyspark.sql import functions as F
43 from pyspark.sql.types import *
44 from pyspark.sql.types import StructType, IntegerType, DateType, StringType, DoubleType, StructField,
   FloatType
45 from pyspark.sql.functions import pandas_udf, PandasUDFType, collect_list, struct, create_map
46 from pyspark.sql.functions import sum, max, col, avg, concat, lit, isnan, when, count, first, last
47 from pyspark.ml.regression import LinearRegression
48 from pyspark.ml.feature import VectorAssembler
49 from pyspark.ml.linalg import Vectors
50 from pyspark.sql import Row
51 from pyspark.sql.functions import row_number, lit
52 from pyspark.sql.window import Window
53 from pyspark.ml.regression import GBRegressor
54 from pyspark.mllib.linalg import Vectors
55 from pyspark.mllib.linalg.distributed import RowMatrix
56 from pyspark.ml.feature import PCA
57 from pyspark.ml.evaluation import RegressionEvaluator
58 from pyspark.ml.tuning import CrossValidator, ParamGridBuilder, TrainValidationSplit
59 from pyspark.ml import Pipeline

```



```

60
61 # import warnings, sys, os, gc, time
62 from itertools import chain
63 # import missingno as msno
64 import numpy as np
65 import pandas as pd
66 # from fbprophet import Prophet
67 # from statsmodels.tsa.stattools import adfuller, acf, pacf
68
69
70
71 # load data; dbfs:/FileStore/shared_uploads/gbatteries_team8@outlook.com/dataset.csv
72 schema = StructType([
73     StructField("_c0", IntegerType()),
74     StructField("battery_model", StringType()),
75     StructField("battery_id", StringType()),
76     StructField("cycle", IntegerType()),
77     StructField("cycle_type", StringType()),
78     StructField("charge_duration_sec", IntegerType()),
79     StructField("discharge_capacity_pct", FloatType()),
80     StructField("soc_region", IntegerType()),
81     StructField("feature_1", FloatType()),
82     StructField("feature_2", FloatType()),
83     StructField("feature_3", FloatType())
84 ])
85
86 df = spark.read.format("csv").load("dbfs:/FileStore/shared_uploads/gbatteries_team8@outlook.com/dataset.csv"
87     , header = True, schema = schema)
88
89 print(df.show(2))
90
91 # get schema of our data
92 # df.printSchema()
93
94 #####
95 # 1. DATA CLEANING, PREPROCESSING
96 #####
97
98 # display count of nans; the data has been thoroughly cleaned and prepared by the GBatteries point
99 # of contact
100 print('validate that there are no null records:')
101 # df.select([count(when(isnan(c), c)).alias(c) for c in df.columns]).show()
102 df.select([count(when(isnan(c) | col(c).isNull(), c)).alias(c) for c in df.columns]).show()
103
104 '''
105 Let's move onwards to EDA and preprocessing. Seran Thirugnanam from GBatteries has outlined in the
106 project brief + introductory touch point meeting that the data has been thoroughly cleaned, and
107 has requested the project team to instead focus on improving modeling capabilities
108 '''
109
110 print('sample frequency of RPT vs regular cycles for a sample battery (id)')
111 df.filter('battery_id == "A0047B2").select('cycle', 'cycle_type').distinct().orderBy('cycle').show(10)
112
113 '''
114 Note that some cycles have missing information (e.g. battery A0047B2 is missing cycle 3 in the below example
115 )
116 '''
117
118 # show distinct (unique) values for soc_region
119 print('unique soc_region values')
120 df.select("soc_region").distinct().show()
121
122 # show distinct (unique) values for battery_model
123 print('unique battery models contained in dataset')
124 df.select("battery_model").distinct().show()
125
126 # remove records for cycles with soc less than 4
127 # df = df.filter(df['battery_id'] != 'A2A5795').filter((df['battery_id'] != 'ACF5ADD') | (df['cycle'] !=
128     215))

```

```

125 # xxx = df.groupBy("battery_id").agg(countDistinct("soc_region"))
126 # xxx.filter(xxx['count(soc_region)'] == 2).show()
127 # df.filter(df['battery_id'] != 'ACF5ADD').show(444)
128
129 '''
130 Note that the GBatteries representative (Seran) commented on cycle types varying between charges, with RPT
    cycle types
131 being exceedingly infrequent as compared to Regular cycles. RPT is a standardized score; RPT cycles are what
    we care about
132 here, but they are measured once every ~50 cycles or so.
133
134 We can interpolate RPT metrics for every charge cycle, and Regular metrics for every charge cycle...
135 * I.e. have rows for cycle 1 ... n for both RPT and Regular
136 * These features can then be used for modelling (e.g. Random Forest requires populated values for every
    cycle)
137
138 As a first step, we also need to pivot the table so that SOC region 1...4 and their corresponding features
    are represented
139 as columns (objective: one row per battery cycle, with all features captured at a column level).
140
141 Objective:
142 | SOC1_f1 | SOC2_f1 | SOC3_f1 | SOC4_f1 | SOC2_f1 | SOC2_f2 | ... | SOC4_f3 |
143
144 Let's start with pivoting our data:
145 '''
146
147 # replace soc_region values ahead of pivoting
148 soc_dict = {1:'SOC1', 2:'SOC2', 3:'SOC3', 4:'SOC4'}
149 mapping_expr = create_map([lit(x) for x in chain(*soc_dict.items())])
150 df = df.withColumn('soc_region', mapping_expr[df['soc_region']])
151
152 # pivot soc_region with all three features
153 dff = df.groupBy('battery_model', 'battery_id', 'cycle', 'cycle_type', 'charge_duration_sec', '
    discharge_capacity_pct').pivot('soc_region').agg(F.first('feature_1').alias('feature_1'), F.first('
    feature_2').alias('feature_2'), F.first('feature_3').alias('feature_3'))
154
155 # order by battery_id and cycle (ascending)
156 order_cols = ['battery_id', 'cycle']
157 dff = dff.orderBy(order_cols, ascending = True)
158
159 dff.show(10)
160
161 df.show()
162
163 dff_2 = dff
164 window_next = Window.partitionBy('battery_id').orderBy('cycle')
165 dff_2 = dff_2.withColumn('next_cycle', F.lead('cycle', 1).over(window_next))
166 dff_2 = dff_2.withColumn('next_cycle_type', F.lead('cycle_type', 1).over(window_next))
167 dff_2 = dff_2.filter(dff_2['next_cycle'].isNull() | ((dff_2['cycle_type'] == 'Regular') | (dff_2['next_cycle']
    - dff_2['cycle'] != 1) | (dff_2['next_cycle_type'] != 'RPT')))
168
169 # get dff count to compare to dff_2 (after operations)
170 print(dff.count())
171
172 # expected 8 less, correct
173 print(dff_2.count()) # > 478
174
175 dff = dff_2.drop('next_cycle', 'next_cycle_type')
176
177 dff_2.show(200)
178
179 dff.coalesce(1).write.format("com.databricks.spark.csv").option("header", "true").save("dbfs:/FileStore/
    uninterpolated_1219v1.csv")
180
181 dff.show()
182
183 '''
184 Next, we move onto interpolating RPT AND Regular metrics for every charge cycle... i.e. have rows for

```

```

185 cycle 1 ... n for both RPT and Regular cycle types
186 '''
187
188 #####
189 # 2. INTERPOLATING ALL REGULAR + RPT CYCLES
190 #####
191
192 # Descriptive statistics for both RPT and Regular cycles
193 print('RPT charge_duration_sec descriptive statistics:')
194 dff.filter(df["cycle_type"] == "RPT").select('charge_duration_sec').describe().show()
195
196 print('Regular charge_duration_sec descriptive statistics:')
197 dff.filter(df["cycle_type"] == "Regular").select('charge_duration_sec').describe().show()
198
199 # we now divide our dataset into two DataFrames: one for RPT cycles and the other for Regular cycles
200 reg = dff.filter(df['cycle_type'] == 'Regular')
201 rpt = dff.filter(df['cycle_type'] == 'RPT')
202
203 #data.select([count(when(isnan(c), c)).alias(c) for c in data.columns]).show()
204 # the records for selected batteries all have 4 soc, this section can be commented
205 # get batteries
206 lt_4_soc = df.groupBy("battery_id", "cycle_type", "cycle").agg(count("*"))
207 lt_4_soc = lt_4_soc.filter(F.col('count(1)').isin([0,1,2,3]))
208 # lt4_soc_ids = [x.battery_id for x in lt_4_soc.select('battery_id').distinct().collect()]
209
210 # filter out batteries with <4 SOC_regions
211 # reg = reg.filter(~col('battery_id').isin(lt4_soc_ids))
212 # rpt = rpt.filter(~col('battery_id').isin(lt4_soc_ids))
213
214 '''
215 Because we're interested in forecasting out RPT cycles, it doesn't make sense to
216 analyze batteries (battery_id) where we have no RPT cycle values. The operation
217 below will exclude batteries that don't have any RPT cycle records.
218 '''
219 # the batteries we selected have RPT cycles, this part is commented
220 # get DataFrame of unique batteries that DO have RPT cycle records
221 unique_rpt_ids = rpt.groupBy("battery_id").agg(countDistinct("cycle"))
222
223 # apply a filter on the Regular cycles DataFrame to exclude only the
224 # aforementioned batteries
225 rpt_ids_arr = [str(row['battery_id']) for row in unique_rpt_ids.collect()]
226 reg = reg.filter(F.col('battery_id').isin(rpt_ids_arr))
227
228 # compare to confirm we have the same count of battery ids in both DataFrames
229 unique_reg_ids = reg.groupBy("battery_id").agg(countDistinct("cycle"))
230 print('# reg batteries:', unique_reg_ids.count(), '# rpt batteries:', unique_rpt_ids.count())
231 battery_id_list = [str(i) for i in np.array(unique_reg_ids.collect())[:,0]]
232
233 '''
234 The sample data generally consists of batteries undergoing a few hundred cycles of charge-discharge (max
235 cycle = 563).
236 We need cycles 1 ... n_i for each battery i; to do this, we construct new DataFrames for Regular and RPT
237 cycles, then
238 interpolate/fill NaN values (i.e. rows in the new DataFrame corresponding to missing cycles in the original
239 DataFrame).
240
241 Note here that the dataset exclusively looks at battery_model == 'A'; we can exclude this column since it
242 conveys no
243 information.
244 '''
245 #reg.groupBy('battery_id').agg(F.max('cycle')).show()
246 #rpt.groupBy('battery_id').agg(F.max('cycle')).show()
247
248 # get max cycle counts per battery
249 max_cycle_counts = dff.groupBy('battery_id').agg(F.max('cycle'))
250 max_cycle_counts = max_cycle_counts.filter(col("battery_id").isin(battery_id_list))
251 max_counts_arr = np.array(max_cycle_counts.collect())

```

```

249 # initialize lists we will use to construct our new DataFrame (new DataFrame to account for missing cycles)
250 id_arr = []
251 cyc_arr = []
252
253 # create lists to initialize "all cycle" DataFrames with
254 for id, count in max_counts_arr:
255     for cycle_no in range(1, int(count) + 1):
256         id_arr.append(str(id))
257         cyc_arr.append(cycle_no)
258
259 # create full DataFrame with all cycle information (f => stands for full)
260 f_reg = sqlContext.createDataFrame(zip(id_arr, cyc_arr), schema=['battery_id', 'cycle'])
261 f_rpt = sqlContext.createDataFrame(zip(id_arr, cyc_arr), schema=['battery_id', 'cycle'])
262
263 # populate existing data; we now have
264 f_reg = f_reg.join(reg, on = ['battery_id', 'cycle'], how = 'left')
265 f_rpt = f_rpt.join(rpt, on = ['battery_id', 'cycle'], how = 'left')
266
267 # sort; easier on the eyes
268 f_reg = f_reg.orderBy(order_cols, ascending = True)
269 f_rpt = f_rpt.orderBy(order_cols, ascending = True)
270
271 # display results
272 # print('Regular cycle DataFrame with missing values:')
273 # f_reg.show(5)
274 # print('RPT cycle DataFrame with missing values:')
275 # f_rpt.show(5)
276
277 # print missingness of RPT cycle DataFrame (recall from before, RPT cycles are far and many between; these
278 # are only measure every ~50 cycles according to Seran)
279 print('RPT cycle DataFrame missingness pct:', (1 - rpt.count()/f_rpt.count())*100,'%')
280
281 # drop battery_model, which are all 'A'
282 f_reg = f_reg.drop('battery_model')
283 f_rpt = f_rpt.drop('battery_model')
284
285 '''
286 There aren't any out-of-the-box solutions for advanced imputation (e.g. MICE, KNN) solutions in PySpark.
287 At the recommendation of GBatteries, we'll look at the following for both f_rpt (RPT cycles) and f_reg
288 (Regular cycles):
289
290 1. Linear interpolation to fill in NULL discharge_capacity_pct values
291 2. Forward filling in NULL charge_duration_sec - i.e. propagating the last valid observation forward
292 (rationale: the charge_duration_sec values reasonable similar within all RPT cycle records, and within
293 all Regular cycle records).
294
295 '''
296
297 def coalesce(*cols):
298     """
299     Returns the first non-null argument if exists. Otherwise, null.
300     source: https://spark.apache.org/docs/2.1.0/api/python/_modules/pyspark/sql/functions.html
301
302     :param *cols: a variable number of columns
303
304     :returns: returns the first column that is not null.
305     """
306     sc = SparkContext._active_spark_context
307     jc = sc._jvm.functions.coalesce(_to_seq(sc, cols, _to_java_column))
308     return Column(jc)
309
310 # linear interpolation window transformation to impute discharge_capacity_pct
311 def fill_linear_interpolation(df, id_cols, order_col, value_col):
312     """
313     Apply linear interpolation to dataframe to fill gaps.
314
315     :param df: Spark DataFrame
316     :param id_cols: string or list of column names to partition by the window function

```

```

317 :param order_col: column to use to order by the window function
318 :param value_col: column to be filled
319
320 :returns: spark dataframe updated with interpolated values
321 """
322 # create row number over window and a column with row number only for non missing values
323 w = Window.partitionBy(id_cols).orderBy(order_col)
324 new_df = df.withColumn('rn', F.row_number().over(w))
325 new_df = new_df.withColumn('rn_not_null', F.when(F.col(value_col).isNotNull(), F.col('rn')))
326
327 # create relative references to the start value (last value not missing)
328 w_start = Window.partitionBy(id_cols).orderBy(order_col).rowsBetween(Window.unboundedPreceding, -1)
329 new_df = new_df.withColumn('start_val', F.last(value_col, True).over(w_start))
330 new_df = new_df.withColumn('start_rn', F.last('rn_not_null', True).over(w_start))
331
332 # create relative references to the end value (first value not missing)
333 w_end = Window.partitionBy(id_cols).orderBy(order_col).rowsBetween(0, Window.unboundedFollowing)
334 new_df = new_df.withColumn('end_val', F.first(value_col, True).over(w_end))
335 new_df = new_df.withColumn('end_rn', F.first('rn_not_null', True).over(w_end))
336
337 # create references to gap length and current gap position
338 new_df = new_df.withColumn('diff_rn', F.col('end_rn') - F.col('start_rn'))
339 new_df = new_df.withColumn('curr_rn', F.col('diff_rn') - (F.col('end_rn') - F.col('rn')))
340
341 # calculate linear interpolation value
342 lin_interp_func = (F.col('start_val') + (F.col('end_val') - F.col('start_val')) / F.col('diff_rn')) * F.col('curr_rn')
343 new_df = new_df.withColumn(value_col, F.when(F.col(value_col).isNull(), lin_interp_func).otherwise(F.col(value_col)))
344
345 if not isinstance(id_cols, list):
346     id_cols = [id_cols]
347
348 keep_cols = id_cols + [order_col, value_col]
349 new_df = new_df.select(keep_cols)
350 return new_df
351
352
353 # forward fill (ffill) interpolation window transformation to impute charge_duration_sec
354 def forward_fill_interpolation(df, id_col, value_col):
355     """
356     Apply forward fill to interpolate missing values in dataframe.
357
358     :param df: Spark DataFrame
359     :param id_cols: string column name to partition by the window function
360     :param value_col: column to be filled
361
362     :returns: spark dataframe updated with interpolated values
363     """
364
365     w1 = Window.partitionBy(id_col).rowsBetween(Window.unboundedPreceding, Window.CurrentRow)
366     w2 = Window.partitionBy(id_col).rowsBetween(Window.CurrentRow, Window.unboundedFollowing)
367
368     new_df = df.withColumn('previous', last(value_col, ignorenulls = True).over(w1))\
369         .withColumn('next', first(value_col, ignorenulls = True).over(w2))\
370         .withColumn('new_score', (coalesce(F.col('previous'), F.col('next')) + coalesce(F.col('next'), F.col('previous')) / 2))\
371         .drop('next', 'previous')
372
373     new_df = new_df.select('battery_id', 'cycle', 'new_score')
374     new_df = new_df.withColumnRenamed('new_score', 'charge_duration_sec')
375
376     return new_df
377
378 # pipeline that we can apply to both the RPT and Regular cycle DataFrames
379 def imputation_pipeline(df, target_feature, ffill_feature, soc_features):
380     """
381     Apply imputation functions (linear, forward fill) to interpolate

```

```

382     missing values in our Regular cycle and RPT cycle DataFrames.
383
384     :param df: input Spark DataFrame
385     :param target_feature: target; linear interpolation
386     :param ffill_feature: feature; forward fill interpolation
387     :param soc_features: list of other features to be linearly interpolated
388
389     :returns: full spark dataframe updated with interpolated values
390     """
391
392     # interpolate
393     cap_pct = fill_linear_interpolation(df = df, id_cols = 'battery_id', order_col = 'cycle', value_col = '
        discharge_capacity_pct')
394
395     # interpolate
396     chg_dur = forward_fill_interpolation(df = df, id_col = 'battery_id', value_col = 'charge_duration_sec')
397
398     # construct output DataFrame
399     new_df = cap_pct.join(chg_dur, on = ['battery_id', 'cycle'], how = 'left')
400
401     # interpolate SOC region features
402     for soc_f in soc_features:
403         intermediate = fill_linear_interpolation(df = df, id_cols = 'battery_id', order_col = 'cycle',
            value_col = soc_f)
404         new_df = new_df.join(intermediate, on = ['battery_id', 'cycle'], how = 'left')
405
406     return new_df
407
408 # features/target to be imputed
409 soc_features = ['SOC1_feature_1', 'SOC1_feature_2', 'SOC1_feature_3',
410                 'SOC2_feature_1', 'SOC2_feature_2', 'SOC2_feature_3',
411                 'SOC3_feature_1', 'SOC3_feature_2', 'SOC3_feature_3',
412                 'SOC4_feature_1', 'SOC4_feature_2', 'SOC4_feature_3']
413 target_feature = 'discharge_capacity_pct'
414 ffill_feature = 'charge_duration_sec'
415
416 # new DataFrame for RPT cycles; all NULL values have been imputed, as detailed above (i => stands for
    interpolated)
417 # then sort by battery_id and cycle
418 i_rpt = imputation_pipeline(f_rpt, target_feature = 'discharge_capacity_pct', ffill_feature = '
        charge_duration_sec', soc_features = soc_features)
419 i_rpt = i_rpt.orderBy(order_cols, ascending = True)
420 i_rpt.cache()
421
422 # new DataFrame for Regular cycles; all NULL values have been imputed, as detailed above (i => stands for
    interpolated)
423 # then sort by battery_id and cycle
424 i_reg = imputation_pipeline(f_reg, target_feature = 'discharge_capacity_pct', ffill_feature = '
        charge_duration_sec', soc_features = soc_features)
425 i_reg = i_reg.orderBy(order_cols, ascending = True)
426 i_reg.cache()
427
428 """
429 Now that we've generated full datasets (containing records for all cycles) for both RPT + Regular cycles
430 (note again we've aligned on this approach with the GBatteries team), we can generate a unified dataset by
431 joining the Regular DataFrame to the RPT dataset (on 'battery_id' and 'cycle').
432
433 Regular cycle features provide valuable information that may improve model performance (source: GBatteries
    team).
434 We'll accomplish this by:
435
436 1. Applying column name labels (prepending) on our features to distinguish between RPT and Regular cycle
    features.
437 2. Performing a join operation on 'battery_id' and 'cycle'.
438
439 Then, we'll move into our modeling phase.
440 """
441

```

```

442 # prepend feature names with 'REG_' and 'RPT_'
443 old_feature_names = ['discharge_capacity_pct', 'charge_duration_sec'] + soc_features
444 rpt_feature_names = ['RPT_' + old_name for old_name in old_feature_names]
445 reg_feature_names = ['REG_' + old_name for old_name in old_feature_names]
446
447 # create column name mapping to apply to i_rpt and i_reg
448 new_rpt_dict = dict(zip(old_feature_names, rpt_feature_names))
449 new_reg_dict = dict(zip(old_feature_names, reg_feature_names))
450
451 def rename_columns(df, columns):
452     if isinstance(columns, dict):
453         for old_name, new_name in columns.items():
454             df = df.withColumnRenamed(old_name, new_name)
455         return df
456     else:
457         raise ValueError("'columns' should be a dict, like {'old_name_1':'new_name_1', 'old_name_2':'new_name_2'}")
458
459 # prepend column names with 'REG_' and 'RPT_', respectively
460 i_reg = rename_columns(i_reg, new_reg_dict)
461 i_rpt = rename_columns(i_rpt, new_rpt_dict)
462
463 # get final resultant DataFrame
464 data = i_rpt.join(i_reg, on = ['battery_id', 'cycle'], how = 'left')
465 data = data.orderBy(order_cols, ascending = True)
466
467 # data.coalesce(1).write.format("com.databricks.spark.csv").option("header", "true").save("dbfs:/FileStore/
    unextrapolated_1219v1.csv")
468
469 data.count()
470
471 # data.coalesce(1).write.format("com.databricks.spark.csv").option("header", "true").save("dbfs:/FileStore/
    unextrapolated_1219v1.csv")
472
473 '''
474 Note here that our linear interpolation function only captures the 'inside' of
475 two bounding values. It does not interpolate in the case that the first or last
476 records for a battery_id are NULL. We need to create a custom function to
477 address this.
478
479 Once this is done, we'll have our final DataFrame 'data' that we can start modeling
480 with. We'll save this as a temp view/table (for the purposes of providing code to
481 the client, as a csv file), which we'll then read as our model input.
482 '''
483
484 # linear interpolation window transformation to impute discharge_capacity_pct
485 def beginning_extrapolation(df, id_cols, order_col, value_col):
486     """
487     Apply linear interpolation to dataframe to fill gaps.
488
489     :param df: Spark DataFrame
490     :param id_cols: string or list of column names to partition by the window function
491     :param order_col: column to use to order by the window function
492     :param value_col: column to be filled
493
494     :returns: spark dataframe updated with interpolated values
495     """
496     # create row number over window and a column with row number only for non missing values
497     lead_window = Window.partitionBy(id_cols).orderBy(order_col)
498     df_2 = df.withColumn("next_value", F.when(F.col(value_col).isNotNull(), F.lead(value_col, 1).over(
        lead_window)))
499
500     w = Window.partitionBy(id_cols).orderBy(order_col)
501     new_df = df_2.withColumn('rn', F.row_number().over(w))
502     new_df = new_df.withColumn('rn_not_null', F.when(F.col(value_col).isNotNull(), F.col('rn')))
503
504     # create relative references to the end values (first 2 value not missing)
505     w_end = Window.partitionBy(id_cols).orderBy(order_col).rowsBetween(0, Window.unboundedFollowing)

```

```

506 new_df = new_df.withColumn('end_val', F.first(value_col, True).over(w_end))
507 new_df = new_df.withColumn('end_rn', F.first('rn_not_null', True).over(w_end))
508 new_df = new_df.withColumn('end_val_2', F.first('next_value', True).over(w_end))
509
510 # create references to gap length and current gap position
511 new_df = new_df.withColumn('diff_rn', F.col('end_rn') - F.col('rn'))
512
513 # calculate linear interpolation value
514 lin_interp_func = (F.col('end_val') - (F.col('end_val_2') - F.col('end_val')) * F.col('diff_rn'))
515 new_df = new_df.withColumn(value_col, F.when(F.col(value_col).isNull(), lin_interp_func).otherwise(F.col(
    value_col)))
516
517 if not isinstance(id_cols, list):
518     id_cols = [id_cols]
519
520 new_df = new_df.drop('next_value', 'rn', 'rn_not_null', 'end_val', 'end_rn', 'end_val_2', 'diff_rn')
521 return new_df
522
523 def ending_extrapolation(df, id_cols, order_col, value_col):
524     """
525     Apply linear interpolation to dataframe to fill gaps.
526
527     :param df: Spark DataFrame
528     :param id_cols: string or list of column names to partition by the window function
529     :param order_col: column to use to order by the window function
530     :param value_col: column to be filled
531
532     :returns: spark dataframe updated with interpolated values
533     """
534     # create row number over window and a column with row number only for non missing values
535     lag_window = Window.partitionBy(id_cols).orderBy(order_col)
536     df_2 = df.withColumn("last_value", F.when(F.col(value_col).isNotNull(), F.lag(value_col, 1).over(
        lag_window)))
537
538     w = Window.partitionBy(id_cols).orderBy(order_col)
539     new_df = df_2.withColumn('rn', F.row_number().over(w))
540     new_df = new_df.withColumn('rn_not_null', F.when(F.col(value_col).isNotNull(), F.col('rn')))
541
542     # create relative references to the start value (last value not missing)
543     w_start = Window.partitionBy(id_cols).orderBy(order_col).rowsBetween(Window.unboundedPreceding, -1)
544     new_df = new_df.withColumn('start_val', F.last(value_col, True).over(w_start))
545     new_df = new_df.withColumn('start_rn', F.last('rn_not_null', True).over(w_start))
546     new_df = new_df.withColumn('start_val_2', F.first('last_value', True).over(w_start))
547
548     # create references to gap length and current gap position
549     new_df = new_df.withColumn('diff_rn', F.col('rn') - F.col('start_rn'))
550
551     # calculate linear interpolation value
552     lin_interp_func = (F.col('start_val') + (F.col('start_val') - F.col('start_val_2')) * F.col('diff_rn'))
553     new_df = new_df.withColumn(value_col, F.when(F.col(value_col).isNull(), lin_interp_func).otherwise(F.col(
        value_col)))
554
555     if not isinstance(id_cols, list):
556         id_cols = [id_cols]
557
558     new_df = new_df.drop('last_value', 'rn', 'rn_not_null', 'start_val', 'start_rn', 'start_val_2', 'diff_rn')
559     return new_df
560
561 # get final df
562 df_extrapolated = data
563
564 # apply changes per the two functions above
565 for col_name in ['RPT_discharge_capacity_pct', 'RPT_SOC1_feature_1', 'RPT_SOC1_feature_2', 'RPT_SOC1_feature_3',
    'RPT_SOC2_feature_1', 'RPT_SOC2_feature_2', 'RPT_SOC2_feature_3', 'RPT_SOC3_feature_1', 'RPT_SOC3_feature_2',
    'RPT_SOC3_feature_3', 'RPT_SOC4_feature_1', 'RPT_SOC4_feature_2', 'RPT_SOC4_feature_3', '
    REG_discharge_capacity_pct', 'REG_SOC1_feature_1', 'REG_SOC1_feature_2', 'REG_SOC1_feature_3', '
    REG_SOC2_feature_1', 'REG_SOC2_feature_2', 'REG_SOC2_feature_3', 'REG_SOC3_feature_1', 'REG_SOC3_feature_2',
    'REG_SOC3_feature_3', 'REG_SOC4_feature_1', 'REG_SOC4_feature_2', 'REG_SOC4_feature_3']:

```



```

566     df_extrapolated = beginning_extrapolation(df = df_extrapolated, id_cols = 'battery_id', order_col = '
        cycle', value_col = col_name)
567
568 # save as csv (alternatively, save as view ... but we're working with Databricks on a client project; we'll
569 # produce the cleaned dataset as a .csv file for easier handoff to the GBatteries team)
570 df_extrapolated = ending_extrapolation(df = df_extrapolated, id_cols = 'battery_id', order_col = 'cycle',
        value_col = col_name)
571
572 df_extrapolated.cache()
573 df_extrapolated.count()
574
575 # save as csv
576 df_extrapolated.coalesce(1).write.format("com.databricks.spark.csv").option("header", "true").save("dbfs:/
        FileStore/extrapolated_output_v7.csv")
577
578 # select only the data of each battery before end of life. Battery A0047B2 has last RPT cycle at 96, too few
        data, and battery A35D29D has only 2 RPT cycles before reaching eol. These 2 batteries are not used for
        modeling.
579 df_eol = df_extrapolated.filter((col("battery_id") != "A8C2C39")|(col("cycle") < 445)).filter((col("
        battery_id") != "A65D895")|(col("cycle") < 368)).filter((col("battery_id") != "A2C0E4A")|(col("cycle") <
        329)).filter((col("battery_id") != "AFFE915")|(col("cycle") < 382)).filter((col("battery_id") != "
        A947A53")|(col("cycle") < 284)).filter((col("battery_id") != "A1D35B0")|(col("cycle") < 280)).filter((
        col("battery_id") != "A74A962")|(col("cycle") < 182)).filter((col("battery_id") != "AFD5A8F")|(col("
        cycle") < 226)).filter((col("battery_id") != "AEBE15E")|(col("cycle") < 138)).filter((col("battery_id")
        != "A716092")|(col("cycle") < 213)).filter((col("battery_id") != "A18242C")|(col("cycle") < 220)).filter
        ((col("battery_id") != "A4E5D13")|(col("cycle") < 240)).filter((col("battery_id") != "AEE8EF6")|(col("
        cycle") < 247)).filter((col("battery_id") != "A231712")|(col("cycle") < 209)).filter((col("battery_id")
        != "A6AD931")|(col("cycle") < 231)).filter((col("battery_id") != "A35D29D").filter((col("battery_id") !=
        "AC3C95D")|(col("cycle") < 119)).filter((col("battery_id") != "A83B987")|(col("cycle") < 211)).filter((
        col("battery_id") != "A73CF24")|(col("cycle") < 209)).filter((col("battery_id") != "A0047B2")|(col("
        cycle") < 127)).filter((col("battery_id") != "AD5B491")|(col("cycle") < 185)).filter((col("battery_id")
        != "A43046B"))
580
581 df_eol.coalesce(1).write.format("com.databricks.spark.csv").option("header", "true").save("dbfs:/FileStore/
        cleaned_data.csv")
582
583 '''
584 Before we move onto the modeling phase. We'll take a look at feature importance. Here we'll
585 also read in the temp view/table/csv that we saved above. To avoid running all the data processing
586 steps again.
587
588 Here, our cleaned dataset contains only two batteries -- we've disregarded the other battery cycles with
589 lower cycle counts (we cannot evaluate sMAPE on higher cycle counts if the raw data does not contain this
590 information). GBatteries has supported this approach (selecting a small subset of batteries to analyze).
        Since
591 are selecting batteries that are representative of the dataset, our methodology can easily be applied to any
592 of the other batteries in the dataset as well. Citation: GBatteries representative (see meeting notes)
593 '''
594
595 # load cleaned data; dbfs:/FileStore/shared_uploads/gbatteries_team8@outlook.com/extrapolated_v6.csv
596 schema_cleaned = StructType([
597     StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
598     StructField("RPT_discharge_capacity_pct", FloatType()), StructField("RPT_charge_duration_sec", FloatType
        ()),
599     StructField("RPT_SOC1_feature_1", FloatType()), StructField("RPT_SOC1_feature_2", FloatType()),
        StructField("RPT_SOC1_feature_3", FloatType()),
600     StructField("RPT_SOC2_feature_1", FloatType()), StructField("RPT_SOC2_feature_2", FloatType()),
        StructField("RPT_SOC2_feature_3", FloatType()),
601     StructField("RPT_SOC3_feature_1", FloatType()), StructField("RPT_SOC3_feature_2", FloatType()),
        StructField("RPT_SOC3_feature_3", FloatType()),
602     StructField("RPT_SOC4_feature_1", FloatType()), StructField("RPT_SOC4_feature_2", FloatType()),
        StructField("RPT_SOC4_feature_3", FloatType()),
603     StructField("REG_discharge_capacity_pct", FloatType()), StructField("REG_charge_duration_sec", FloatType
        ()),
604     StructField("REG_SOC1_feature_1", FloatType()), StructField("REG_SOC1_feature_2", FloatType()),
        StructField("REG_SOC1_feature_3", FloatType()),
605     StructField("REG_SOC2_feature_1", FloatType()), StructField("REG_SOC2_feature_2", FloatType()),
        StructField("REG_SOC2_feature_3", FloatType()),

```

```

606     StructField("REG_SOC3_feature_1", FloatType()), StructField("REG_SOC3_feature_2", FloatType()),
        StructField("REG_SOC3_feature_3", FloatType()),
607     StructField("REG_SOC4_feature_1", FloatType()), StructField("REG_SOC4_feature_2", FloatType()),
        StructField("REG_SOC4_feature_3", FloatType())
608 ])
609
610 # cleaned dataset for 2 batteries that we will model ...
611 dff = spark.read.format("csv").load("dbfs:/FileStore/shared_uploads/gbatteries_team8@outlook.com/
        cleaned_data.csv", header = True, schema = schema_cleaned)
612
613 #####
614 # 2.1 SEASONALITY AND STATIONARY
615 #####
616
617 '''
618 Now we'll move onto the modeling phase. In this section, we'll test several (>=5) time series models
619 on our processed data set and compare results. Specifically, we'll touch upon the following for EACH
620 model...
621
622 * Notes on what affects performance and sMAPE for each model
623 * Basic model information and how the model was built
624 * Model construction and optimization process
625 * Challenges and limitations
626
627 We'll cover the following 5 models in this notebook...
628 * Linear Regression
629 * Simple Moving Average
630 * ARIMA
631 * Random Forest
632 * XGBoost
633
634 '''
635
636 #####
637 # 3. MODELING, MODEL EVALUATION (5 MODELS)
638 #####
639
640 '''
641 The RPT discharge capacity has a linear-like relationship with cycle number and linearity is an
642 important concept throughout the project (e.g. linear interpolation, linear extrapolation, and
643 first order differencing). Linear regression is definitely a model needed to be attempted.
644 While performance of regular linear regression suffers from outliers, enhanced linear regression has
645 outlier rejection algorithm involved.
646 '''
647
648 #####
649 # 3.a) LINEAR REGRESSION
650 #####
651 schema_cleaned = StructType([
652     StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
653     StructField("discharge_capacity_pct", DoubleType()), StructField("consecutive_rpt", IntegerType())
654 ])
655
656 data = spark.read.format("csv").load("dbfs:/FileStore/tables/rpt_actual_with_consec.csv", header = True,
        schema = schema_cleaned)
657
658 # get battery id list
659 unique_id = data.select("battery_id").distinct()
660 unique_id_list = [str(i) for i in np.array(unique_id.collect())[:,0]]
661
662 # create empty result dataframe
663 schema = StructType([
664     StructField('battery_id', StringType(), True),
665     StructField('cycle', IntegerType(), True),
666     StructField('discharge_capacity_pct', DoubleType(), True),
667     StructField('regular_prediction', IntegerType(), True),
668     StructField('enhanced_prediction', IntegerType(), True)
669 ])

```

```

670 df_result = spark.createDataFrame([], schema)
671
672 # data dataframe setup for applying ml
673 assembler = VectorAssembler(inputCols=["cycle"],outputCol="features")
674 data_ml = assembler.transform(data)
675 data_ml = data_ml.withColumnRenamed("discharge_capacity_pct",'label')
676
677 # set up window
678 max_window = Window.partitionBy(col("battery_id")).orderBy("battery_id")
679 max_diff = max(col("difference")).over(max_window)
680
681 for battery_id in unique_id_list:
682     # regular linear regression
683     lr = LinearRegression()
684     train = data_ml.filter((data_ml["battery_id"]==battery_id) & (data_ml["cycle"]<=110))
685     lr_model = lr.fit(train)
686     test = lr_model.transform(data_ml.filter((data_ml["battery_id"]==battery_id) & (data_ml["cycle"]>110))).
        withColumnRenamed("prediction","regular_prediction")
687
688     # regression with weighted data points and outlier removing
689     train = train.filter((col("label") < lit(0.95)) & (col("label") > lit(0.80)))
690     train_consec = train.filter(col("consecutive_rpt") == 1)
691     train = train.union(train_consec).union(train_consec)
692     for i in range(3):
693         lr = LinearRegression()
694         lr_model_1 = lr.fit(train)
695         train = lr_model_1.transform(train).withColumn("difference",F.abs(col("label") - col("prediction")))
696         count_1 = train.count()
697         train = train.select(col("*"), max_diff.alias("max_difference")).filter((col("difference") != col("
        max_difference")) | (col("max_difference") < lit(0.045))).drop("difference","max_difference","prediction
        ")
698         count_2 = train.count()
699         # no more outlier
700         if (count_1 == count_2):
701             break
702
703     lr = LinearRegression()
704     lr_enhanced_model = lr.fit(train)
705     test = lr_enhanced_model.transform(test).withColumnRenamed("prediction","enhanced_prediction").drop("
        consecutive_rpt","features").withColumnRenamed("label","discharge_capacity_pct")
706     df_result = df_result.union(test)
707
708     battery_id = unique_id_list[0]
709
710 # regular linear regression
711 lr = LinearRegression()
712 train = data_ml.filter((data_ml["battery_id"]==battery_id) & (data_ml["cycle"]<=110))
713 lr_model = lr.fit(train)
714 test = lr_model.transform(data_ml.filter((data_ml["battery_id"]==battery_id) & (data_ml["cycle"]>110))).
        withColumnRenamed("prediction","regular_prediction")
715
716 # regression with weighted data points and outlier removing
717 train = train.filter((col("label") < lit(0.95)) & (col("label") > lit(0.80)))
718 train_consec = train.filter(col("consecutive_rpt") == 1)
719 train = train.union(train_consec).union(train_consec)
720 for i in range(3):
721     lr = LinearRegression()
722     lr_model_1 = lr.fit(train)
723     train = lr_model_1.transform(train).withColumn("difference",F.abs(col("label") - col("prediction")))
724     count_1 = train.count()
725     train = train.select(col("*"), max_diff.alias("max_difference")).filter((col("difference") != col("
        max_difference")) | (col("max_difference") < lit(0.045))).drop("difference","max_difference","prediction
        ")
726     count_2 = train.count()
727     # no more outlier
728     if (count_1 == count_2):
729         break
730

```

```

731 lr = LinearRegression()
732 lr_enhanced_model = lr.fit(train)
733 test = lr_enhanced_model.transform(test).withColumnRenamed("prediction", "enhanced_prediction").drop("
    consecutive_rpt", "features").withColumnRenamed("label", "discharge_capacity_pct")
734
735 # export result table as .csv
736 df_result.coalesce(1).write.format("com.databricks.spark.csv").option("header", "true").save("dbfs:/
    FileStore/lin_reg_result_table.csv")
737
738 # read in resultant csv just exported
739 res_schema = StructType([
740     StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
741     StructField("discharge_capacity_pct", DoubleType()), StructField("regular_prediction",
    DoubleType()),
742     StructField("enhanced_prediction", DoubleType())
743 ])
744
745 lin_reg_res = spark.read.format("csv").load("dbfs:/FileStore/lin_reg_result_table.csv", header = True, schema
    = res_schema)
746
747 # compute smape
748 lin_reg_res = lin_reg_res.withColumn("time_step", F.when(col("cycle") < 161, 50).when(col("cycle") < 211,
    100).when(col("cycle") < 211, 100).when(col("cycle") < 261, 150).when(col("cycle") < 311, 200).when(col(
    "cycle") < 361, 250).when(col("cycle") < 411, 300).when(col("cycle") < 461, 350).otherwise(0))
749 lin_reg_res = lin_reg_res.withColumn("regular_expre", F.abs(col("discharge_capacity_pct") - col("
    regular_prediction")) * 200 / (col("discharge_capacity_pct") + col("regular_prediction")))
750 lin_reg_res = lin_reg_res.withColumn("enhanced_expre", F.abs(col("discharge_capacity_pct") - col("
    enhanced_prediction")) * 200 / (col("discharge_capacity_pct") + col("enhanced_prediction")))
751 enhanced_smape = lin_reg_res.groupBy("time_step").agg(F.mean("enhanced_expre").alias("smape"))
752 regular_smape = lin_reg_res.groupBy("time_step").agg(F.mean("regular_expre").alias("smape"))
753
754 '''
755 Simple Moving Average is a method of time series smoothing and is a very basic forecasting technique.
756 It does not need estimation of parameters, but rather is based on order selection. This will serve as
757 a good baseline model and highlight certain areas to address with improved, more complex models.
758
759 Specifically in this problem domain, SMA models give a lot of weight to old data -- this is something
760 that exponential moving averages (EMA) can address. SMAs and EMAs are used in similar ways: to identify
761 trends and find potential areas of 'support' or 'resistance'.
762 '''
763
764 #####
765 # 3.b) Simple Moving Average
766 #####
767
768 schema_cleaned = StructType([
769     StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
770     StructField("RPT_discharge_capacity_pct", FloatType()), StructField("RPT_charge_duration_sec", FloatType
    ()),
771     StructField("RPT_SOC1_feature_1", FloatType()), StructField("RPT_SOC1_feature_2", FloatType()),
    StructField("RPT_SOC1_feature_3", FloatType()),
772     StructField("RPT_SOC2_feature_1", FloatType()), StructField("RPT_SOC2_feature_2", FloatType()),
    StructField("RPT_SOC2_feature_3", FloatType()),
773     StructField("RPT_SOC3_feature_1", FloatType()), StructField("RPT_SOC3_feature_2", FloatType()),
    StructField("RPT_SOC3_feature_3", FloatType()),
774     StructField("RPT_SOC4_feature_1", FloatType()), StructField("RPT_SOC4_feature_2", FloatType()),
    StructField("RPT_SOC4_feature_3", FloatType()),
775     StructField("REG_discharge_capacity_pct", FloatType()), StructField("REG_charge_duration_sec", FloatType
    ()),
776     StructField("REG_SOC1_feature_1", FloatType()), StructField("REG_SOC1_feature_2", FloatType()),
    StructField("REG_SOC1_feature_3", FloatType()),
777     StructField("REG_SOC2_feature_1", FloatType()), StructField("REG_SOC2_feature_2", FloatType()),
    StructField("REG_SOC2_feature_3", FloatType()),
778     StructField("REG_SOC3_feature_1", FloatType()), StructField("REG_SOC3_feature_2", FloatType()),
    StructField("REG_SOC3_feature_3", FloatType()),
779     StructField("REG_SOC4_feature_1", FloatType()), StructField("REG_SOC4_feature_2", FloatType()),
    StructField("REG_SOC4_feature_3", FloatType())
780 ])

```

```

781
782 data = spark.read.format("csv").load("dbfs:/FileStore/tables/final_data.csv",header = True,schema =
    schema_cleaned)
783
784 sma_data = data.select([col("battery_id"),col("cycle"),col("RPT_discharge_capacity_pct").alias("label")])
785 # get battery id list
786 unique_id = sma_data.select("battery_id").distinct()
787 unique_id_list = [str(i) for i in np.array(unique_id.collect())[:,0]]
788
789 # create empty result dataframe
790 schema = StructType([
791     StructField('battery_id', StringType(), True),
792     StructField('cycle', IntegerType(), True),
793     StructField('sma_prediction',IntegerType(),True)
794 ])
795 df_result = spark.createDataFrame([], schema)
796
797 # this section of code computes predictions forward with window size of 5, but takes much less time (run
    time of several seconds)
798 for id in unique_id_list:
799     sma_data_id = sma_data.filter(col("battery_id") == id).drop("battery_id")
800     sma_data_id_count = sma_data_id.count()
801     sma_data_id_fill0 = sma_data_id.withColumn("label",F.when(F.col("cycle") < 101,col("label")).otherwise(lit
        (0)))
802     sma_predicted_id = sma_data_id_fill0.collect()
803     for i in range(101,sma_data_id_count+1):
804         new_row = Row(cycle = i, label = (sma_predicted_id[i-2][1] + sma_predicted_id[i-3][1] + sma_predicted_id
            [i-4][1] + sma_predicted_id[i-5][1] + sma_predicted_id[i-6][1])/5)
805         sma_predicted_id[i-1] = new_row
806     sma_data_id_forecast = spark.createDataFrame(sma_predicted_id).withColumn("battery_id",lit(id)).select("
        battery_id","cycle","label").withColumnRenamed("label","sma_prediction")
807     df_result = df_result.union(sma_data_id_forecast)
808
809     schema = StructType([
810         StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
811         StructField("discharge_capacity_pct", DoubleType()), StructField("consecutive_rpt", IntegerType())
812     ])
813
814 # this section of code also computes predictions forward, takes long time to run (not for running purpose,
    or variable name needs to be adjusted)
815 # sma_data_1_forecast = sma_data_1.filter(col("cycle") < 111)
816 # sma_data_2_forecast = sma_data_2.filter(col("cycle") < 79)
817 # create a one-row dataframe with cycle and predicted label then union to the actual data
818 # for i in range(111,sma_data_1_count+1):
819 #     sma_data_1_append = sma_data_1_forecast.filter(col("cycle") > (i-1)-5).drop("cycle").agg(avg(col("label"))
        .alias("label")).withColumn("cycle",lit(i)).select("cycle","label")
820 #     sma_data_1_forecast = sma_data_1_forecast.union(sma_data_1_append).orderBy("cycle")
821
822 # for i in range(79,sma_data_2_count+1):
823 #     sma_data_2_append = sma_data_2_forecast.filter(col("cycle") > (i-1)-5).drop("cycle").agg(avg(col("label"))
        .alias("label")).withColumn("cycle",lit(i)).select("cycle","label")
824 #     sma_data_2_forecast = sma_data_2_forecast.union(sma_data_2_append).orderBy("cycle")
825
826
827 actual_data = spark.read.format("csv").load("dbfs:/FileStore/tables/rpt_actual_with_consec.csv",header =
    True,schema = schema).filter(col("cycle") > 100)
828
829 df_result = actual_data.join(df_result,["battery_id","cycle"],"left")
830
831 df_result.coalesce(1).write.format("com.databricks.spark.csv").option("header", "true").save("dbfs:/
    FileStore/sma_result_table.csv")
832
833 # read in stored resultant data from sma
834 schema = StructType([
835     StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
836     StructField("discharge_capacity_pct", DoubleType()), StructField("consecutive_rpt", IntegerType()),
837     StructField("sma_prediction", DoubleType())
838 ])

```

```

839
840 sma_res = spark.read.format("csv").load("dbfs:/FileStore/sma_result_table.csv",header = True,schema = schema
      )
841
842 # compute smape
843 sma_res = sma_res.select("battery_id","cycle","discharge_capacity_pct","sma_prediction")
844 sma_res = sma_res.withColumn("time_step",F.when(col("cycle") < 151, 50).when(col("cycle") < 201, 100).when(
      col("cycle") < 211, 100).when(col("cycle") < 251, 150).when(col("cycle") < 301, 200).when(col("cycle") <
      351, 250).when(col("cycle") < 401, 300).when(col("cycle") < 451, 350).otherwise(0))
845 sma_res = sma_res.withColumn("sma_expre",F.abs(col("discharge_capacity_pct")-col("sma_prediction")) * 200 /
      (col("discharge_capacity_pct")+col("sma_prediction")))
846 sma_smape = sma_res.groupBy("time_step").agg(F.mean("sma_expre").alias("smape"))
847
848 import matplotlib.pyplot as plt
849
850 fig, ax = plt.subplots(figsize=(20,5))
851 ax.plot(actl, 'black', label='test data actuals')
852 ax.plot(pred, 'r', label='test data forecast')
853 #ax.plot(ts_log_diff, 'blue', label='training data')
854 legend = ax.legend(loc='upper left')
855 legend.get_frame().set_facecolor('w')
856 display(fig.figure)
857
858 schema_gt_ARIMA = StructType([StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
      StructField("discharge_capacity_pct", FloatType())])
859 gt = spark.read.format("csv").load("dbfs:/FileStore/tables/rpt_actual_v8.csv",header = True,schema =
      schema_gt_ARIMA)
860
861 def SMAPE_ARIMA (pred, battery_id, num):
862     ground_truth_df = gt.filter(col("battery_id")==battery_id)
863     ground_truth = np.array([ground_truth_df.select("discharge_capacity_pct").collect()[i][
      "discharge_capacity_pct"] for i in range(ground_truth_df.filter(col("cycle")<=100).count(),
      ground_truth_df.filter(col("cycle")<=(100+num)).count())])
864     predictions = []
865     cycles = np.array([ground_truth_df.select("cycle").collect()[i]["cycle"] for i in range(ground_truth_df.
      filter(col("cycle")<=100).count(), ground_truth_df.filter(col("cycle")<=(100+num)).count())])
866
867     for i in cycles:
868         index = int(i)
869         predictions.append(pred[index-1])
870     predictions = np.array(predictions)
871
872     return 100*np.sum(np.abs(ground_truth-predictions)/((np.abs(ground_truth)+np.abs(predictions))/2))/len(
      ground_truth)
873
874 #####
875 # 3.c) ARIMA
876 #####
877
878 !pip install pmdarima
879 from pmdarima import auto_arima
880 from statsmodels.tsa.arima_model import ARIMA
881 import matplotlib.pyplot as plt
882 import math
883
884 def ARIMA_model(battery_id):
885
886     # Filter out one battery from the cleaned_data.csv
887     train = dff.filter(col("battery_id") == battery_id).select('battery_id','cycle','
      RPT_discharge_capacity_pct')
888     train_pandas = train.toPandas()
889
890     # Create test array and sort in descending order, also define a sub-array to include only the first 50
      datapoints
891     train_array = np.array(train_pandas.RPT_discharge_capacity_pct)
892     train_array100 = train_array[:100]
893
894     # Use auto_arima to tune the parameters p,d,q for the ARIMA model

```

```

895 stepwise_fit = auto_arima(train_array100,
896                           start_p = 1,
897                           start_q = 1,
898                           max_p = 3,
899                           max_q = 3,
900                           m = 1,
901                           seasonal = None,
902                           d = 1,
903                           trace = False,
904                           random = True,
905                           error_action = 'ignore',
906                           suppress_warnings = True,
907                           stepwise = True)
908
909 # Use the sub-array to fit the model, the best parameters are set according to the result from auto_ARIMA
910 model = ARIMA(train_array100, order=stepwise_fit.order)
911 results_ARIMA = model.fit()
912
913 # Apply the training set to forecast the rest of the test set
914 pred = results_ARIMA.predict(1, len(train_array)-1, typ='levels')
915
916 # Plot the actual test result and the prediction result
917 fig, ax = plt.subplots(figsize=(8,6))
918 ax.set(title='Prediction for battery'+str(battery_id), xlabel='Cycles', ylabel='RPT discharge capacity')
919 ax.plot(train_array, 'black', label='test data actuals')
920 ax.plot(pred, 'r', label='test data forecast')
921 legend = ax.legend(loc='upper right')
922 legend.get_frame().set_facecolor('w')
923 display(fig.figure)
924
925 # SMAPE calculations
926 cycles_list = [50,100,150,200,250,300,350]
927 smape_val = []
928 for cycles in cycles_list:
929     smape_val.append(sMAPE_ARIMA(pred, battery_id, cycles))
930 smape_val = [float(val) for val in smape_val]
931 sqlContext.createDataFrame(zip(cycles_list, smape_val), schema=['no_cycles(prediction)', 'smape']).show()
932 smape_val = np.array(smape_val)
933 return smape_val
934
935 # 20 batteries
936 battery_id = np.array(['A0047B2', 'A18242C', 'A1D35B0', 'A231712', 'A2C0E4A', 'A4E5D13', 'A65D895', 'A6AD931',
937                        'A716092', 'A73CF24', 'A74A962', 'A83B987',
938                        'A8C2C39', 'A947A53', 'AC3C95D', 'AD5B491', 'AEBE15E', 'AEE8EF6', 'AFD5A8F', 'AFFE915'])
939
940 # For each battery, calculate its SMAPE
941 smape_sum_ARIMA = np.zeros(7)
942 for i in battery_id:
943     a = ARIMA_model(i)
944     if not (np.isnan(a[3])) :
945         smape_sum_ARIMA = smape_sum_ARIMA + a
946
947 # Obtain the average SMAPE for all 20 batteries
948 smape_avg_ARIMA = smape_sum_ARIMA/len(battery_id)
949 print(smape_avg_ARIMA)
950
951 '''
952 Random Forest is a popular and effective ensemble machine learning algorithm, widely used for
953 classification and regression predictive modeling problems. It can also be used for time series
954 forecasting, provided that our data is transformed using a sliding-window representation
955 (source: https://machinelearningmastery.com/random-forest-for-time-series-forecasting/).
956
957 Multivariate Time Series Forecasting Using Random Forest
958
959 '''
960 #####
961 # 3.d) RANDOM FOREST

```



```

962 #####
963
964 from pyspark.mllib.tree import RandomForest, RandomForestModel
965 from pyspark.mllib.util import MLUtils
966 from pyspark.sql import Row
967 from pyspark.ml.linalg import Vectors
968 from pyspark.ml import Pipeline
969 from pyspark.ml.regression import LinearRegression
970 from pyspark.ml.feature import VectorIndexer
971 from pyspark.ml.evaluation import RegressionEvaluator
972 from pyspark.ml.regression import RandomForestRegressor
973 import matplotlib.pyplot as plt
974 from pyspark.sql.window import Window
975 import pyspark.sql.functions as func
976 from pyspark.sql.functions import lit
977 from pyspark.sql.functions import expr
978
979
980
981 # Load the data file
982 dff_RF = spark.read.format("csv").load("dbfs:/FileStore/tables/extrapolated_v8.csv", header = True, schema =
    schema_cleaned)
983
984 # a function used to convert features to dense vectors
985 def transData(data):
986     return data.rdd.map(lambda r: [Vectors.dense(r[1]),r[-1]]).toDF(['features','label'])
987
988 def RF (batt_id):
989     data1 = dff_RF.filter(col("battery_id") == batt_id ).withColumn('user',lit('group8'))
990     #data2 is for generating test data later on
991     data2 = dff_RF.filter(col("battery_id") == batt_id ).select(['battery_id','cycle','RPT_SOC1_feature_1','
        RPT_SOC1_feature_2','RPT_SOC1_feature_3','RPT_SOC2_feature_1','RPT_SOC2_feature_2','RPT_SOC2_feature_3',
        'RPT_SOC3_feature_1','RPT_SOC3_feature_2','RPT_SOC3_feature_3','RPT_SOC4_feature_1','RPT_SOC4_feature_2',
        'RPT_SOC4_feature_3','REG_discharge_capacity_pct','REG_SOC1_feature_1','REG_SOC1_feature_2','
        REG_SOC1_feature_3','REG_SOC2_feature_1','REG_SOC2_feature_2','REG_SOC2_feature_3','REG_SOC3_feature_1',
        'REG_SOC3_feature_2','REG_SOC3_feature_3','REG_SOC4_feature_1','REG_SOC4_feature_2','REG_SOC4_feature_3',
        'RPT_discharge_capacity_pct'])
992
993     # generate the new label
994     data1 = data1.withColumn('RPT_discharge_capacity_pct_lag', func.lag(data1['RPT_discharge_capacity_pct']).
        over(Window.orderBy("user"))))
995     data1 = data1.withColumn('Difference', data1['RPT_discharge_capacity_pct'] - data1['
        RPT_discharge_capacity_pct_lag'] )
996
997     # generate the new set of features
998     for i in ('RPT_SOC1_feature_1','RPT_SOC1_feature_2','RPT_SOC1_feature_3','RPT_SOC2_feature_1',
        RPT_SOC2_feature_2','RPT_SOC2_feature_3','RPT_SOC3_feature_1','RPT_SOC3_feature_2','RPT_SOC3_feature_3',
        'RPT_SOC4_feature_1','RPT_SOC4_feature_2','RPT_SOC4_feature_3','REG_discharge_capacity_pct',
        REG_SOC1_feature_1','REG_SOC1_feature_2','REG_SOC1_feature_3','REG_SOC2_feature_1','REG_SOC2_feature_2',
        'REG_SOC2_feature_3','REG_SOC3_feature_1','REG_SOC3_feature_2','REG_SOC3_feature_3','REG_SOC4_feature_1',
        'REG_SOC4_feature_2','REG_SOC4_feature_3'):
999         data1 = data1.withColumn(f"lag_{i}", func.lag(data1[i]).over(Window.orderBy("user")))
1000         data1 = data1.withColumn(f"difference_{i}", data1[i] - data1[f"lag_{i}"])
1001
1002
1003     # create a new df consisting of the new features and labels
1004     data = data1.select(['battery_id','cycle','difference_RPT_SOC1_feature_1','difference_RPT_SOC1_feature_2',
        'difference_RPT_SOC1_feature_3','difference_RPT_SOC2_feature_1','difference_RPT_SOC2_feature_2',
        'difference_RPT_SOC2_feature_3','difference_RPT_SOC3_feature_1','difference_RPT_SOC3_feature_2',
        'difference_RPT_SOC3_feature_3','difference_RPT_SOC4_feature_1','difference_RPT_SOC4_feature_2',
        'difference_RPT_SOC4_feature_3','difference_REG_discharge_capacity_pct','difference_REG_SOC1_feature_1',
        'difference_REG_SOC1_feature_2','difference_REG_SOC1_feature_3','difference_REG_SOC2_feature_1',
        'difference_REG_SOC2_feature_2','difference_REG_SOC2_feature_3','difference_REG_SOC3_feature_1',
        'difference_REG_SOC3_feature_2','difference_REG_SOC3_feature_3','difference_REG_SOC4_feature_1',
        'difference_REG_SOC4_feature_2','difference_REG_SOC4_feature_3','difference'])
1005
1006     # Convert the data to dense vector
1007     transformed = transData(data)

```



```

1008 transformed2 = transData(data2)
1009
1010 # Deal with categorical variables
1011 featureIndexer = VectorIndexer(inputCol="features", \
1012                                outputCol="indexedFeatures", \
1013                                maxCategories=4).fit(transformed)
1014
1015 data_transformed = featureIndexer.transform(transformed)
1016 data_transformed2 = featureIndexer.transform(transformed2)
1017 w = Window().partitionBy(lit('a')).orderBy(lit('a'))
1018
1019 data_transformed = data_transformed.withColumn("row_num", row_number().over(w))
1020 data_transformed2 = data_transformed2.withColumn("row_num", row_number().over(w))
1021
1022 # train-test split, the first 100 cycles are used for training.
1023 # The average feature decreasing rates are called testData, but it is used for prediction
1024 trainingData = data_transformed.filter(col("row_num").between(2,100))
1025 testData = (data_transformed2.filter(col('row_num') == 100).subtract(data_transformed2.filter(col('row_num'
1026                                     ') == 1))).toPandas()
1027 testData = testData.div(100)
1028 testData = sqlContext.createDataFrame(testData)
1029
1030 # Define RandomForest algorithm
1031 rf = RandomForestRegressor()
1032
1033 # Pipeline Architecture
1034 pipeline = Pipeline(stages=[featureIndexer, rf])
1035
1036 model= pipeline.fit(trainingData)
1037 predictions = model.transform(testData)
1038
1039 #set the decrease rate for the label after the first 100 training cycles
1040 num_rows_needed = data1.count()-100
1041 result_rows = predictions.select('prediction').withColumn('n',lit(num_rows_needed))
1042 result_rows = result_rows.withColumn('n', expr('explode(array_repeat(n,int(n))))')
1043 result_rows = result_rows.drop('n')
1044 first_row = data1.selectExpr("RPT_discharge_capacity_pct as prediction")
1045
1046 #use cumSum to calculate the future capacities
1047 w = Window().partitionBy(lit('a')).orderBy(lit('a'))
1048 first_row = first_row.withColumn("row_num", row_number().over(w))
1049 first_row = first_row.filter(col("row_num").between(100,100))
1050 first_row = first_row.drop('row_num')
1051 result_rows = first_row.union(result_rows)
1052 result_rows = result_rows.withColumn("row_num",row_number().over(w))
1053 result_rows = result_rows.withColumn('RPT_prediction', sum('prediction').over(Window.orderBy('row_num'))))
1054 result_rows = result_rows.select('RPT_prediction')
1055 training_rows = data1.selectExpr("RPT_discharge_capacity_pct as RPT_prediction").limit(99)
1056 result = training_rows.union(result_rows)
1057
1058
1059
1060
1061 # Plot
1062 pred = np.array([val.RPT_prediction for val in result.select('RPT_prediction').collect()])
1063 actual = np.array([val.RPT_discharge_capacity_pct for val in data1.select('RPT_discharge_capacity_pct').
1064                                     collect()])
1065
1066 fig, ax = plt.subplots(figsize=(8,6))
1067 ax.set(title='Prediction for battery ' + str(batt_id), xlabel='Cycles', ylabel='RPT discharge capacity')
1068 ax.plot(actual, label='Actual')
1069 ax.plot(pred, label='Prediction')
1070 legend = ax.legend(loc='upper right')
1071 display(fig.figure)
1072
1073 # compute sMAPE
1074 num_forecast = [50,100,150,200,250,300,350]

```

```

1074 smape_val = []
1075 for num in num_forecast:
1076     smape_val.append(sMAPE_RF(pred,batt_id,num))
1077 smape_val = [float(val) for val in smape_val]
1078 sqlContext.createDataFrame(zip(num_forecast, smape_val), schema=['no_cycles(prediction)', 'smape']).show()
1079 smape_val = np.array(smape_val)
1080 return smape_val
1081
1082 # import ground truth
1083 schema_gt_rf = StructType([StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
1084                               StructField("discharge_capacity_pct", FloatType())])
1084 gt = spark.read.format("csv").load("dbfs:/FileStore/tables/rpt_actual_v8.csv",header = True,schema =
1085                               schema_gt_rf)
1086
1087 #calculate sMAPE
1088 def sMAPE_RF (pred, batt_id, num):
1089     ground_truth_df = gt.filter(col("battery_id")==batt_id)
1090     ground_truth = np.array([ground_truth_df.select("discharge_capacity_pct").collect()[i]["
1091                               discharge_capacity_pct"] for i in range(ground_truth_df.filter(col("cycle")<=100).count(),
1092                               ground_truth_df.filter(col("cycle")<=(100+num)).count())])
1093     predictions = []
1094     cycles = np.array([ground_truth_df.select("cycle").collect()[i]["cycle"] for i in range(ground_truth_df.
1095                               filter(col("cycle")<=100).count(), ground_truth_df.filter(col("cycle")<=(100+num)).count())])
1096
1097     for i in cycles:
1098         index = int(i)
1099         predictions.append(pred[index-1])
1100     predictions = np.array(predictions)
1101
1102     return 100*np.sum(np.abs(ground_truth-predictions)/((np.abs(ground_truth)+np.abs(predictions))/2))/len(
1103           ground_truth)
1104
1105 #RF calculate average for all batteries
1106 batt_id = np.array(dff_RF.select("battery_id").distinct().collect())
1107 smape_sum = np.zeros(7)
1108 b = 0
1109 for i in batt_id:
1110     #two of the abandoned batteries are still in the df, removed
1111     if i[0] != 'A43046B' and i[0] != 'A35D29D':
1112         a = RF(i[0])
1113         if not (np.isnan(a[3])) :
1114             smape_sum = smape_sum + a
1115             b = b+1
1116 smape_avg = smape_sum/b
1117
1118 batt_id = np.array(['AEBE15E', 'A716092', 'A4E5D13', 'A83B987', 'AEE8EF6', 'A8C2C39', 'AC3C95D', 'A6AD931',
1119                    'A18242C', 'A1D35B0', 'AFD5A8F', 'A947A53', 'A74A962'])
1120 smape_sum = np.zeros(7)
1121 for i in batt_id:
1122     a = RF(i)
1123     if not (np.isnan(a[3])) :
1124         smape_sum = smape_sum + a
1125 smape_avg = smape_sum/len(batt_id)
1126
1127 # get the avg smape
1128 print(smape_avg)
1129
1130 '''
1131 GBT
1132 '''
1133
1134 #####
1135 # 3.e) GBT
1136 #####
1137 schema_cleaned = StructType([
1138     StructField("battery_id", StringType()), StructField("cycle", IntegerType()),
1139     StructField("RPT_discharge_capacity_pct", FloatType()), StructField("RPT_charge_duration_sec", FloatType
1140     ()),

```

```

1134     StructField("RPT_SOC1_feature_1", FloatType()), StructField("RPT_SOC1_feature_2", FloatType()),
1135     StructField("RPT_SOC1_feature_3", FloatType()), StructField("RPT_SOC2_feature_1", FloatType()), StructField("RPT_SOC2_feature_2", FloatType()),
1136     StructField("RPT_SOC2_feature_3", FloatType()), StructField("RPT_SOC3_feature_1", FloatType()), StructField("RPT_SOC3_feature_2", FloatType()),
1137     StructField("RPT_SOC3_feature_3", FloatType()), StructField("RPT_SOC4_feature_1", FloatType()), StructField("RPT_SOC4_feature_2", FloatType()),
1138     StructField("RPT_SOC4_feature_3", FloatType()), StructField("REG_discharge_capacity_pct", FloatType()), StructField("REG_charge_duration_sec", FloatType
1139     ()), StructField("REG_SOC1_feature_1", FloatType()), StructField("REG_SOC1_feature_2", FloatType()),
1140     StructField("REG_SOC1_feature_3", FloatType()), StructField("REG_SOC2_feature_1", FloatType()), StructField("REG_SOC2_feature_2", FloatType()),
1141     StructField("REG_SOC2_feature_3", FloatType()), StructField("REG_SOC3_feature_1", FloatType()), StructField("REG_SOC3_feature_2", FloatType()),
1142     StructField("REG_SOC3_feature_3", FloatType()), StructField("REG_SOC4_feature_1", FloatType()), StructField("REG_SOC4_feature_2", FloatType()),
1143     StructField("REG_SOC4_feature_3", FloatType())
1144 ]])
1145 schema_gt = StructType([StructField("battery_model", StringType()), StructField("battery_id", StringType()),
1146     StructField("cycle", IntegerType()), StructField("RPT_discharge_capacity_pct", FloatType())])
1147 df = spark.read.format("csv").load("dbfs:/FileStore/tables/extrapolated_v6-1.csv", header = True, schema =
1148     schema_cleaned)
1149 gt = spark.read.format("csv").load("dbfs:/FileStore/tables/cleaned_data_RPT-2.csv", header = True, schema =
1150     schema_gt).filter((col("battery_id")=="A8C2C39") | (col("battery_id")=="AFFE915")).drop("battery_model")
1151 # list of battery and their eol
1152 bad_batt = ["AEBE15E", "AEE8EF6", "A231712", "A35D29D", "AC3C95D", "AD5B491", "A43046B"]
1153 bad_eol = [254, 249, 238, 179, 191, 196, 221]
1154 battery = ["A716092", "A4E5D13", "A83B987", "A73CF24", "A6AD931", "A18242C", "A1D35B0", "AFD5A8F", "A947A53",
1155     "A74A962", "A2C0E4A", "A0047B2", "A65D895", "A8C2C39", "AFFE915"]
1156 eol = [212, 239, 210, 208, 230, 219, 298, 225, 375, 181, 328, 146, 403, 444, 381]
1157 all_battery = battery + bad_batt
1158 all_eol = eol + bad_eol
1159 # STEP 1 Linear Regression of both Regular and RPT Batteries
1160 def LR(df, battery_list, eol_list):
1161     assembler = VectorAssembler(
1162         inputCols=["cycle"],
1163         outputCol="features")
1164     lr = LinearRegression()
1165     coef_rpt = []
1166     inter_rpt = []
1167     coef_reg = []
1168     inter_reg = []
1169     for i in range(len(battery_list)):
1170         train = df.filter((df["battery_id"]==battery_list[i]) & (df["cycle"]<=eol_list[i])).withColumnRenamed("
1171             RPT_discharge_capacity_pct", "label").select("cycle", "label")
1172         a = assembler.transform(train)
1173         lrModel = lr.fit(a)
1174         coef_rpt.append(lrModel.coefficients.values[0].item())
1175         inter_rpt.append(lrModel.intercept)
1176     for i in range(len(battery_list)):
1177         train = df.filter((df["battery_id"]==battery_list[i]) & (df["cycle"]<=eol_list[i])).withColumnRenamed("
1178             REG_discharge_capacity_pct", "label").select("cycle", "label")
1179         a = assembler.transform(train)
1180         lrModel = lr.fit(a)
1181         coef_reg.append(lrModel.coefficients.values[0].item())

```

```

1187     inter_reg.append(lrModel.intercept)
1188
1189     return coef_rpt, inter_rpt, coef_reg, inter_reg
1190
1191
1192 # Step 2 Prepare data for GBT
1193 def dataPrepare_gbt(df, battery_list, coef_rpt, inter_rpt, coef_reg, inter_reg, threshold=50):
1194     w = Window().orderBy(lit('A'))
1195
1196     feature = df.columns[18:] # only regular features
1197
1198     mean_df = df.filter((df["cycle"]<=threshold) & (df["battery_id"].isin(battery_list))).select(feature+["
        battery_id"]).groupBy("battery_id").agg(*[avg(c) for c in feature]).withColumn("row_num", row_number().
        over(w))
1199
1200     c_rpt_df = sc.parallelize(coef_rpt).map(Row("coef_rpt")).toDF().withColumn("row_num", row_number().over(w)
        )
1201     i_rpt_df = sc.parallelize(inter_rpt).map(Row("inter_rpt")).toDF().withColumn("row_num", row_number().over(
        w))
1202     c_reg_df = sc.parallelize(coef_reg).map(Row("coef_reg")).toDF().withColumn("row_num", row_number().over(w)
        )
1203     i_reg_df = sc.parallelize(inter_reg).map(Row("inter_reg")).toDF().withColumn("row_num", row_number().over(
        w))
1204
1205     final_df = c_rpt_df.join(i_rpt_df, c_rpt_df.row_num == i_rpt_df.row_num)\
1206         .join(c_reg_df, c_rpt_df.row_num == c_reg_df.row_num)\
1207         .join(i_reg_df, c_rpt_df.row_num == i_reg_df.row_num)\
1208         .join(mean_df, c_rpt_df.row_num == mean_df.row_num).drop("row_num")
1209
1210     return final_df
1211
1212
1213 # Step 3a GBT & Hyperparameter tuning
1214 def GBT(df, label = 'coef_rpt'):
1215     # assembler
1216     feature = df.columns[2:4] + df.columns[5:]
1217     assembler = VectorAssembler(
1218         inputCols= feature,
1219         outputCol="features")
1220     #gbt_input = assembler.transform(df)
1221     #features_df = gbt_input.select("features")
1222     #pca = PCA(k = num_feature, inputCol="features", outputCol="pca_features")
1223     #pca_model = pca.fit(features_df)
1224     #gbt_df = pca_model.transform(gbt_input)
1225     pca = PCA(inputCol="features", outputCol="pca_features")
1226
1227     # data split
1228     trainSet = df.filter((col("battery_id")!= "A8C2C39") & (col("battery_id")!= "AFFE915"))
1229     testSet = df.filter((col("battery_id")== "A8C2C39") | (col("battery_id")== "AFFE915"))
1230
1231     # gbt
1232     gbt = GBTRegressor(featuresCol = 'pca_features', labelCol = label)
1233     #gbt_model = gbt.fit(trainSet)
1234
1235     pipeline = Pipeline(stages=[assembler,pca,gbt])
1236
1237     # parameter grid
1238     paramGrid = ParamGridBuilder() \
1239         .addGrid(pca.k, [5, 8, 10, 14])\
1240         .addGrid(gbt.maxIter, [3, 5, 10]) \
1241         .addGrid(gbt.maxDepth, [3, 6, 9]) \
1242         .addGrid(gbt.stepSize, [0.01, 0.05, 0.1]) \
1243         .build()
1244
1245     # evaluator
1246     evaluator_reg = RegressionEvaluator(labelCol=label, predictionCol="prediction", metricName="rmse")
1247
1248

```

```

1249 #train-val splid for parameter tuning
1250 tvs = TrainValidationSplit(estimator=pipeline,
1251                             estimatorParamMaps=paramGrid,
1252                             evaluator=evaluator_reg,
1253                             trainRatio = 0.7)
1254
1255 # train the model
1256 tvsModel = tvs.fit(trainSet)
1257
1258 # prediction
1259 gbt_predictions_train = tvsModel.transform(trainSet)
1260 gbt_predictions_test = tvsModel.transform(testSet)
1261
1262 #loss
1263 rmse_train = evaluator_reg.evaluate(gbt_predictions_train)
1264 rmse_test = evaluator_reg.evaluate(gbt_predictions_test)
1265
1266 print("Train Loss: {}\nTest Loss: {}".format(rmse_train,rmse_test))
1267
1268 return tvsModel, gbt_predictions_train, gbt_predictions_test
1269
1270
1271
1272 # Step 3b GBT after Hyperparameter tuning
1273 def GBT_after(df, label, num_feature,ite,depth,step):
1274     # assembler
1275     feature = df.columns[2:4] + df.columns[5:]
1276     assembler_gbt = VectorAssembler(
1277         inputCols= feature,
1278         outputCol="features")
1279     gbt_input = assembler_gbt.transform(df)
1280
1281     # feature reduction
1282     features_df = gbt_input.select("features")
1283     pca = PCA(k = num_feature, inputCol="features", outputCol="pca_features")
1284     pca_model = pca.fit(features_df)
1285     gbt_df = pca_model.transform(gbt_input)
1286
1287     # data split
1288     trainSet = gbt_df.filter((col("battery_id")!= "A8C2C39") & (col("battery_id")!= "AFFE915"))
1289     testSet = gbt_df.filter((col("battery_id")== "A8C2C39") | (col("battery_id")== "AFFE915"))
1290
1291     # gbt
1292     gbt = GBTRegressor(featuresCol = 'pca_features', labelCol = label, maxIter=ite,maxDepth=depth,stepSize=
1293         step)
1294     gbt_model = gbt.fit(trainSet)
1295
1296     # evaluator
1297     evaluator_reg = RegressionEvaluator(labelCol=label, predictionCol="prediction", metricName="rmse")
1298
1299     # prediction
1300     gbt_predictions_train = gbt_model.transform(trainSet)
1301     gbt_predictions_test = gbt_model.transform(testSet)
1302
1303     #loss
1304     rmse_train = evaluator_reg.evaluate(gbt_predictions_train)
1305     rmse_test = evaluator_reg.evaluate(gbt_predictions_test)
1306
1307     print("Train Loss: {}\nTest Loss: {}".format(rmse_train,rmse_test))
1308
1309     return gbt_model, gbt_predictions_train, gbt_predictions_test
1310
1311
1312 # Step 4 calculate sMAPE
1313 def sMAPE_calculation(pred_coef_df, pred_inter_df, battery_id, num):
1314     coef_test = pred_coef_df.select("battery_id", "prediction").withColumnRenamed("prediction", "pred_coef")
1315     inter_test = pred_inter_df.select("battery_id", "prediction").withColumnRenamed("prediction", "pred_inter")

```

```

1316 coef = gbt_pred_test_coef.filter((col("battery_id")==battery_id)).select("prediction").collect()[0]["
1317 prediction"]
1318 inter = gbt_pred_test_inter.filter((col("battery_id")==battery_id)).select("prediction").collect()[0]["
1319 prediction"]
1320 compare = gt.filter(col("battery_id")==battery_id).withColumn("prediction",col("cycle")*coef+inter)
1321
1322 a = np.array([compare.select("RPT_discharge_capacity_pct").collect()[i]["RPT_discharge_capacity_pct"] for
1323 i in range(compare.filter((col("cycle")>50) & (col("cycle")<(50+num))).count())])
1324 b = np.array([compare.select("prediction").collect()[i]["prediction"] for i in range(compare.filter((col("
1325 cycle")>50) & (col("cycle")<(50+num))).count())])
1326
1327 SMAPE = np.sum(np.abs(a-b)/((np.abs(a)+np.abs(b))/2))/len(a)
1328
1329 return compare, SMAPE
1330
1331 # Implementation
1332 # Linear Regression to get coefficient and constant for each battery id
1333 # coef_rpt, inter_rpt, coef_reg, inter_reg = LR(df, battery, eol)
1334 coef_rpt, inter_rpt, coef_reg, inter_reg = LR(df, all_battery, all_eol)
1335
1336 # prepare data for gbt
1337 threshold = 100
1338 lr_df = dataPrepare_gbt(df, all_battery, coef_rpt, inter_rpt, coef_reg, inter_reg, threshold)
1339
1340 # gbt & hyperparameter tuning # running time = 30 mins
1341 #model_coef, train_pred_coef, test_pred_coef = GBT(lr_df, 'coef_rpt')
1342 ##threshold=50: numFeature = 14, maxIter = 3, maxDepth = 6, stepSize = 0.01
1343 ##threshold=100: numFeature = 14, maxIter = 10, maxDepth = 3, stepSize = 0.1
1344
1345 #model_inter, train_pred_inter, test_pred_inter = GBT(lr_df, 'inter_rpt')
1346 ##threshold=50: numFeature = 5, maxIter = 10, maxDepth = 6, stepSize = 0.1
1347 ##threshold=100: numFeature = 8, maxIter = 10, maxDepth = 3, stepSize = 0.1
1348
1349
1350 # use GBT_after() with the best combination of hyperparameters for different threshold or different number
1351 # of forecast
1352 gbt_model_coef, gbt_pred_train_coef, gbt_pred_test_coef = GBT_after(lr_df, 'coef_rpt', 14, 10, 3, 0.1)
1353 gbt_model_inter, gbt_pred_train_inter, gbt_pred_test_inter = GBT_after(lr_df, 'inter_rpt', 8, 10, 3, 0.1)
1354
1355 # SMAPE
1356 num_forecast = [50,100,150,200,250,300,350,400]
1357 for num in num_forecast:
1358     com_df1, SMAPE_1 = SMAPE_calculation(gbt_pred_test_coef, gbt_pred_test_inter, "A8C2C39",num)
1359     com_df2, SMAPE_2 = SMAPE_calculation(gbt_pred_test_coef, gbt_pred_test_inter, "AFFE915",num)
1360     print("number of forecast: {}, A8C2C39 SMAPE: {}, AFFE915 SMAPE: {}, Average SMAPE: {}".format(num, SMAPE_1
1361     , SMAPE_2, (SMAPE_1+SMAPE_2)/2))
1362
1363 # display(gbt_pred_test_coef.select("battery_id","coef_rpt","prediction"))
1364 # display(gbt_pred_test_inter.select("battery_id","inter_rpt","prediction"))
1365 # display(com_df1.select("cycle","RPT_discharge_capacity_pct","prediction"))
1366 # display(com_df2.select("cycle","RPT_discharge_capacity_pct","prediction"))
1367
1368 # Important
1369 # This section is some Scala code utilizing school cluster for this project, not for running purpose because
1370 # it does not run on databricks
1371 # Include data cleaning and processing processes with some similar as the main script in pyspark and some
1372 # different (side individual analysis)
1373 # Linear regression tuning is also included
1374 // the output data set has no null or NaN, and each record represents a cycle 12 features (3 features from
1375 // each soc)
1376 // only cycles with 4 soc regions are included, 37 records from raw_data are removed
1377 import org.apache.spark.sql.Row

```

```

1375 import org.apache.spark.sql.functions._
1376 import org.apache.spark.sql.expressions.Window
1377 import org.apache.spark.ml.regression.LinearRegression
1378 import org.apache.spark.ml.linalg.{Vector, Vectors}
1379 import org.apache.spark.ml.param.ParamMap
1380 import org.apache.spark.sql.Row
1381 import org.apache.spark.ml.feature.VectorAssembler
1382
1383 // the given data set from client
1384 val df_raw = spark.table("gbatteries_team_8.raw_data")
1385
1386 val df_1 = df_raw.groupBy("battery_id", "cycle").agg(sum("soc_region") as "soc_region_sum", count(lit(1)) as
    "soc_region_count")
1387 val df_2 = df_1.filter($"soc_region_sum" != 10 || $"soc_region_count" != 4)
1388
1389 // soc_region_sum
1390 df_1.groupBy("soc_region_sum", "soc_region_count").agg(count(lit(1))).show()
1391 df_2.groupBy("soc_region_sum", "soc_region_count").agg(count(lit(1))).show()
1392
1393 df_2.orderBy("battery_id", "cycle").show(df_2.count().asInstanceOf[Int], false)
1394
1395 // no null or NaN values
1396 val df_3 = df_raw.na.drop("any") // drop all records with any entry being null or NaN
1397 df_3.count()
1398
1399 df_raw.select("battery_id").distinct().count() // 48 batteries
1400 val df_cleaned = df_raw.filter($"battery_id" != "A2A5795").filter($"battery_id" != "ACF5ADD" || $"cycle"
    != 215)
1401
1402 // df_cleaned has no null or NaN, and all cycles have 4 soc
1403 df_cleaned.count() // 31760
1404 df_raw.count() // 32266
1405 df_raw.filter($"battery_id" === "A2A5795").count() // 503
1406 df_raw.filter($"battery_id" === "ACF5ADD" && $"cycle" === 215).count() // 3
1407
1408 // check if all batteries have RPT cycles
1409 df_cleaned.select("battery_id").distinct().count() // 47
1410 df_cleaned.filter($"cycle_type" === "RPT").select("battery_id").distinct().count() // 22
1411 // there are some batteries with no RPT cycles
1412
1413 // number of cycles
1414 df_cleaned.select("battery_id", "cycle").distinct().count() // 7940, match 31760/4
1415
1416 //
1417 val df_cycles = df_cleaned.select("battery_model", "battery_id", "cycle", "cycle_type", "charge_duration_sec", "
    discharge_capacity_pct").distinct()
1418
1419 val df_soc_1 = df_cleaned.select("battery_model", "battery_id", "cycle", "feature_1", "feature_2", "feature_3").
    filter($"soc_region" === 1).withColumnRenamed("feature_1", "soc1_feature_1").withColumnRenamed("feature_2",
    "soc1_feature_2").withColumnRenamed("feature_3", "soc1_feature_3")
1420
1421 val df_soc_2 = df_cleaned.select("battery_model", "battery_id", "cycle", "feature_1", "feature_2", "feature_3").
    filter($"soc_region" === 2).withColumnRenamed("feature_1", "soc2_feature_1").withColumnRenamed("feature_2",
    "soc2_feature_2").withColumnRenamed("feature_3", "soc2_feature_3")
1422
1423 val df_soc_3 = df_cleaned.select("battery_model", "battery_id", "cycle", "feature_1", "feature_2", "feature_3").
    filter($"soc_region" === 3).withColumnRenamed("feature_1", "soc3_feature_1").withColumnRenamed("feature_2",
    "soc3_feature_2").withColumnRenamed("feature_3", "soc3_feature_3")
1424
1425 val df_soc_4 = df_cleaned.select("battery_model", "battery_id", "cycle", "feature_1", "feature_2", "feature_3").
    filter($"soc_region" === 4).withColumnRenamed("feature_1", "soc4_feature_1").withColumnRenamed("feature_2",
    "soc4_feature_2").withColumnRenamed("feature_3", "soc4_feature_3")
1426
1427 val df_combined_cycles = df_cycles.join(df_soc_1, Seq("battery_model", "battery_id", "cycle")).join(df_soc_2,
    Seq("battery_model", "battery_id", "cycle")).join(df_soc_3, Seq("battery_model", "battery_id", "cycle")).
    join(df_soc_4, Seq("battery_model", "battery_id", "cycle"))
1428
1429 val df_output = df_combined_cycles.orderBy($"battery_model", $"battery_id", $"cycle")

```

```

1430 df_output.coalesce(1).write.saveAsTable("gbatteries_team8.cleaned_data")
1431 df_output.coalesce(1).write.option("header", "true").csv("cleaned_data.csv")
1432
1433
1434 // there are some data weird
1435 // step 1: check if a battery has enough data
1436 val df_battery_cycles_count = df_combined_cycles.groupBy("battery_id").agg(count(lit(1)) as "battery_count",
    count(when($"cycle_type"=="RPT",1)) as "RPT_count", count(when($"cycle_type"=="Regular",1)) as "
    Regular_count").orderBy("battery_id")
1437
1438 df_battery_cycles_count.show()
1439
1440 df_battery_cycles_count.coalesce(1).write.saveAsTable("gbatteries_team8.battery_cycles_count")
1441 df_battery_cycles_count.coalesce(1).write.option("header", "true").csv("battery_cycles_count.csv")
1442
1443 // step 2: separate the data set in terms of cycle type to observe trend
1444 val df_output_RPT = df_combined_cycles.filter($"cycle_type" === "RPT").select("battery_model", "battery_id", "
    cycle", "discharge_capacity_pct").orderBy($"battery_model", $"battery_id", $"cycle")
1445 val df_output_Regular = df_combined_cycles.filter($"cycle_type" === "Regular").select("battery_model", "
    battery_id", "cycle", "discharge_capacity_pct").orderBy($"battery_model", $"battery_id", $"cycle")
1446
1447 df_output_RPT.coalesce(1).write.saveAsTable("gbatteries_team8.cleaned_data_rpt")
1448 df_output_RPT.coalesce(1).write.option("header", "true").csv("cleaned_data_RPT.csv")
1449
1450 df_output_Regular.coalesce(1).write.saveAsTable("gbatteries_team8.cleaned_data_regular")
1451 df_output_Regular.coalesce(1).write.option("header", "true").csv("cleaned_data_Regular.csv")
1452
1453 val rpt_raw = spark.table("gbatteries_team8.cleaned_data_rpt").select("battery_id", "cycle", "
    discharge_capacity_pct")
1454
1455 // select only batteries having RPT cycles
1456 val reg_raw = spark.table("gbatteries_team8.cleaned_data_regular").select("battery_id", "cycle").filter($"
    battery_id".isin("A8C2C39", "A65D895", "A2C0E4A", "AFFE915", "A947A53", "A1D35B0", "A74A962", "AFD5A8F", "
    AEBE15E", "A716092", "A18242C", "A4E5D13", "AEE8EF6", "A231712", "A6AD931", "A35D29D", "AC3C95D", "A83B987", "
    A73CF24", "A0047B2", "AD5B491", "A43046B"))
1457
1458 // take last RPT cycle of consecutive RPT cycles as measurement (and label whether it is from a consecutive
    set)
1459 // create column of next RPT cycle number and column of last RPT cycle number
1460 val battery_id_window = Window.partitionBy($"battery_id").orderBy("cycle")
1461 val lead_cycle = lead("cycle", 1, 0).over(battery_id_window)
1462 val rpt_temp = rpt_raw.select($"*", lead_cycle as "next_RPT_cycle")
1463 val last_cycle = lag("cycle", 1, 0).over(battery_id_window)
1464 rpt_temp = rpt_temp.select($"*", last_cycle as "last_RPT_cycle")
1465 // only keep records if next cycle is not RPT, and label if its last cycle is RPT
1466 rpt_temp = rpt_temp.filter($"next_RPT_cycle" - $"cycle" != 1)
1467 rpt_temp = rpt_temp.withColumn("consecutive_rpt", when(($"cycle" - $"last_RPT_cycle" === 1) && ($"cycle" != 1), 1).otherwise(0))
1468
1469 // eliminate regular data after reaching last RPT cycle
1470 val rpt_last_cycle = rpt_temp.groupBy("battery_id").agg(max("cycle") as "last_rpt_cycle")
1471 val reg_temp = reg_raw.join(rpt_last_cycle, "battery_id")
1472 reg_temp = reg_temp.filter($"cycle" < $"last_rpt_cycle")
1473
1474 // exporting data
1475 val rpt_final = rpt_temp.select("battery_id", "cycle", "discharge_capacity_pct", "consecutive_rpt").orderBy(
    "battery_id", "cycle")
1476 val reg_final = reg_temp.select("battery_id", "cycle").orderBy("battery_id", "cycle")
1477
1478 // rpt actual data for computing smape
1479 val rpt_eol_3 = rpt_final.select("battery_id", "cycle", "discharge_capacity_pct")
1480
1481 rpt_eol_3.coalesce(1).write.saveAsTable("gbatteries_team8.rpt_actual_v8")
1482 rpt_eol_3.coalesce(1).write.option("header", "true").csv("rpt_actual_v8.csv")
1483
1484
1485 val df_rpt_raw_data = spark.table("gbatteries_team8.cleaned_data_rpt")
1486 val df_rpt_selected_data = df_rpt_raw_data.filter($"battery_id".isin("AEBE15E", "AFD5A8F", "A74A962", "A1D35B0")

```



```

    , "A947A53", "AFFE915", "A2C0E4A", "A65D895", "A8C2C39"))
1487
1488 // verify data counts
1489 val check_1 = df_rpt_selected_data.groupBy("battery_id").agg(count(lit(1)) as "RPT_count").orderBy("
    battery_id")
1490
1491 // create column of next RPT cycle number
1492 val lead_window = Window.partitionBy($"battery_id").orderBy("cycle")
1493 val lead_cycle = lead("cycle", 1, 0).over(lead_window)
1494 val df1 = df_rpt_selected_data.select($"*", lead_cycle as "next_RPT_cycle")
1495 df1.select("cycle", "next_RPT_cycle").show(25)
1496
1497 // next_RPT_cycle is 0 for final row, which is different from the null in Pyspark
1498 // remove consecutive RPT cycles
1499 val df2 = df1.filter($"next_RPT_cycle" - $"cycle" != 1).drop("next_RPT_cycle")
1500 df2.show(25)
1501
1502 // split into training and testing set
1503 val df_output = df2.select("battery_id", "cycle", "discharge_capacity_pct")
1504 val df_output_training = df2.filter($"battery_id".isin("AEBE15E", "AFD5A8F", "A74A962", "A1D35B0", "A947A53", "
    A2C0E4A", "A65D895")).filter($"cycle" < 120)
1505 val df_output_testing = df2.filter($"battery_id".isin("A8C2C39", "AFFE915"))
1506
1507 // export data sets
1508 df_output_training.coalesce(1).write.saveAsTable("gbatteries_team.8.lin_reg_training")
1509 df_output_training.coalesce(1).write.option("header", "true").csv("lin_reg_training.csv")
1510
1511 df_output_testing.coalesce(1).write.saveAsTable("gbatteries_team.8.lin_reg_testing")
1512 df_output_testing.coalesce(1).write.option("header", "true").csv("lin_reg_testing.csv")
1513
1514
1515 // This part is tuning for threshold value for step 2 (deviation from line of training points) of enhanced
    linear regression
1516 val df_training = spark.table("gbatteries_team.8.lin_reg_training")
1517
1518 // perform linear regression to each set of early RPT points to see how much the points deviate from the
    line generated
1519
1520 // boundary 0.05 eliminate a minor outlier but 0.043 would make the slope +
1521 val df1 = df_training.filter($"battery_id" === "AEBE15E").select("cycle", "discharge_capacity_pct")
1522
1523 // boundary 0.04 cycle65 for df2 would make perfect line
1524 val df2 = df_training.filter($"battery_id" === "AFD5A8F").select("cycle", "discharge_capacity_pct").filter($"
    cycle" != 65)
1525
1526 val assembler = new VectorAssembler().setInputCols(Array("cycle")).setOutputCol("features")
1527
1528 val df1_input = assembler.transform(df1).drop("cycle").withColumnRenamed("discharge_capacity_pct", "label")
1529 val lr = new LinearRegression()
1530 val model1 = lr.fit(df1_input)
1531 val df1_output = model1.transform(df1_input).withColumn("label - prediction", $"label" - $"prediction")
1532 df1_output.show()
1533
1534 val df2_input = assembler.transform(df2).drop("cycle").withColumnRenamed("discharge_capacity_pct", "label")
1535 val lr = new LinearRegression()
1536 val model2 = lr.fit(df2_input)
1537 val df2_output = model2.transform(df2_input).withColumn("label - prediction", $"label" - $"prediction")
1538 df2_output.show()

```