

VIETNAM NATIONAL UNIVERSITY, HO CHI MINH CITY
UNIVERSITY OF TECHNOLOGY
FACULTY OF COMPUTER SCIENCE AND ENGINEERING



**BIG DATA ANALYTICS
AND BUSINESS INTELLIGENCE
Semester 221**

**Build a model
to predict asteroid size using
Spark and Google Cloud Service**

Teacher: Assoc. Prof. Nguyen Thanh Binh
Student: Le Tien Loc - 1952328
Phan Tuan Khai - 1952780
Nguyen Van Quoc Chuong - 1950004

HO CHI MINH CITY, DECEMBER 2022



Contents

1	Project Description	2
1.1	Requirement	2
1.2	Our topic	2
1.3	Big data tools	2
1.4	Notice	2
2	Google Cloud Service - Data preparation	3
2.1	Dataset	3
2.2	Uploading dataset to Google Cloud Storage	4
2.3	Create a Service Account to get access to the bucket	5
3	SparkSession - Google Cloud Storage connection	5
4	SparkSQL - data visualization and processing	6
4.1	Data visualization	6
4.1.1	Bar Chart for counting number of asteroids base on its diameter	6
4.1.2	Bar Chart for average diameter by asteroid's classes	7
4.2	Data processing	8
5	SparkML - ML model and evaluation	10
6	Summary	13



1 Project Description

1.1 Requirement

This project has to meet the following requirements as what the teacher Thanh Binh said in class:

- Make use of big data tools/frameworks to process, visualize data
- Datasets could be any size but should be as large as possible

1.2 Our topic

In this project, our group will train a model to predict asteroids size using a dataset of asteroids' observations provided by NASA. (about 1.2 million records)

1.3 Big data tools

After several discussions, we decided to use Google Cloud Services, which is a storage for big data and Spark, a very popular framework for data analyst to process data. Moreover, Spark also provides us a library named SparkML to build some basic ML models, hence, we will make use of it for building model in our project. The tools we used will be summarized in the following table:

Tools/Frameworks	Purposes
Google Cloud Storage	Store Data
SparkSQL	Data visualization and processing
SparkML	Build ML model and evaluation

Table 1: Big data tools

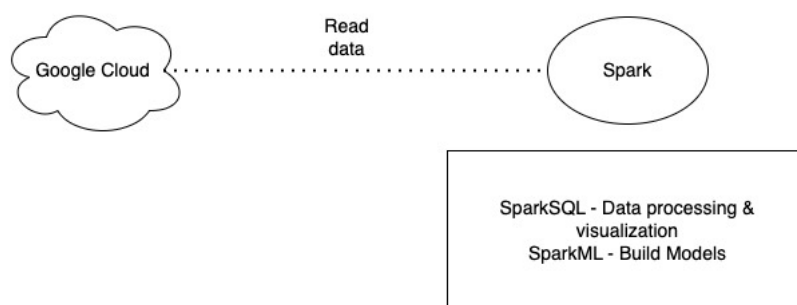


Figure 1

1.4 Notice

In this project, the path we use is absolute path which may not work in your workspace. To make the code run you have to make sure the path leading to correct required files. The secrets.JSON is the file for authentication so it will not be uploaded or the value will be changed. You can generate the secrets.JSON by yourself by creating a key for your Service Account in Google Cloud Service.

2 Google Cloud Service - Data preparation

2.1 Dataset

We use the datasets provided by NASA [NAS], which comprises of about 1.2 millions records with 31 fields in total. We think this datasets is large enough for us to practice working with big data tools/frameworks

The 31 fields information of dataset:

- a: semi-major axis (au)
- e: eccentricity
- i: inclination with respect to ecliptic plane
- om: longitude of the ascending node
- w: argument of perihelion
- q: perihelion distance (au)
- ad: aphelion distance (au)
- per_y: orbital period (years)
- data_arc: span of recorded data (days)
- condition_code: orbital condition code
- n_obs_used: number of observations used
- H: absolute magnitude parameter
- neo: near-earth object
- pha: physically hazardous object
- diameter: diameter (target variable)
- extent: Object bi/tri axial ellipsoid dimensions(Km)
- albedo: geometric albedo
- rot_per: rotation period (hours)
- GM: gravitational parameter. Product of mass and gravitational constant
- BV: Color index B-V magnitude difference
- UB: Color index U-B magnitude difference
- IR: Color index I-R magnitude difference
- spec_B: Spectral taxonomic type(SMASSII)
- spec_T: Spectral taxonomic type (Tholen)
- G: Magnitude slope parameter
- moid: Earth minimum orbit intersection distance

- class: asteroid orbit class
- n: mean motion (degrees/day)
- per: orbital period (days)
- ma: mean anomaly (degrees)

2.2 Uploading dataset to Google Cloud Storage

To use Google Cloud Service, we have to sign up with a credit card. Luckily, Google will give us \$300 free credit for 3 months. After successful registry, we have to create a bucket (figure 2), which is a place to store data in Google Cloud Storage. We then create a folder and upload our dataset to that folder in bucket.

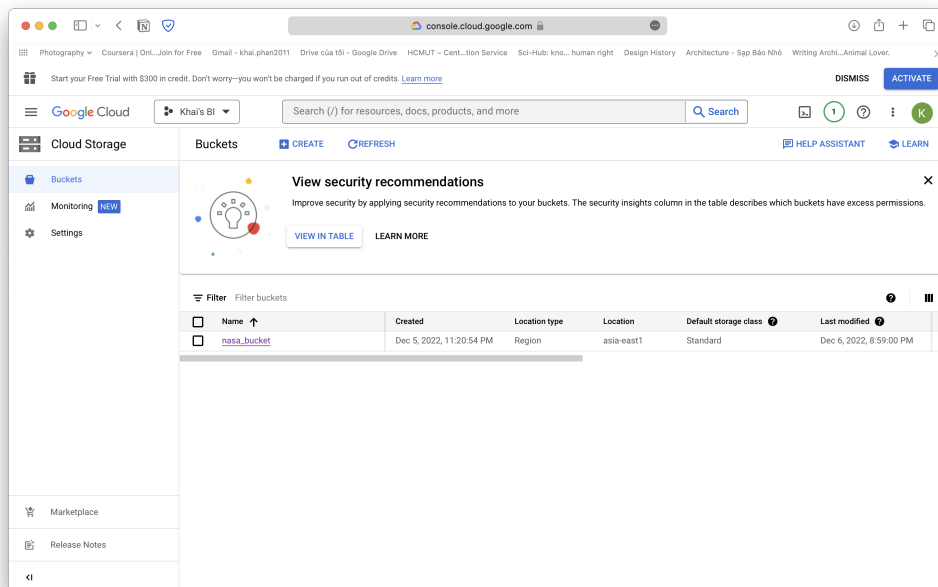


Figure 2: Created bucket



2.3 Create a Service Account to get access to the bucket

After creating a new bucket, we have to create a Service Account which will be use to get access to the data bucket.

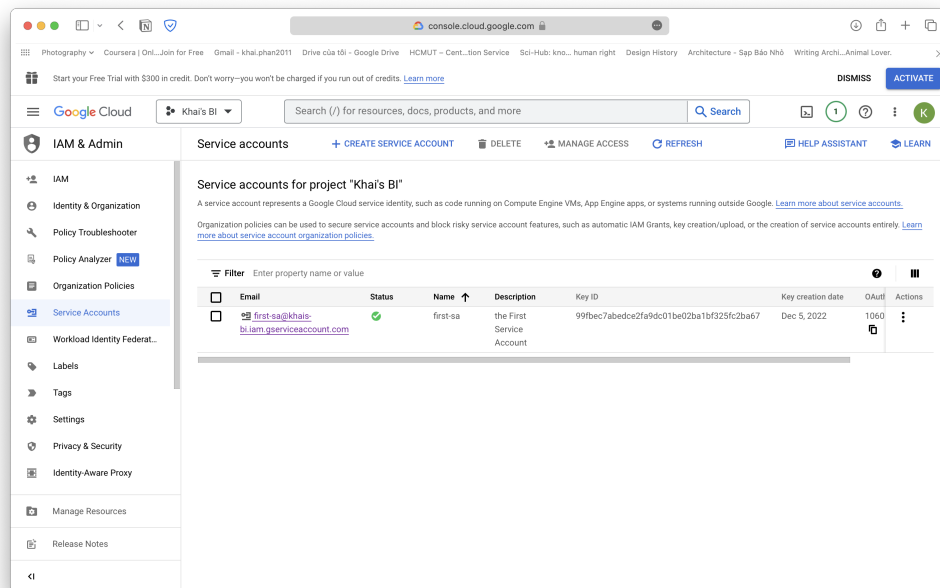


Figure 3: Service Account

With the created Service Account, we will create a key stored in a JSON file which will be renamed as secrets.JSON and we will not upload this file due to security. We will use this JSON file to for authentication and get access to the data bucket.

3 SparkSession - Google Cloud Storage connection

We now have the data in Google Cloud Storage but Spark have not had a connection to it yet. We have to configure the Spark Session so that it can connect to Google Cloud Storage. We will need a connector to connect Spark to Google Cloud Service which will be store in file gcs-connector-hadoop3-latest.jar. When we have the connector file, we have 2 main Spark configurations to make Spark work with Google Cloud.

- Setup hadoop fs configuration for schema gs://
- Setup hadoop fs configuration for Google Cloud authentication using secrets.JSON

```
1 #create a Spark Session with Google Cloud connector
2 spark = SparkSession \
3     .builder \
4     .appName("GCS_NASA") \
5     .config("spark.jars", "drive/Shareddrives/Big_Data_and_BI/connector/
6     gcs-connector-hadoop3-latest.jar") \
7     .getOrCreate()
```

```
8 # Setup hadoop fs configuration for schema gs://
9 conf = spark.sparkContext._jsc.hadoopConfiguration()
10 conf.set("fs.gs.impl", "com.google.cloud.hadoop.fs.gcs.
    GoogleHadoopFileSystem")
11 conf.set("fs.AbstractFileSystem.gs.impl", "com.google.cloud.hadoop.fs.
    gcs.GoogleHadoopFS")
12
13 # Google Cloud Authentication
14 conf.set("fs.gs.auth.type", "SERVICE_ACCOUNT_JSON_KEYFILE")
15 conf.set("fs.gs.auth.service.account.json.keyfile", "drive/Shareddrives/
    Big_Data_and_BI/key/secrets.json")
```

Now, we can read data from Google Cloud Storage bucket. If it works, the following code will not throw any error.

```
1 bucket_name="nasa_bucket"
2 path=f"gs://{bucket_name}/asteroid/Asteroid_Updated.csv"
3
4 df=spark.read.csv(path, sep=',', inferSchema=True, header=True)
```

4 SparkSQL - data visualization and processing

4.1 Data visualization

In this part, we will visualize data focusing on the asteroid diameter, we do not have much knowledge on this field so the way we make classes for diameter size may not practical in real life. This part is done after the data is cleaned by dropping na, which is on the Data processing part.

4.1.1 Bar Chart for counting number of asteroids base on its diameter

Firstly, I check for the range of diameter.

```
1 astro_df.agg(min(col("diameter")), max(col("diameter"))).show()
```

And we got the following result, the range is from 0.0025 to 939.4

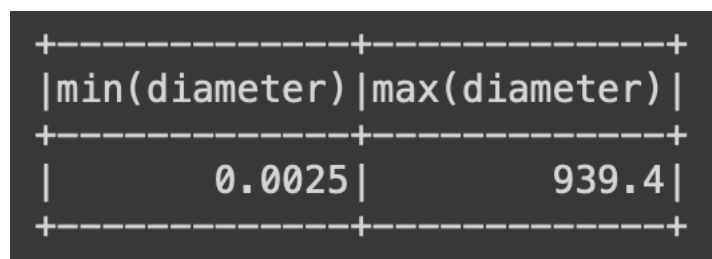


Figure 4: Diameter Range

With the range above, we tried to divide diameter to 5 class (0,200), (200, 400), (400, 600), (600, 800), (800, 1000). However, we face the problem that the number of asteroid with diameter from 0 to 200 is considerable larger than the other classes. After several trial, we have specify another class as following.

```
1 count_list = []
2 label = []
3 for i in range (0, 10, 2):
4     count_list += [astro_df.filter( (col("diameter") > (i)) & (col("
        diameter") < (i+ 2)) ).count()]
```

```
5 label += [f"{i} - {i + 2}"]
6
7 for i in range (10, 50, 10):
8     count_list += [astro_df.filter( (col("diameter") > (i)) & (col("
9         diameter") < (i+ 10)) ).count()]]
10    label += [f"{i} - {i + 10}"]
11
12 for i in range (50, 200, 50):
13     count_list += [astro_df.filter( (col("diameter") > (i)) & (col("
14         diameter") < (i+ 50)) ).count()]]
15    label += [f"{i} - {i + 50}"]
16
17 for i in range (200, 1000, 400):
18     count_list += [astro_df.filter( (col("diameter") > (i)) & (col("
19         diameter") < (i+ 400)) ).count()]]
20    label += [f"{i} - {i + 400}"]
21
22 diameter_data = {'label': label, 'data': count_list}
23 diameter_df = pd.DataFrame(data=diameter_data)
24 diameter_df
```

With the data frame, we use Pandas to plot a bar chart

```
1 # set width of bar
2 barWidth = 0.5
3 fig = plt.subplots(figsize =(16, 8))
4
5 # Set position of bar on X axis
6 X_axis = np.arange(len(diameter_df["label"]))
7
8 # Make the plot
9 count = 0
10 color_list = ["#3981BF", "#38A649", "#F2AE30", "#F2594B", "#F2594C"]
11
12 plt.bar(X_axis + barWidth*(0.5+ count), diameter_df["data"], color =
13     color_list[count], width = barWidth,
14     edgecolor = 'grey')
15
16 # Adding Xticks
17 plt.xlabel('Diameter', fontweight = 'bold', fontsize = 15)
18 plt.ylabel('Number of asteroids', fontweight = 'bold', fontsize = 15)
19 plt.xticks(X_axis + 0.5*barWidth,diameter_df["label"])
20
21 plt.legend()
22 plt.show()
```

And we got:

4.1.2 Bar Chart for average diameter by asteroid's classes

Firstly, we calculate the average diameter grouped by class column.

```
1 avrg_pd = astro_df.groupBy(col("class")).agg( avg(col("diameter")).alias
2     ("average")).toPandas()
```

With the data frame, we use Pandas to plot a bar chart

```
1 # set width of bar
2 barWidth = 0.5
3 fig = plt.subplots(figsize =(16, 8))
4
5 # Set position of bar on X axis
6 X_axis = np.arange(len(avrg_pd["class"]))
7
8 # Make the plot
9 count = 0
```

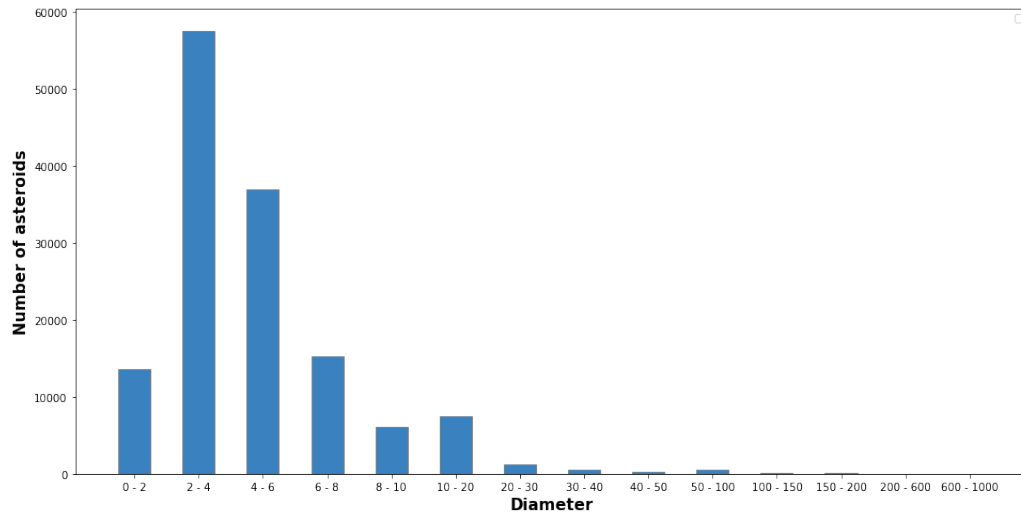



Figure 5: Bar Chart showing number of asteroids by Diameter ranges

```

10
11 plt.bar(X_axis + barWidth*(0.5+ count), avrg_pd["average"], color = "#3981BF", width = barWidth,
12         edgcolor = 'grey')
13
14
15
16 # Adding Xticks
17 plt.xlabel('Class', fontweight = 'bold', fontsize = 15)
18 plt.ylabel('Average Diameter', fontweight = 'bold', fontsize = 15)
19 plt.xticks(X_axis + 0.5*barWidth, avrg_pd["class"])
20
21 plt.legend()
22 plt.show()

```

And we got:

From this chart, we know that the diameter of asteroid in TNO class is considerably larger than the others

4.2 Data processing

First of all, we need to handle the missing data by using `dropna()` function in pandas.

```

1
2 astro_df = df.dropna(subset=['diameter'])
3 drop_list = ["extent",
4             "rot_per",
5             "GM", "BV", "UB", "IR",
6             "spec_B",
7             "spec_T",
8             "G", 'data_arc', 'H', 'albedo']
9 astro_df = astro_df.drop(*drop_list)

```

Take note that we modified the True/False values in the 'neo' and 'pha' columns before converting them to boolean values and then casting them as integers. Simply put, this is binary encoding. These columns' True/False values are now shown as 1s

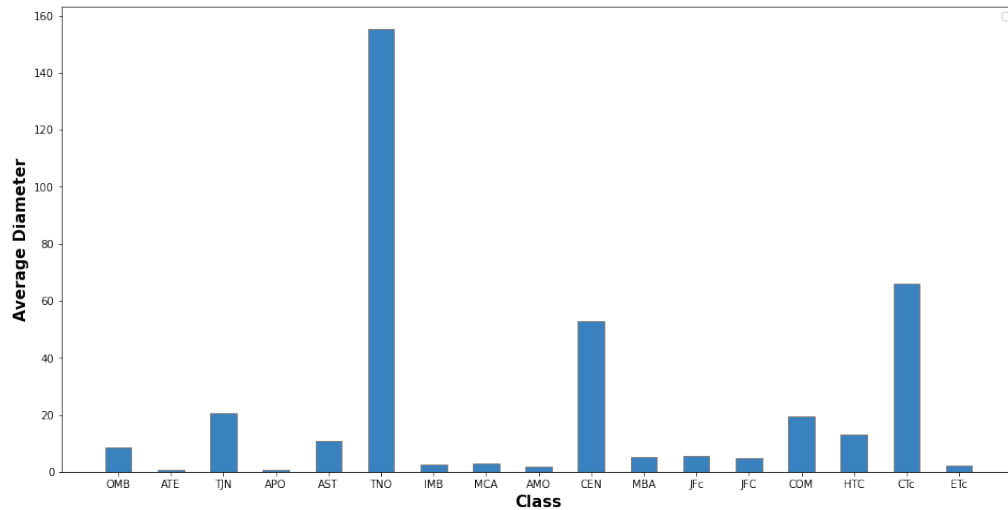


Figure 6: Bar Chart showing average diameter by Asteroids' Class

and 0s. Since many models in SparkML need 'double' datatypes as input, we also updated floating in this code to doubles.

```

1 astro_df = astro_df.withColumn('neo', regexp_replace('neo', 'Y', 'True')
2 )
3 astro_df = astro_df.withColumn('neo', regexp_replace('neo', 'N', 'False'
4 )
5 astro_df = astro_df.withColumn('pha', regexp_replace('pha', 'Y', 'True')
6 )
7 astro_df = astro_df.withColumn('pha', regexp_replace('pha', 'N', 'False'
8 )
9 astro_df = astro_df.withColumn('n_obs_used', astro_df['n_obs_used'].cast(
10 'double'))
11 astro_df = astro_df.withColumn('diameter', astro_df['diameter'].cast('
12 double'))
13
14 for column in ['neo', 'pha']:
15     astro_df = astro_df.withColumn(column, astro_df[column].cast('
16     boolean').cast('int'))
17
18 astro_df = astro_df.dropna(subset=['diameter'])

```

Utilizing Spark's StringIndexer, OneHotEncoder, and Pipeline classes, we have to encode categorical variables, then, we can drop the original columns and re-order our dataframe.

```

1 condition_code_indexer = StringIndexer(inputCol="condition_code",
2     outputCol="condition_codeIndex").setHandleInvalid("skip")
3 class_indexer = StringIndexer(inputCol="class", outputCol="classIndex").
4     setHandleInvalid("skip")
5 onehotencoder_condition_code_vector = OneHotEncoder(inputCol="
6     condition_codeIndex", outputCol="condition_code_vec")
7 onehotencoder_class_vector = OneHotEncoder(inputCol="classIndex",
8     outputCol="class_vec")
9
10 encoding_pipeline = Pipeline(stages=[condition_code_indexer,
11     class_indexer,
12     onehotencoder_condition_code_vector,

```

```
10         onehotencoder_class_vector
11     ])
12
13 astro_df = encoding_pipeline.fit(astro_df).transform(astro_df)
14
15 astro_df = astro_df.drop('condition_code', 'class', 'condition_codeIndex',
16                          'classIndex')
17
18 astro_df = astro_df.select('full_name', "a","e","i","om","w","q","ad", '
19                          per_y', 'n_obs_used', 'neo', 'pha', 'moid', 'n', 'per', 'ma', '
20                          condition_code_vec', 'class_vec', 'diameter')
21
22 features = astro_df.schema.names[1:-1]
```

After that, we split our data into training, validation, and testing sets. All of our features must be consolidated into a single vector as well. One of the stranger features of building models with Spark is this. To input all of the feature columns into a single vector for a Spark model, we need to use the `VectorAssembler()` class.

```
1 train, val, test = astro_df.randomSplit([0.6, 0.2, 0.2], seed=42)
2
3 train_df = train.drop('name')
4 val_df = val.drop('name')
5 test_df = test.drop('name')
6
7 assembler = VectorAssembler(inputCols=features, outputCol='features').
8     setHandleInvalid("skip")
9
10 test_pack = assembler.transform(test_df)
11 train_pack = assembler.transform(train_df)
12 val_pack = assembler.transform(val_df)
13
14 for field in features:
15     test_pack = test_pack.drop(field)
16     train_pack = train_pack.drop(field)
17     val_pack = val_pack.drop(field)
```

We use the `MinMaxScaler` to normalize our features for training section.

```
1 scaler = MinMaxScaler(inputCol="features", outputCol="scaledFeatures")
2 scaler_model = scaler.fit(train_pack)
3
4 train_pack = scaler_model.transform(train_pack)
5 val_pack = scaler_model.transform(val_pack)
6 test_pack = scaler_model.transform(test_pack)
```

5 SparkML - ML model and evaluation

Now that we're all ready for modelling, let's instantiate, fit, and make predictions with a few different regression methods available in SparkML:

```
1 gbt = GBTRegressor(featuresCol='features', labelCol='diameter', maxIter
2     =100, maxDepth=5, seed=42, lossType='squared', stepSize=.1)
3
4 gbt_model = gbt.fit(train_pack)
5 gbt_pred = gbt_model.transform(val_pack)
6
7 rf = RandomForestRegressor(featuresCol='features', labelCol='diameter',
8     maxDepth=5, seed=42, bootstrap=True, numTrees=100)
9
10 rf_model = rf.fit(train_pack)
11 rf_pred = rf_model.transform(val_pack)
```

```
10
11 lr = LinearRegression(featuresCol='features', labelCol='diameter',
12                        maxIter=100, loss='squaredError', elasticNetParam=0.5, regParam=0.1,
13                        fitIntercept=True, standardization=True, solver='auto', tol=.1)
14 lr_model = lr.fit(train_pack)
15 lr_pred = lr_model.transform(val_pack)
```

	rmse	r2
LR	13.708287	-1.376717
RF	6.683033	0.435117
GBT	6.590959	0.450575

Now let's evaluate the performance of these three models:

```
1 rmse = RegressionEvaluator(
2     labelCol="diameter", predictionCol="prediction", metricName="rmse")
3
4 r2 = RegressionEvaluator(
5     labelCol="diameter", predictionCol="prediction", metricName="r2")
6
7 metrics = [rmse, r2]
8 metric_labels = ['rmse', 'r2']
9
10 predictions = [lr_pred, rf_pred, gbt_pred]
11 predict_labels = ['LR', 'RF', 'GBT']
12
13 eval_list = list()
14
15 for pred in zip(predict_labels, predictions):
16     name = pred[0]
17     predict = pred[1]
18
19     metric_vals = pd.Series(dict([(x[0], x[1].evaluate(predict))
20                                for x in zip(metric_labels, metrics)]),
21                             name=name)
22     eval_list.append(metric_vals)
23
24 eval_df = pd.concat(eval_list, axis=1).T
25 eval_df = eval_df[metric_labels]
26 eval_df
```

Based on the results, Gradient-Boosted Regression performed best with R2 score of .47, and a Root Mean Squared Error of 7.1 km, so let's try to optimize the model with GridSearch and cross-validation:

```
1 from pyspark.ml.tuning import ParamGridBuilder, CrossValidator
2
3 # the GradientBoostedTree model
4 gbt2 = GBRegressor(featuresCol='features',
5                    labelCol='diameter',
6                    predictionCol='prediction',
7                    seed=13,
8                    lossType='squared',
9                    maxIter=100,
10                    )
11
```

```
12 pipeline = Pipeline(stages=[gbt2])
13
14 paramgrid = (ParamGridBuilder().addGrid(gbt2.maxDepth, [2,4,6])
15         .addGrid(gbt2.stepSize, [0.0001, 0.001,
16         0.01, 0.1, 1.0]).build())
17
18 evaluator = RegressionEvaluator(labelCol='diameter',
19         predictionCol='prediction'
20         ,
21         metricName='r2')
22
23 crossval = CrossValidator(estimator=pipeline,
24         estimatorParamMaps=paramgrid,
25         evaluator=evaluator,
26         numFolds=3)
27 gbt2_model = crossval.fit(train_pack)
```

Now we need to extract our best model from the cross-validator. We can also see the feature importance of our model with the code below:

```
1
2 best_pipeline = gbt2_model.bestModel
3 best_gbt_model = best_pipeline.stages[0]
4
5 feature_importances = best_gbt_model.featureImportances.toArray()
6
7
8 feature_names = astro_df.columns[:-1]
9
10 feature_series = (pd.Series(dict(zip(features, feature_importances)))
11         .sort_values(ascending=True))
12
13 feature_series
```

```
class_vec      0.000000
condition_code_vec  0.000000
pha            0.000085
w              0.001119
neo            0.001324
e              0.002504
om             0.005929
ma             0.005932
moid           0.007812
q              0.030492
n              0.031186
ad             0.051679
per            0.058555
per_y          0.069070
a              0.074710
i              0.089795
n_obs_used     0.481063
dtype: float64
```

Finally, let's make predictions on our test set using the optimized GBT model and see how it compares with our initial model results:

```
1
2 gbt_pred_test = best_pipeline.transform(test_pack)
3
4 predictions = [lr_pred, rf_pred, gbt_pred, gbt_pred_test]
5 predict_labels = ['LR', 'RF', 'GBT', 'GBT_GridSearch']
6
7 eval_list = []
```



```
8 for pred in zip(predict_labels, predictions):
9     name = pred[0]
10    predict = pred[1]
11
12    metric_vals = pd.Series(dict([(x[0], x[1].evaluate(predict))
13                                for x in zip(metric_labels, metrics)]),
14                             name=name)
15    eval_list.append(metric_vals)
16
17
18 eval_df = pd.concat(eval_list, axis=1).T
19 eval_df = eval_df[metric_labels]
20 eval_df
```

	rmse	r2
LR	13.708287	-1.376717
RF	6.683033	0.435117
GBT	6.590959	0.450575
GBT_GridSearch	5.170815	0.482980

6 Summary

In this mini-project, I set out to provide a conceptual understanding of Spark and Google Cloud Storage, and a practical example of creating a machine learning model using both.

As the results of the model show, some additional work could be done on the data (hint: outliers), but the point was not to create an optimal model, merely to show the workflow for someone who might be starting out.

The way we approach working with large datasets is different than we might be used to if we're just making the transition from learning to practitioner. As we make that change, it's important that we start thinking about the optimal solution path for your specific problem — there are plenty of other ways we could have approached this task, such as sampling the data. Often, the tools you will have at your disposal will depend on the company you're working for, so it pays to be flexible, and have a working understanding of a wide range of tools



References

- [NAS] NASA. *Orbital and/or physical parameters for asteroids and/or comets of interest from our small-body database (SBDB)*. URL: https://ssd.jpl.nasa.gov/tools/sbdb_query.html.