

COMP9313 Big Data Management

Project 2

T2, 2020

Name : Yuchen Yang

Student number : z5189310

1. Evaluation of your stacking model on the test data.

This is how I think about the design of the entire model. This project mainly investigates the use of pipeline and how to better use Estimator.

For the data preprocessing of Test1.1, I mainly used Estimator function(NaiveBayes, StringIndexer and CountVectorizer) to train the whole original DataFrame with FIT operation, changing the word in descript to a single and listing a new column according to features, The transform operation with Tokenizer is then used to generate New DataFrame.

After the operation of `gen_binary_labels`, I got the training_set of 5 groups and 10 columns of data among Task1.2 input. According to the generate meta feature concept, I divided the original data into 5 groups again by using the 5 groups that have been classified by the project2. Taking group1 as an example, I divide group2, group3, group4 and group5 into one group as base_model, and group1 as generate meta feature. The same is true of other groups. I use base_model to train the features and label in group2-5 again to generate the new pipeline. Then, generate a meta feature into a new pipeline to generate a new prediction column. Finally, the label `nb_pred` and `svm_pred` were joined to get the final result.

As for Task1.3, I passed the meta_features and label I had created into the Meta Model Logistic Regression for training and prediction model generation.

Reflecting on the whole process of the project, I tried my best to avoid reducing or avoiding join operations in the code. There are two reasons. The first one is to double-calculate each sequence and then join the pipeline, which requires repeated traversal and takes a long time. The second reason is that join operation is too easy to shuffle for big data.

According to the test Data, my running time is 9.5367.

The screenshot shows the PyCharm IDE with a project named '9313'. The code in the main editor defines two functions: `gen_meta_features` and `test_prediction`. `gen_meta_features` iterates over 5 groups, training a pipeline on each and generating meta-features. `test_prediction` uses these meta-features to train a final meta-classifier. The Run window at the bottom shows the execution output, including a warning about the native hadoop library and the final running time of 9.5367431640625e-07.

```

20  return pipeline
21
22 def gen_meta_features(training_df, nb_0, nb_1, nb_2, svm_0, svm_1, svm_2):
23     i = 0
24     while i < 5:
25         condition = training_df['group'] == i
26         train_data = training_df.filter(~condition)
27         test_data = training_df.filter(condition)
28
29         new_pipeline = Pipeline(stages=[nb_0, nb_1, nb_2, svm_0, svm_1, svm_2])
30         if i == 0:
31             DF = new_pipeline.fit(train_data).transform(test_data)
32         else:
33             DF = DF.union(new_pipeline.fit(train_data).transform(test_data))
34         i += 1
35
36     DF = DF.withColumn("joint_pred_0", DF["nb_pred_0"] * 2 + DF["svm_pred_0"] * 1)
37     DF = DF.withColumn("joint_pred_1", DF["nb_pred_1"] * 2 + DF["svm_pred_1"] * 1)
38     DF = DF.withColumn("joint_pred_2", DF["nb_pred_2"] * 2 + DF["svm_pred_2"] * 1)
39
40     return DF
41
42
43 def test_prediction(test_df, base_features_pipeline_model, gen_base_pred_pipeline_model, gen_meta_feature_pipeline_model, meta_classifier):
44     DataFrame_1 = base_features_pipeline_model.transform(test_df)
45     DataFrame_2 = gen_base_pred_pipeline_model.transform(DataFrame_1)
46     final_result = DataFrame_2.withColumn("joint_pred_0", DataFrame_2["nb_pred_0"] * 2 + DataFrame_2["svm_pred_0"] * 1) \
47         .withColumn("joint_pred_1", DataFrame_2["nb_pred_1"] * 2 + DataFrame_2["svm_pred_1"] * 1) \
48         .withColumn("joint_pred_2", DataFrame_2["nb_pred_2"] * 2 + DataFrame_2["svm_pred_2"] * 1)
49
50     df_features = gen_meta_feature_pipeline_model.transform(final_result)
51     prediction = meta_classifier.transform(df_features).select("id", "label", "final_prediction")
52

```

Run: Executed (2) x

```

/Users/yangyuchen/opt/anaconda3/envs/COMP9313/bin/python "/Users/yangyuchen/PycharmProjects/9313/venv/9313 project2/Executed.py"
2020-08-09 17:46:03 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
2020-08-09 17:46:16 WARN BLAS:61 - Failed to load implementation from: com.github.fommil.netlib.NativeSystemBLAS
2020-08-09 17:46:16 WARN BLAS:61 - Failed to load implementation from: com.github.fommil.netlib.NativeRefBLAS
running time: 9.5367431640625e-07
0.7483312619309965
Process finished with exit code 0

```

2. How would you improve the performance (e.g., F1) of the stacking model ?

1) Remove the join mode and use the mode of pipeline

When I wrote the project for the first time, I calculated the number of each column separately and then joined the original data one by one. But, a shuffle is very easy to happen in a join operation. In the final testing stage, I found that this was too time-consuming.

```

25 nb_pred_0 = nb_0.fit(train_data).transform(test_data)
26 nb_pred_1 = nb_1.fit(train_data).transform(test_data)
27 nb_pred_2 = nb_2.fit(train_data).transform(test_data)
28
29 svm_pred_0 = svm_0.fit(train_data).transform(test_data)
30 svm_pred_1 = svm_1.fit(train_data).transform(test_data)
31 svm_pred_2 = svm_2.fit(train_data).transform(test_data)
32
33
34 if count == 1:
35     nb1 = nb_pred_0
36     nb2 = nb_pred_1
37     nb3 = nb_pred_2
38     svm1 = svm_pred_0
39     svm2 = svm_pred_1
40     svm3 = svm_pred_2
41     count = 0
42 else:
43     nb1 = nb1.union(nb_pred_0)
44     nb2 = nb2.union(nb_pred_1)
45     nb3 = nb3.union(nb_pred_2)
46     svm1 = svm1.union(svm_pred_0)
47     svm2 = svm2.union(svm_pred_1)
48     svm3 = svm3.union(svm_pred_2)
49
50 result = training_df.join(nb1[['id', 'nb_pred_0']], on='id')
51 result = result.join(nb2[['id', 'nb_pred_1']], on='id')
52 result = result.join(nb3[['id', 'nb_pred_2']], on='id')
53
54 result = result.join(svm1[['id', 'svm_pred_0']], on='id')
55 result = result.join(svm2[['id', 'svm_pred_1']], on='id')
56 result = result.join(svm3[['id', 'svm_pred_2']], on='id')

```

First-version in project2 (join function)

So I tried using a Pipeline, and in my second code, I created a new pipeline so that all NB and SVM computations could be done in the 4 lines of code instead of joining 6 times.

```

new_pipeline = Pipeline(stages=[nb_0, nb_1, nb_2, svm_0, svm_1, svm_2])
if i == 0:
    DF = new_pipeline.fit(train_data).transform(test_data)
else:
    DF = DF.union(new_pipeline.fit(train_data).transform(test_data))

```

2) More needs to be done in the early data preprocessing

When looking at descript data, I found that there were a lot of punctuation that greatly affected our implementation efficiency. It would be more efficient to get rid of the punctuation early on. tuning hyper para of meta model.

3) Tuning hyper para of meta model

Logistic regression is a statistical model that in its basic form uses a logistic function to model a binary dependent variable, although many more complex extensions exist. I want to improve F1 by adjusting the parameters here.

```

the meta classifier
= LogisticRegression(featuresCol='meta_features', labelCol='label', predictionCol='final_prediction', maxIter=20, regParam=1., elasticNetParam=0)
sifier = lr_model.fit(meta_features)

```

- maxIter

Max number of iterations (≥ 0). (default: 100, current: 20). Logistic regression is optimized by iterative methods like gradient descent. I used the Development Dataset to adjust the data.

When I increase the maximum number of iterations to 100, the running time increases but the accuracy increases to 0.7502.

By reducing the maximum number of iterations to 2, the running time is reduced but the accuracy is reduced to 0.73506.

```

60 |ler pipeline
61 |r(inputCols=['nb_pred_0', 'nb_pred_1', 'nb_pred_2', 'svm_pred_0', 'svm_pred_1', 'svm_pred_2', 'joint_pred_0', 'joint_pred_1', 'joint_pred_2'], outputCols=['
62 |utCols=['vec()'.format(i) for i in range(9)], outputCol='meta_features')
63 |tages=[onehot_encoder, vector_assembler])
64 |meta_feature_pipeline.fit(meta_features)
65 |line_model.transform(meta_features)
66 |
67 |
68 |Col='meta_features', labelCol='label', predictionCol='final_prediction', maxIter=2, regParam=1., elasticNetParam=0)
69 |atures)
70 |
71 |
72 |base_features_pipeline_model, gen_base_pred_pipeline_model, gen_meta_feature_pipeline_model, meta_classifier)
73 |
74 |
75 |
76 |
77 |t_time)
78 |
79 |
80 |

```

Executed (2) x

```

/Users/yangyuchen/opt/anaconda3/envs/COMP9313/bin/python "/Users/yangyuchen/PycharmProjects/9313/venv/9313_project2/Executed.py"
2020-08-09 19:41:42 WARN NativeCodeLoader:62 - Unable to load native-hadoop library for your platform... using builtin-java classes where applicable
Setting default log level to "WARN".
To adjust logging level use sc.setLogLevel(newLevel). For SparkR, use setLogLevel(newLevel).
2020-08-09 19:41:59 WARN BLAS:61 - Failed to load implementation from: com.github.fommil.netlib.NativeSystemBLAS
2020-08-09 19:41:59 WARN BLAS:61 - Failed to load implementation from: com.github.fommil.netlib.NativeRefBLAS
running time: 0.6
0.735067638062769
Process finished with exit code 0

```

maxIter = 2 with 0.73506 prediction

- regParam
regularization parameter (≥ 0). (default: 0.0)
- elasticNetParam

The ElasticNet mixing parameter, in range $[0, 1]$. For $\alpha = 0$, the penalty is an L2 penalty.

For $\alpha = 1$, it is an L1 penalty. (default: 0.0)

For the purpose of hyperparameter tuning we will consider the following parameters:

- maxIter $[10, 100, 1000]$
- regParam $[0.01, 0.5, 2.0]$
- elasticNetParam $[0.0, 0.5, 1.0]$