

# 工作总结

💡 信者行之基，行者人之本。—— 刘昼

汇报人：@管中港

部门：人工智能部

日期：2022年3月24日

## 关于java端调用sklearn模型的总结

### 1.sklearn 模型保存为pmml文件

一般可以通过sklearn2pmml或者nyoka这两个包来将sklearn模型的打包为pmml的格式，在由java进行调用。[pmml](#)模型具有平台无关性，可以实现跨平台部署。

java调用sklearn模型的一般步骤如下：

- 将数据（文本或非文本）转化为表格数据的形式，即col\_name:data
- 使用特定的pipeline包裹每步操作，并在该管道下进行fit，其中两种方式的包裹如下：

## Python

```
1 from sklearn.pipeline import Pipeline
2 from sklearn_pandas import DataFrameMapper
3
4 lr_pmm1 = Pipeline([
5     ("mapper", DataFrameMapper([
6
7     ("query", TfidfVectorizer(analyzer='char', min_df=5, max_features=1024, dtype=np.float32)),
8     ])),
9     ('lr', LogisticRegression())])
10
11 print("lr_pmm1:", lr_pmm1)
12 lr_pmm1.fit(train_df, train_y)
13
14
15
16 from sklearn2pmm1.pipeline import PMMLPipeline
17 from sklearn_pandas import DataFrameMapper
18
19 lr_pmm1_ = PMMLPipeline([
20     ("mapper", DataFrameMapper([
21
22     ("query", TfidfVectorizer(analyzer='char', min_df=5, max_features=1024, dtype=np.float32)),
23     ])),
24     ('lr', LogisticRegression())])
25
26 print("lr_pmm1_:", lr_pmm1_)
27 lr_pmm1_.fit(train_df, train_y)
```

- 可以看到，两种方式的定义格式基本相同只是基础包不一样。其中DataFrameMapper主要用来包裹对特征的处理方式，比如：标准化、归一化等。我在这里是将每条query当作了col\_name为query下的一个表格数据，然后就是对query整列数据进行TfidfVectorizer。需要注意的有以下三点：
  - "mapper"是固定的
  - "query"表示表格数据中的query一列
  - TfidfVectorizer中dtype=np.float32表示半精度，默认为64，当精度再下降时TfidfVectorizer会报错
- 然后就是保存，两种定义的保存方式也不一样

## Python

```
1 from nyoka import skl_to_pmml
2
3 skl_to_pmml(lr_pmml,col_names=
  ['query'],pmml_f_name='classifier_f1_%.3f.pmml'%f1_score)
4
5
6 from sklearn2pmml import sklearn2pmml
7
8 sklearn2pmml(lr_pmml_,'classifier_f1_%.3f.pmml'%f1_score,with_repr = True)
```

- 对于使用skl\_to\_pml的保存没问题，但是使用sklearn2pmml的保存会报以下错误，可以看出是运行时间错误，考虑到进行TfidfVectorizer会使得参数量相当多，所以猜测sklearn2pmml的保存可能有一定的参数上限，相比skl\_to\_pmml则可以容纳更多的参数量

## Python

```
1 RuntimeError: The JPMML-SkLearn conversion application has failed.
2 The Java executable should have printed more information
3 about the failure into its standard output and/or standard error streams
```

## 2.java调用pmml文件

java对pmml文件的调用有以下几步：

- 首先，添加以下依赖

## Java

```
1 "org.jpmmml" % "pmml-evaluator" % "1.4.15",
2 "org.jpmmml" % "pmml-evaluator-extension" % "1.4.15",
3 "com.huaban" % "jieba-analysis" % "1.0.2" //如果需要分词的话
```

- 导入以下类

## Java

```
1 import javax.xml.bind.JAXBException;
2 import java.io.File;
3 import java.io.FileInputStream;
4 import java.io.IOException;
5 import java.io.InputStream;
6 import java.util.*;
7
8 import org.dmg.pmml.FieldName;
9 import org.dmg.pmml.PMML;
10 import org.jpmmml.evaluator.*;
11 import org.xml.sax.SAXException;
```

- 读入pmml文件，这里有一点需要注意
  - 第27行modelEvaluatorFactory.newModelEvaluator(pmml)的输入只有pmml，但是这在1.5.x版本中是会报错的，因为1.5.x版本的输入为modelEvaluatorFactory.newModelEvaluator(pmml,model)，所以安装依赖时需要严格控制pmml-evaluator的版本

## Java

```
1 public static Evaluator loadPmml(String path){
2     PMML pmml=new PMML();
3     InputStream inputStream=null;
4     File file=new File(path);
5     //建立文件实例
6     try{
7         inputStream=new FileInputStream(file);
8     }catch (Exception e){
9         e.printStackTrace();
10    }
11    if(inputStream==null){
12        return null;
13    }
14    //读取文件
15    try{
16        pmml=org.jpmmml.model.PMMLUtil.unmarshal(inputStream);
17    }catch (JAXBException | SAXException e1){
18        e1.printStackTrace();
19    }
20    //关闭文件
21    try{
22        inputStream.close();
23    }catch (IOException e){
24        e.printStackTrace();
25    }
26    ModelEvaluatorFactory
    modelEvaluatorFactory=ModelEvaluatorFactory.newInstance();
27    Evaluator evaluator = modelEvaluatorFactory.newModelEvaluator(pmml);
28    pmml=null;
29    return evaluator;
30 }
```

### · 导入模型参数

## Java

```
1 Map<FieldName, FieldValue> arguments = new LinkedHashMap<>(); //存储模型参数
2 //导入模型参数
3 for (InputField inputField:inputFields){
4     FieldName inputFieldName=inputField.getName(); //参数名
5     FieldValue
    inputValue=inputField.prepare(map.get(inputFieldName.getValue())); //句子输入
6     arguments.put(inputFieldName,inputValue);
7 }
```

- java调用pmml文件进行预测时就是以其中的参数进行  $y = w^t x + b$  形式的计算

Java

```
1 Map<FieldName,?> result=model.evaluate(arguments);
```

- 取出result中分类出的类别

Java

```
1 Map<String, Object> resultMap = new HashMap<>();
2 for(TargetField targetField : targetFields) {
3     FieldName targetFieldName = targetField.getName();
4     // System.out.println("targetFieldName"+targetFieldName);
5     Object targetFieldValue = result.get(targetFieldName);
6     if (targetFieldValue instanceof Computable) {
7         Computable computable = (Computable) targetFieldValue;
8         resultMap.put(targetFieldName.getValue(), computable.getResult());
9     }else {
10         resultMap.put(targetFieldName.getValue(), targetFieldValue);
11     }
12 }
```

- 整个过程还有一个重要问题，在调用pmml模型的时候，通常会报错，找不到[http://www.dmg.org/PMML-4\\_4](http://www.dmg.org/PMML-4_4)，这一般时pmml版本与java的对应没做好。原始的pmml开头是上面代码这样，只用将4.4改为4.3就可以正常运行

HTML

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <PMML xmlns="http://www.dmg.org/PMML-4_4" version="4.4.1">
3     <Header copyright="Copyright (c) 2021 Software AG" description="Default
4         description">
5         <Application name="Nyoka" version="5.2.0"/>
6         <Timestamp>2022-03-22 17:54:52.470393</Timestamp>
7     </Header>
8
9
10 <?xml version="1.0" encoding="UTF-8"?>
11 <PMML xmlns="http://www.dmg.org/PMML-4_3" version="4.3.1">
12     <Header copyright="Copyright (c) 2021 Software AG" description="Default
13         description">
14         <Application name="Nyoka" version="5.2.0"/>
15         <Timestamp>2022-03-22 17:54:52.470393</Timestamp>
16     </Header>
```

### 3.碰到的问题

使用java调用sklearn模型可以成功运行，但是在预测时发现输入中文句子，输出的概率很多都一样，输入英文时可以这场预测。查看pmml文件后，发现如下问题：

#### HTML

```
1 <NumericPredictor name="tfidf@[query](s)" exponent="1"
  coefficient="-1.6606425890343206"/>
2 <NumericPredictor name="tfidf@[query](t)" exponent="1"
  coefficient="-4.6560710174221311"/>
3 <NumericPredictor name="tfidf@[query](u)" exponent="1"
  coefficient="-2.0655156213307730"/>
4 <NumericPredictor name="tfidf@[query](v)" exponent="1"
  coefficient="-2.3017803101764049"/>
5 <NumericPredictor name="tfidf@[query](w)" exponent="1"
  coefficient="-3.6962048914338621"/>
6 <NumericPredictor name="tfidf@[query](x)" exponent="1"
  coefficient="-4.2844472787296892"/>
7 <NumericPredictor name="tfidf@[query](y)" exponent="1"
  coefficient="-1.0379443563439172"/>
8 <NumericPredictor name="tfidf@[query](\xe4\xb8\x81)" exponent="1"
  coefficient="-2.7510269805589518"/>
9 <NumericPredictor name="tfidf@[query](\xe4\xb8\x87)" exponent="1"
  coefficient="-1.8688661028053180"/>
10 <NumericPredictor name="tfidf@[query](\xe4\xb8\x91)" exponent="1"
  coefficient="1.1914348323726289"/>
11 <NumericPredictor name="tfidf@[query](\xe4\xb8\x93)" exponent="1"
  coefficient="-3.3635811569164296"/>
12 <NumericPredictor name="tfidf@[query](\xe4\xb8\x96)" exponent="1"
  coefficient="-2.2035012627585639"/>
13 <NumericPredictor name="tfidf@[query](\xe4\xb8\x9a)" exponent="1"
  coefficient="2.2454507982333558"/>
```

- 英文字母是可以正常表示的，但是中文被表示为16进制的格式

想到的解决方案有：

- 保存pmml是否有参数可以控制中文编码
- 读取pmml的时候，是否有参数可以控制中文编码
- 将java端的输入转换为16进制

但是，考虑到这个问题解决后也无法保证没有其他问题的出现，而且，时间上也不允许，所以就中断这个方案了。

### 4.思考总结

在网上搜集资料时，使用的数据基本都是表格数据，其中的数据类型有数值和英文（用tfidf来处理英文作为其中的一个特征来训练模型），对于这种纯nlp的任务没有找到例子。考虑到，在业务上，对于nlp任务很少使用到sklearn，所以，猜测pmml这种格式不一定支持中文的处理。