Published on InterSystems Developer Community (https://community.intersystems.com)

記事

Toshihiko Minamoto · 2022年2月3日 12m read

Open Exchange

PandasデータフレームをIRISに保存する - 簡易メモ

キーワード: Pandasデータフレーム、IRIS、Python、JDBS

目的

PandasデータフレームはEDA(探索的データ分析)に一般的に使用されるツールです。 MLタスクは通常、データをもう少し理解することから始まります。

先週、私はKaggleにあるこちらのCovid19データセットを試していました。基本的に、このデータは1925件の遭 遇の行と231列で構成されており、タスクは、患者(1つ以上の遭遇レコードにリンク)がICUに入室するかどうか を予測するものです。 つまりこれは、いつものようにpandas.DataFrameを使用して、まず簡単にデータを確認す る、通常の分類タスクです。

現在では、<u>IRIS IntegratedML</u>

が提供されています。これには強力な「AutoML」のオプションに関する洗練されたSQLラッパーがあるため、従 来型のMLアルゴリズムに対抗して、多様なデータフレームのステージをIRISデータベーステーブルに保存してか ら、IntegratedMLを実行する方法を頻繁に採用しています。

ただし、<u>dataframe.tosql()</u>

はまだIRISで機能しないため、実際には、ほとんどの時間を他のデータ保存手段をいじることに充てていました。 いわば、土曜の朝の楽しい朝食の時間に素敵なオムレツを作ろうとしていたのに、一日中コンロの下で、ガスとシ ンクの配管作業をしていたような状況です。

さて、完璧ではありませんが、数週間後に忘れてしまわないように、簡単なメモを残しておくことにします。

範囲

IRISでdataframe.tosql()

を作成しませんでした。残念ながら、まだそこにはたどり着いていませんが、JDBC(JayDeBeApi)を介してデ ータフレームを動的に直接IRISに保存する簡単なPython関数をできるだけ単純かつ生の状態を維持してさくせい しました。 既知の問題(「MemoryError」)により、PyODBCではまだ機能しない可能性があります。

環境

以下のスクリプトのベースとして、<u>単純なdocker-coposeを介したこちらのIntegratedMLテンプレート</u> を使って、テストしています。 <u>環境トポロジー</u>はGitHubリポジトリに含まれます。 IRISコンテナーに接続するには、こちらの<u>JDBC Jupyter ノートブック</u>を使用しています。

テスト

1. dataframe.to<u>s</u>ql() をエミュレートするPython**関数を定義する**

ノートブックのセルで以下を実行しました。

```
def to_sql_iris(cursor, dataFrame, tableName, schemaName='SQLUser', drop_table=False
):
<span style="color:#999999;"> """"
```

```
Published on InterSystems Developer Community (https://community.intersystems.com)
```

```
Dynamically insert dataframe into an IRIS table via SQL by "excutemany"
        Inputs:
            cursor:
                         Python JDBC or PyODBC cursor from a valid and establised DB
connection
            dataFrame:
                         Pandas dataframe
            tablename: IRIS SQL table to be created, inserted or apended
            schemaName: IRIS schemaName, default to "SQLUser"
            drop table: If the table already exsits, drop it and re-
create it if True; othrewise keep it and appen
        Output:
            True is successful; False if there is any exception.
        """</span>
        if drop_table:
            try:
                curs.execute("DROP TABLE %s.%s" %(schemaName, tableName))
            except Exception:
                pass
        try:
            dataFrame.columns = dataFrame.columns.str.replace("[() -]", "_")
            curs.execute(pd.io.sql.get schema(dataFrame, tableName))
        except Exception:
            pass
        curs.fast_executemany = True
        cols = ", ".join([str(i) for i in dataFrame.columns.tolist()])
        wildc =''.join('?, ' * len(dataFrame.columns))
        wildc = '(' + wildc[:-2] + ')'
        sql = "INSERT INTO " + tableName + " ( " + cols.replace('-', '_') + " ) VALUE
S" + wildc
        #print(sql)
        curs.executemany(sql, list(dataFrame.itertuples(index=False, name=None)) )
        return True
```

基本的に、上記はIRISテーブルにデータフレームを動的に挿入しようとしています。 テーブルがすでに存在する 場合は、完全なデータフレームがその最後にアペンドされますが、存在しない場合は、データフレームの次元(列 名と列の型)に基づく新しいテーブルが作成され、その全コンテンツが挿入されます。 executemanyメソッドを使用しているだけです。

2. テスト - 生データファイルをデータフレームに読み込む

ノートブックで以下のコードを実行し、ローカルのドライブからデータフレームに生データを読み込みます。 生データは、<u>こちらのKaggleサイトからダウンロード</u>可能です。

```
import numpy as np
import pandas as pd
from sklearn.impute import SimpleImputer
import matplotlib.pyplot as plt
from sklearn.linear_model import LogisticRegression
from sklearn.model_selection import train_test_split
from sklearn.metrics import classification_report, roc_auc_score, roc_curve
import seaborn as sns
sns.set(style="whitegrid")
```

PandasデータフレームをIRISに保存する - 簡易メモ Published on InterSystems Developer Community (https://community.intersystems.com)

```
import os
for dirname, _, filenames in os.walk('./input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))
```

```
df = pd.read_excel("./input/raw_data_kaggle_covid_icu.xlsx")
```

df

 $./input/datasets \underline{6}05991\underline{1}272346\underline{K}aggle \underline{S}irio \underline{L}ibanes \underline{I}CU\underline{P}rediction.xlsx$

Out[2]:

	PA	AG	AG	GE	DIS	DIS	DIS	DIS	DIS	DIS	 TE	OX	BL	BL	HE	RE	TE	OX	WI	ICU
	TIEN		EP	NDE	EAS	EAS	EAS	EAS	EAS	EAS	MPE	YGE			ART	SPI	MPE	YGE	NDO	
	I <u>V</u> I ⊂IT	BOV		к	EG	EG	EG	EG	EG	EG		N <u>S</u> ^ T I I			KA TE	RAI		N <u>S</u>	vv	
	511_ IDE	E05										ΑΙU ΒλΤΙ	550	55U DF						
	IUE NTIE		┡											к⊑_ еіет						
						62	0.5	64	33		1		STO							
															-	RF		RF		
														FF		Ľ		Ľ		
													RE	REL						
													L							
		1			0.0	0.0	0.0	0.0	1.0	1.0	 -1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	0-2	
			60th								0000	0000	0000	0000	0000	0000	0000	0000		
4		4			0.0	0.0	0.0	0.0	1.0	1.0			0	0	0		0		2.4	
		I	60th		0.0	0.0	0.0	0.0	1.0	1.0	 0000	0000	0000	0000	0000	0000	0000	h000	2-4	
											0	0	0	0	0	0	0	0		
2		1			0.0	0.0	0.0	0.0	1.0	1.0	 -	-	-	-	-	-	-	Ē	4-6	
			60th								NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN		
3		1			0.0	0.0	0.0	0.0	1.0	1.0	 -1.0	-1.0					-1.0	-1.0		
			60th								0000	0000	NaN	NaN	NaN	NaN	0000	0000	6-12	
4		1			0.0	0.0	0.0	0.0	10	10			0.2	0.4	0.0	0.0				1
4		I	60th		0.0	0.0	0.0	0.0	1.0	1.0	 -0.2 2800	1919	-0.3 8006	0.4	-0.2 2016	0.0	1-0.2 1228	1113		I
											5003 5	2	0330 7	8 8	2040 2	4	b 220	3	12	
											 	- 	· 		- 	·			<u> </u>	
	384			1	0.0	0.0	0.0	0.0	0.0	0.0	 -1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	0-2	
1920			50th								0000	0000	0000	0000	0000	0000	0000	0000		
											0	0	0	0	0	0	0	0		
	384			1	0.0	0.0	0.0	0.0	0.0	0.0	 -1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	2-4	
1921			50th								0000	0000	0000	0000	0000	0000	0000	0000		
	384			1	0.0	0.0	0.0	0.0	0.0	0.0	U _1 0	U _1 0	U _1 0	U _1 0	U _1 0	U _1 0	U _1 0	0	4-6	
1922	504		50th	'	0.0	0.0	0.0	0.0	0.0	0.0	 0000	0000	0000	0000	0000	0000	0000	0000	4-0	
											0	0	0	0	0	0	0	0		
	384			1	0.0	0.0	0.0	0.0	0.0	0.0	 -1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0	-1.0		
1923			50th								0000	0000	0000	0000	0000	0000	0000	0000	6-12	
											0	0	0	0	0	0	0	0		
1001	384			1	0.0	0.0	1.0	0.0	0.0	0.0	 -0.5	-0.8	-0.7	-0.5	-0.7	-0.6	-0.5	-0.8	AB	
1924			b0th								4761 h	1838	U186 h	8596	6386 o	1290 h	р133 Н	B202	OVE	
			I I	1	I	1	1		1		M	4	d d	1	Ø	b	1/	ĸ	11Z	

1925 rows \times 231 columns

3. テスト - Python over JDBCでIRIS DBに接続する

```
import jaydebeapi
url = "jdbc:IRIS://irisimlsvr:51773/USER"
driver = 'com.intersystems.jdbc.IRISDriver'
user = "SUPERUSER"
password = "SYS"
jarfile = "./intersystems-jdbc-3.1.0.jar"
conn = jaydebeapi.connect(driver, url, [user, password], jarfile)
```

4. テスト - データフレームをIRISテーブルに保存する

```
iris_schema = 'SQLUser'
iris_table = 'Covid19RawTableFromCSV'
```

curs = conn.cursor()

```
to_sql_iris(curs, df, iris_table, iris_schema, drop_table=True) # save it into a ne
w IRIS table of specified name
#to_sql_iris(curs, df, iris_table) # append dataframe to an exsiting IRIS table
```

Out[4]: True

import pandas as pd
from IPython.display import display

df2 = pd.read_sql("SELECT COUNT(*) from %s.%s" %(iris_schema, iris_table),conn)
display(df2)

	Aggregate <u>1</u>
0	1925

したがって、全データが「Covid19RawTableFromCSV」というIRISテーブルに挿入されました。IRIS管理ポータ ルにログインすると、レコードが含まれるそのテーブルも表示されます。

5. テスト - 簡単なペンチマークを実行する

このデータフレームをたとえば10回挿入して、1つのJDBCセッションでこの基本的なCE dockerに掛かった時間を確認しましょう。

from tqdm import tqdm
import pandas as pd
import time
from IPython.display import display

PandasデータフレームをIRISに保存する - 簡易メモ Published on InterSystems Developer Community (https://community.intersystems.com)

start = time.clock()
for i in tqdm(range(0, 10)):
 to_sql_iris(curs, df, iris_table)

print("Total time elasped: ", time.clock()-start, " for importing total records:")
df2 = pd.read_sql("SELECT COUNT(*) from %s.%s" %(iris_schema, iris_table),conn)
display(df2)

100% ???????? 10/10 [00:14<00:00, 1.42s/it]

Total time elasped: 12.612431999999998 for importing total records:

	Aggregate <u>1</u>
0	19250

以上です。非常に基本的ではありますが、少なくともデータ分析パイプラインに沿って操作されたデータフレーム を保存し、多少本格的なMLを試すために、SQLインターフェースを介してIntegratedMLを呼び出すことができる ようになりました。

言及すべき**警告**:

データフレームの文字列は、「 オブジェクト」として解釈されることがあるため、 df['column'].astype(str)などを使用して、IRISテーブルに挿入される前に文字列に変換する必要があります。 「DROP TABLE」は前のテーブルを上書きするために使用されます。「DROP VIEW」は前のビューを削除するために使用できます。

<u>#JDBC #Python #Machine Learning (ML) #InterSystems IRIS</u> InterSystems Open Exchangeで関連アプリケーションを確認してください

ソースURL:

https://jp.community.intersystems.com/post/pandas%E3%83%87%E3%83%BC%E3%82%BF%E3%83%95%E3%83%AC%E3%83%BC%E3%83%A0%E3%82%92iris%E3%81%AB%E4%BF%9D%E5%AD%98%E3%81%99%E3%82%8B-%E7%B0%A1%E6%98%93%E3%83%A1%E3%83%A2