

```

1 import DataPackageSupport as dps
2
3
4 class DataPackage:
5     __version = 0.1
6
7     def __init__(self,
8                 origData,
9                 uniqueColumn,
10                targetColumn
11                ):
12         self.uniqueColumn = uniqueColumn
13         self.targetColumn = targetColumn
14         self.__setOrigData(origData)
15
16
17     def __setOrigData(self, origData):
18         self.origData = origData
19         self.isOrigDataLoaded = True
20
21         #Get and set features listing, removing unique
22         and target columns
23         self.dataFeatures = list(origData.columns)
24         # Remove unique and target columnm
25         self.dataFeatures.remove(self.uniqueColumn)
26         self.dataFeatures.remove(self.targetColumn)
27
28         # A new dataframe means we need to reset our
29         work
30         self.__resetWork()
31
32         # if the data gets changed then we need to
33         # invalidate all the results/work done previously
34         def __resetWork(self):
35             self.isBalanced = False
36
37             self.__clearTrainTestData()

```

```
37     def __setTrainData(self, trainData):
38         self.trainData = trainData
39         self.isTrainDataLoaded = True
40
41     def getTrainData(self):
42         return self.trainData
43
44     def __setTestData(self, testData):
45         self.testData = testData
46         self.isTestDataLoaded = True
47
48     def getTestData(self):
49         return self.testData
50
51     def __clearOrigData(self):
52         self.origData = None
53         self.isOrigDataLoaded = False
54
55     def __clearTrainTestData(self):
56         self.isTrainTestSplit = False
57         self.isTrainDataLoaded = False
58         self.trainData = None
59
60         self.isTestDataLoaded = False
61         self.testData = None
62
63     def splitTrainTest(self,
64                        stratifyColumn=None,
65                        test_size=0.2,
66                        random_state=765,
67                        shuffle=True
68                        ):
69
70         if stratifyColumn is None:
71             stratifyColumn = self.targetColumn
72
73         train, test = dps.trainTestSplit(dataFrame=self
.getOrigData(),
```

```

74                                     test_size=
    test_size,
75                                     random_state=
    random_state,
76
    stratifyColumn=stratifyColumn,
77                                     shuffle=
    shuffle)
78
79     self.__setTrainData(train)
80     self.__setTestData(test)
81     self.isTrainTestSplit = True
82     self.__clearOrigData()
83
84     def getOrigData(self):
85         if self.isOrigDataLoaded == False:
86             display("Original data frame is not loaded
87 ")
88             return self.origData
89
90     def display(self):
91         emptySpace = '      '
92         indent = emptySpace + '---> '
93
94         print(f'{emptySpace}DataPackage summary:')
95         print(f'{emptySpace}Attributes:')
96         print(f'{indent}uniqueColumn: {self.
uniqueColumn}')
97         print(f'{indent}targetColumn: {self.
targetColumn}')
98
99         print(f'{emptySpace}Process:')
100        print(f'{indent}isBalanced: {self.isBalanced}'
)
101        print(f'{indent}isTrainTestSplit: {self.
isTrainTestSplit}')
102
103        print(f'{emptySpace}Data:')

```

```
103         print(f'{indent}isOrigDataLoaded: {self.  
isOrigDataLoaded}')  
104         print(f'{indent}isTrainDataLoaded: {self.  
isTrainDataLoaded}')  
105         print(f'{indent}isTestDataLoaded: {self.  
isTrainDataLoaded}')  
106  
107     def displayClassBalance(self, columnName=None,  
verbose=False, showRecords=5):  
108         if columnName is None:  
109             columnName = self.targetColumn  
110  
111         dps.displayClassBalance(dataFrame=self.  
getOrigData(),  
112                                 columnName=columnName,  
113                                 showRecords=  
showRecords,  
114                                 verbose=verbose)  
115  
116     def classBalanceUndersample(self,  
117                                 columnName=None):  
118  
119  
120         if columnName is None:  
121             columnName = self.targetColumn  
122  
123         # Needs to be balanced  
124         dfBalanced = dps.classBalanceUndersample(  
dataFrame=self.getOrigData(),  
125         columnName=columnName)  
126  
127         if not self.isBalanced:  
128             self.__setOrigData(dfBalanced)  
129             self.isBalanced = True  
130  
131     def showClassBalance(self,  
132                           columnName=None):
```

```
133         if columnName is None:
134             columnName = self.targetColumn
135
136             # Needs to be balanced
137             dfBalanced = dps.classBalanceUndersample(
138                 dataframe=self.getOrigData(),
139                 columnName=columnName,
140                 alreadyBalanced=True)
141
```