```python
 1  import DataPackage as dp
 2  import DataExperimentSupport as des
 3  import copy
 4  import pickle
 5
 6
 7  class DataExperiment:
 8
 9      def __init__(self,
10                   projectName,
11                   experimentName,
12                   origData,
13                   uniqueColumn,
14                   targetColumn,
15                   classifier,
16                   processDataPackage=False
17                   ):
18          self.projectName = projectName
19          self.experimentName = experimentName
20          self.__setDataPackage(origData=origData,
21                                uniqueColumn=uniqueColumn
    ,
22                                targetColumn=targetColumn
    )
23          self.__setClassifier(classifier)
24
25          # Should really consider putting these into a
    function
26          # Following are default values on init for
    stuff set later
27          self.isBaseModelLoaded = False
28          self.baseModel = None
29          self.isFinalModelLoaded = False
30          self.finalModel = None
31
32          self.isBaseModelPredicted = False
33          self.baseModelPrediction = None
34          self.baseModelAccuracy = None
```

```
35            self.baseModelPrecision = None
36            self.baseModelRecall = None
37            self.baseModelF1 = None
38            self.baseModelCohenKappa = None
39
40            self.isFinalModelPredicted = False
41            self.finalModelPrediction = None
42            self.finalModelPrediction = None
43            self.finalModelAccuracy = None
44            self.finalModelPrecision = None
45            self.finalModelRecall = None
46            self.finalModelF1 = None
47            self.finalModelCohenKappa = None
48
49            self.isBaseModelLearningCurveCreated = False
50            self.baseModel_train_sizes = None
51            self.baseModel_train_scores = None
52            self.baseModel_test_scores = None
53            self.baseModel_fit_times = None
54            self.baseModelROCAUC = None
55
56            self.isFinalModelLearningCurveCreated = False
57            self.finalModel_train_sizes = None
58            self.finalModel_train_scores = None
59            self.finalModel_test_scores = None
60            self.finalModel_fit_times = None
61
62            self.isBaseModelROCAUCCalculated = False
63            self.baseModelROCAUC = None
64            self.isFinalModelROCAUCCalculated = False
65            self.finalModelROCAUC = None
66
67            self.finalFeaturesAll = None
68            self.finalFeatures = None
69
   # ==================================================
70
71        if processDataPackage:
```

```python
72              self.processDataPackage()

73

74          self.display()

75

76      def display(self):
77          indent = '---> '
78          print(f'DataExperiment summary:')
79          print(f'{indent}projectName: {self.projectName
    }')
80          print(f'{indent}experimentName: {self.
    experimentName}')
81          print(f'{indent}isDataPackageLoaded: {self.
    isDataPackageLoaded}')

82

83          print(f'{indent}isBaseModelLoaded: {self.
    isBaseModelLoaded}')
84          print(f'{indent}isBaseModelPredicted: {self.
    isBaseModelPredicted}')
85          print(f'{indent}
    isBaseModelLearningCurveCreated: {self.
    isBaseModelLearningCurveCreated}')

86

87          print(f'{indent}isFinalModelLoaded: {self.
    isFinalModelLoaded}')
88          print(f'{indent}isFinalModelPredicted: {self.
    isFinalModelPredicted}')
89          print(f'{indent}
    isFinalModelLearningCurveCreated: {self.
    isFinalModelLearningCurveCreated}')

90

91          print(f'{indent}isClassifierLoaded: {self.
    isClassifierLoaded}')
92          print(self.getClassifier())
93          print('')
94          self.dataPackage.display()

95

96      def getClassifier(self):
97          return copy.deepcopy(self.classifier)
```

```python
 98
 99    def __setClassifier(self, classifier):
100        self.classifier = classifier
101        self.isClassifierLoaded = True
102
103    def __setDataPackage(self,
104                         origData,
105                         uniqueColumn,
106                         targetColumn):
107
108        self.dataPackage = dp.DataPackage(origData=
   origData,
109                                          uniqueColumn
   =uniqueColumn,
110                                          targetColumn
   =targetColumn)
111        self.isDataPackageLoaded = True
112
113    def processDataPackage(self, verbose=False):
114        self.dataPackage.classBalanceUndersample()
115        self.dataPackage.splitTrainTest()
116
117        if verbose:
118            self.dataPackage.display()
119
120    def createBaseModel(self):
121        model = des.createModel(data=self.dataPackage.
   getTrainData(),
122                                uniqueColumn=self.
   dataPackage.uniqueColumn,
123                                targetColumn=self.
   dataPackage.targetColumn,
124                                classifier=self.
   getClassifier())
125
126        self.__setBaseModel(model)
127
128    def createFinalModel(self,
```

```
129                                    featureImportanceThreshold=0.
     002):
130          # TO DO:
131          # need to get list of features we wish to use
     to subset full data
132
133          impFeatureListFull = self.
     __getFinalModelFeatures(returnAbove=
     featureImportanceThreshold,
134
        includeUniqueAndTarget=True)
135
136          impFeatureList = self.__getFinalModelFeatures(
     returnAbove=featureImportanceThreshold,

137
     includeUniqueAndTarget=False)
138
139          # get full training dataframe
140          df = self.dataPackage.getTrainData()
141
142          model = des.createModel(data=df[
     impFeatureListFull],
143                                     uniqueColumn=self.
     dataPackage.uniqueColumn,
144                                     targetColumn=self.
     dataPackage.targetColumn,
145                                     classifier=self.
     getClassifier())
146          self.__setFinalModel(model=model,
147                               finalFeatures=
     impFeatureList,
148                               finalFeaturesAll=
     impFeatureListFull)
149
150      def __setBaseModel(self, model):
151          self.baseModel = model
152          self.isBaseModelLoaded = True
153
```

```
154              # when you set base model invalidate some
    things
155          self.isBaseModelPredicted = False
156          self.baseModelPrediction = None
157
158      def __setFinalModel(self,
159                            model,
160                            finalFeatures,
161                            finalFeaturesAll):
162          self.finalModel = model
163          self.isFinalModelLoaded = True
164          self.finalFeaturesAll = finalFeaturesAll
165          self.finalFeatures = finalFeatures
166
167          # when you set base model invalidate some
    things
168          self.isFinalModelPredicted = False
169          self.finalModelPrediction = None
170
171      def getBaseModel(self):
172          return self.baseModel
173
174      def getFinalModel(self):
175          return self.finalModel
176
177      def __setBaseModelPrediction(self,
178                                      predictionData,
179                                      colActual,
180                                      colPredict,
181                                      average='weighted',
182                                      sigDigs=2):
183          self.baseModelPrediction = predictionData
184          self.isBaseModelPredicted = True
185          self.baseModelPredictionColActual = colActual
186          self.baseModelPredictionColPredict =
    colPredict
187
188          self.baseModelAccuracy = round(des.
```

```
188  getModelAccuracy(data=predictionData,
189
         colActual=colActual,
190
         colPredict=colPredict), sigDigs)
191
192      self.baseModelPrecision = round(des.
    getModelPrecision(data=predictionData,
193
         colActual=colActual,
194
         colPredict=colPredict,
195
         average=average), sigDigs)
196
197      self.baseModelRecall = round(des.
    getModelRecall(data=predictionData,
198
    colActual=colActual,
199
    colPredict=colPredict,
200
    average=average), sigDigs)
201
202      self.baseModelF1 = round(des.getModelF1(data=
    predictionData,
203
    colActual=colActual,
204
    colPredict=colPredict,
205
    average=average), sigDigs)
206
207      self.baseModelCohenKappa = round(des.
    getModelCohenKappa(data=predictionData,
208
            colActual=colActual,
209
```

```
209                    colPredict=colPredict), sigDigs)
210
211    def __setFinalModelPrediction(self,
212                                  predictionData,
213                                  colActual,
214                                  colPredict,
215                                  average='weighted',
216                                  sigDigs=2):
217        self.finalModelPrediction = predictionData
218        self.isFinalModelPredicted = True
219        self.finalModelPredictionColActual = colActual
220        self.finalModelPredictionColPredict =
    colPredict
221
222        self.finalModelAccuracy = round(des.
    getModelAccuracy(data=predictionData,
223
        colActual=colActual,
224
        colPredict=colPredict), sigDigs)
225
226        self.finalModelPrecision = round(des.
    getModelPrecision(data=predictionData,
227
          colActual=colActual,
228
          colPredict=colPredict,
229
          average=average), sigDigs)
230
231        self.finalModelRecall = round(des.
    getModelRecall(data=predictionData,
232
      colActual=colActual,
233
      colPredict=colPredict,
234
      average=average), sigDigs)
```

```
235
236            self.finalModelF1 = round(des.getModelF1(data=
       predictionData,

237
       colActual=colActual,

238
       colPredict=colPredict,

239
       average=average), sigDigs)

240
241            self.finalModelCohenKappa = round(des.
       getModelCohenKappa(data=predictionData,

242
                colActual=colActual,

243
                colPredict=colPredict), sigDigs)

244
245      def showBaseModelStats(self):
246          print(f'Base Model Stats:')
247          print(f'Accuracy: {self.baseModelAccuracy}')
248          print(f'Precision: {self.baseModelPrecision}')
249          print(f'Recalll: {self.baseModelRecall}')
250          print(f'F1 Score: {self.baseModelF1}')
251          print(f'Cohen kappa:: {self.
       baseModelCohenKappa}')

252
253      def showFinalModelStats(self):
254          print(f'Final Model Stats:')
255          print(f'Accuracy: {self.finalModelAccuracy}')
256          print(f'Precision: {self.finalModelPrecision}'
       )
257          print(f'Recalll: {self.finalModelRecall}')
258          print(f'F1 Score: {self.finalModelF1}')
259          print(f'Cohen kappa:: {self.
       finalModelCohenKappa}')

260
261      def getBaseModelPrediction(self):
262          if self.isBaseModelPredicted:
```

```python
263              return self.baseModelPrediction
264          else:
265              print(f'No base model predictions
     calculated.')
266              return None
267
268      def getFinalModelPrediction(self):
269          if self.isFinalModelPredicted:
270              return self.finalModelPrediction
271          else:
272              print(f'No final model predictions
     calculated.')
273              return None
274
275      def predictBaseModel(self, average='weighted'):
276          if self.isBaseModelPredicted:
277              display("Base model already predicted.
     Displaying results:")
278              self.showBaseModelStats()
279              return
280
281          tDf, colActual, colPredict = des.predictModel(
     model=self.getBaseModel(),
282
     data=self.dataPackage.getTestData(),
283
     uniqueColumn=self.dataPackage.uniqueColumn,
284
     targetColumn=self.dataPackage.targetColumn)
285
286          self.__setBaseModelPrediction(predictionData=
     tDf,
287                                        colActual=
     colActual,
288                                        colPredict=
     colPredict,
289                                        average=average)
290
```

```python
291                self.showBaseModelStats()
292
293        def predictFinalModel(self, average='weighted'):
294            if self.isFinalModelPredicted:
295                display("Final model already predicted.
    Displaying results:")
296                self.showFinalModelStats()
297                return
298
299            testData = self.dataPackage.getTestData()
300            testData = testData[self.finalFeaturesAll].
    copy()
301
302            tDf, colActual, colPredict = des.predictModel(
    model=self.getFinalModel(),
303
    data=testData,
304
    uniqueColumn=self.dataPackage.uniqueColumn,
305
    targetColumn=self.dataPackage.targetColumn)
306
307            self.__setFinalModelPrediction(predictionData=
    tDf,
308                                           colActual=
    colActual,
309                                           colPredict=
    colPredict,
310                                           average=average
    )
311            self.showFinalModelStats()
312
313        def analyzeBaseModelFeatureImportance(self,
314                                               returnAbove=
    0.002,
315                                               startValue=0
    .0001,
316                                               increment=0.
```

```
316 0001,
317                                       upperValue=0
    .01,
318                                       showSummary=
    True,
319                                       showPlot=
    True):
320
321         df, featureLabel, valueLabel = des.
    getModelFeatureImportance(self.getBaseModel())
322
323         retDf = des.analyzeModelFeatureImportance(data
    =df,
324
    valueLabel=valueLabel,
325
    startValue=startValue,
326
    increment=increment,
327
    upperValue=upperValue,
328
    returnAbove=returnAbove,
329
    showSummary=showSummary,
330
    showPlot=showPlot)
331         return retDf
332
333    def analyzeFinalModelFeatureImportance(self,
334                                       returnAbove
    =0.002,
335                                       startValue=
    0.0001,
336                                       increment=0
    .0001,
337                                       upperValue=
    0.01):
```

```python
338
339          df, featureLabel, valueLabel = des.
    getModelFeatureImportance(self.getFinalModel())
340
341          retDf = des.analyzeModelFeatureImportance(data
    =df,
342
    valueLabel=valueLabel,
343
    startValue=startValue,
344
    increment=increment,
345
    upperValue=upperValue,
346
    returnAbove=returnAbove,
347
    showSummary=True)
348          return retDf
349
350      def showBaseModelFeatureImportance(self,
351                                          startValue=0.
    0001,
352                                          increment=0.
    0001,
353                                          upperValue=0.01
    ):
354
355          df, featureLabel, valueLabel = des.
    getModelFeatureImportance(self.getBaseModel())
356
357          des.analyzeModelFeatureImportance(data=df,
358                                            startValue=
    startValue,
359                                            increment=
    increment,
360                                            upperValue=
    upperValue,
```

```
361                                     showSummary=
    False)

362
363         des.showAllModelFeatureImportance(data=df,
364                                     featureLabel
    =featureLabel,
365                                     valueLabel=
    valueLabel
366                                     )
367
368     def showFinalModelFeatureImportance(self,
369                                     startValue=0.
    0001,
370                                     increment=0.
    0001,
371                                     upperValue=0.
    01):
372
373         df, featureLabel, valueLabel = des.
    getModelFeatureImportance(self.getFinalModel())
374
375         des.analyzeModelFeatureImportance(data=df,
376                                     startValue=
    startValue,
377                                     increment=
    increment,
378                                     upperValue=
    upperValue,
379                                     showSummary=
    False)
380
381         des.showAllModelFeatureImportance(data=df,
382                                     featureLabel
    =featureLabel,
383                                     valueLabel=
    valueLabel
384                                     )
385
```

```
386        def showBaseModelReport(self,
387                                axisLabels
388                                ):
389            self.showBaseModelStats()
390
391            des.showReport(data=self.
   getBaseModelPrediction(),
392                           colNameActual=self.
   baseModelPredictionColActual,
393                           colNamePredict=self.
   baseModelPredictionColPredict,
394                           axisLabels=axisLabels,
395                           titleSuffix=self.experimentName
   )
396
397            self.showBaseModelLearningCurve()
398            self.showBaseModelROCAUC(axisLabels=axisLabels
   )
399
400        def showBaseModelROCAUC(self, axisLabels,
   useStored=False):
401            if useStored and self.
   isBaseModelROCAUCCalculated:
402                print('Base model ROCAUC already
   calculated. Displaying stored results')
403                tViz = self.__getBaseModelROCAUC()
404                tViz.show()
405            else:
406                print('Base model ROCAUC not calculated.
   Starting now')
407                viz = des.showROCAUC(dataTrain=self.
   dataPackage.getTrainData(),
408                                     dataTest=self.
   dataPackage.getTestData(),
409                                     classifier=self.
   getClassifier(),
410                                     axisLabels=axisLabels
   ,
```

```python
411                                                 colNameActual=self.
    dataPackage.targetColumn,
412                                                 features=self.
    getBaseFeatures())
413             self.__setBaseModelROCAUC(visualizer=viz)
414             viz.show()
415
416     def __setBaseModelROCAUC(self,
417                             visualizer):
418         self.isBaseModelROCAUCCalculated = True
419         self.baseModelROCAUC = pickle.dumps(visualizer
    )
420
421     def __getBaseModelROCAUC(self):
422         return pickle.loads(self.baseModelROCAUC)
423
424     def showFinalModelROCAUC(self, axisLabels,
    useStored=False):
425         if useStored and self.
    isFinalModelROCAUCCalculated:
426             print('Final model ROCAUC already
    calculated. Displaying stored results')
427             tViz = self.__getBaseModelROCAUC()
428             tViz.show()
429         else:
430             print('Final model ROCAUC not calculated.
    Starting now')
431             viz = des.showROCAUC(dataTrain=self.
    dataPackage.getTrainData(),
432                                  dataTest=self.
    dataPackage.getTestData(),
433                                  classifier=self.
    getClassifier(),
434                                  axisLabels=axisLabels
    ,
435                                  colNameActual=self.
    dataPackage.targetColumn,
436                                  features=self.
```

```
436 getFinalFeatures())
437             self.__setFinalModelROCAUC(visualizer=viz)
438             viz.show()
439
440     def __getFinalModelROCAUC(self):
441         return pickle.loads(self.finalModelROCAUC)
442
443     def __setFinalModelROCAUC(self,
444                               visualizer):
445         self.isFinalModelROCAUCCalculated = True
446         self.finalModelROCAUC = pickle.dumps(
    visualizer)
447
448     def showFinalModelReport(self,
449                              axisLabels
450                              ):
451         self.showFinalModelStats()
452
453         des.showReport(data=self.
    getFinalModelPrediction(),
454                        colNameActual=self.
    finalModelPredictionColActual,
455                        colNamePredict=self.
    finalModelPredictionColPredict,
456                        axisLabels=axisLabels,
457                        titleSuffix=self.experimentName
    )
458
459         self.showFinalModelLearningCurve()
460         self.showFinalModelROCAUC(axisLabels=
    axisLabels)
461
462     def getBaseFeatures(self):
463         return self.dataPackage.dataFeatures
464
465     def getFinalFeatures(self):
466         return self.finalFeatures
467
```

```
468          def createBaseModelLearningCurve(self,
469                                           cv=None,
470                                           n_jobs=None,
471                                           train_sizes=None,
472                                           verbose=4):
473          # If it is already predicted just show it
474          if self.isBaseModelLearningCurveCreated:
475              print('Base model learning curve already
     calculated. Displaying results:')
476              self.showBaseModelLearningCurve()
477          else:
478              df = self.dataPackage.getTrainData()
479              train_sizes, train_scores, test_scores,
     fit_times = des.create_learning_curve(
480                  estimator=self.getClassifier(),
481                  X=df[self.dataPackage.dataFeatures],
482                  y=df[self.dataPackage.targetColumn],
483                  cv=cv,
484                  n_jobs=n_jobs,
485                  train_sizes=train_sizes,
486                  verbose=verbose)
487
488              self.__setBaseModelLearningData(
     train_sizes=train_sizes,
489
     train_scores=train_scores,
490
     test_scores=test_scores,
491                                              fit_times=
     fit_times)
492
493      def createFinalModelLearningCurve(self,
494                                        cv=None,
495                                        n_jobs=None,
496                                        train_sizes=None
     ,
497                                        verbose=4):
498          # If it is already predicted just show it
```

```python
499                 if self.isFinalModelLearningCurveCreated:
500                     print('Final model learning curve already
    calculated. Displaying results:')
501                     self.showFinalModelLearningCurve()
502             else:
503
504                 df = self.dataPackage.getTrainData()
505                 train_sizes, train_scores, test_scores,
    fit_times = des.create_learning_curve(
506                     estimator=self.getClassifier(),
507                     X=df[self.finalFeatures],
508                     y=df[self.dataPackage.targetColumn],
509                     cv=cv,
510                     n_jobs=n_jobs,
511                     train_sizes=train_sizes,
512                     verbose=verbose)
513                 self.__setFinalModelLearningData(
    train_sizes=train_sizes,
514
    train_scores=train_scores,
515
    test_scores=test_scores,
516                                                 fit_times
    =fit_times)
517
518     def __setBaseModelLearningData(self,
519                                     train_sizes,
520                                     train_scores,
521                                     test_scores,
522                                     fit_times):
523         self.isBaseModelLearningCurveCreated = True
524         self.baseModel_train_sizes = train_sizes
525         self.baseModel_train_scores = train_scores
526         self.baseModel_test_scores = test_scores
527         self.baseModel_fit_times = fit_times
528
529     def __setFinalModelLearningData(self,
530                                     train_sizes,
```

```
531                                            train_scores,
532                                            test_scores,
533                                            fit_times):
534          self.isFinalModelLearningCurveCreated = True
535          self.finalModel_train_sizes = train_sizes
536          self.finalModel_train_scores = train_scores
537          self.finalModel_test_scores = test_scores
538          self.finalModel_fit_times = fit_times
539
540      def showBaseModelLearningCurve(self,
541                                     axes=None,
542                                     ylim=(0.0, 1.01)
543                                     ):
544          if self.isBaseModelLearningCurveCreated:
545
546              des.plot_learning_curve(train_sizes=self.
    baseModel_train_sizes,
547                                      train_scores=self.
    baseModel_train_scores,
548                                      test_scores=self.
    baseModel_test_scores,
549                                      fit_times=self.
    baseModel_fit_times,
550                                      title=self.
    experimentName,
551                                      axes=axes,
552                                      ylim=ylim
553                                      )
554          else:
555              display('Base model Learning curve has not
     yet been calculated')
556
557      def showFinalModelLearningCurve(self,
558                                      axes=None,
559                                      ylim=(0.0, 1.01)
560                                      ):
561          if self.isFinalModelLearningCurveCreated:
562              des.plot_learning_curve(train_sizes=self.
```

```
562 finalModel_train_sizes,
563                                        train_scores=self.
    finalModel_train_scores,
564                                        test_scores=self.
    finalModel_test_scores,
565                                        fit_times=self.
    finalModel_fit_times,
566                                        title=self.
    experimentName,
567                                        axes=axes,
568                                        ylim=ylim
569                                        )
570         else:
571             display('Final model Learning curve has
    not yet been calculated')
572
573     def __getFinalModelFeatures(self,
574                                 returnAbove=0.002,
575                                 includeUniqueAndTarget
    =False):
576         # get a list of the features that have been
    deemed important
577         # Get full list of features
578         features = self.dataPackage.dataFeatures
579
580         df, featureLabel, valueLabel = des.
    getModelFeatureImportance(self.getBaseModel())
581
582         retDf = des.analyzeModelFeatureImportance(data
    =df,
583
    valueLabel=valueLabel,
584
    returnAbove=returnAbove,
585
    showSummary=False,
586
    showPlot=False)
```

```
587
588            keepFeatures = retDf[featureLabel].to_list()
589
590            # Initialize important features list
591            features_important = []
592
593            for x in keepFeatures:
594                features_important.append(features[x])
595
596            if includeUniqueAndTarget:
597                # Feature list doesn't include target and
     unique
598                features_important.append(self.dataPackage
     .uniqueColumn)
599                features_important.append(self.dataPackage
     .targetColumn)
600
601            return features_important
602
```