

Programming Assignment #1

1. Use adult.data as your training set and adult.test as your test set.
We used the following data for this project:
Training dataset: adult.data
Test dataset: adult.test
Both files are available from <http://archive.ics.uci.edu/ml/datasets/Adult>
2. You may want to consider creating a validation set from the training set to help you select hyper-parameters. Hyper-parameters you may want to tune are:
 - a) the threshold value for Logistic Regression/Naïve Bayes;
 - b) the depth of the decision tree;
 - c) the value of k in Nearest Neighbors.

The following steps describe the tuning and training process we used for each algorithm:

1. Training dataset (adult.data) was divided into train' and test' partitions (80/20% split)
2. Train' partition was processed using kfold (k=10) to tune hyper parameter for algorithm.
3. Each parameter value was tested for 10 kfold configurations using the corresponding validation dataset to compute score; a mean score was obtained for each parameter value.
4. The best parameter value was thus selected and used against the test' dataset for final reporting algorithm score against the training dataset (please see program output in Appendix A and also in Canvas file "predict_adult.out").
5. Trained function against the training dataset (adult.data).

The following table lists hyper parameters tuned for each algorithm:

Algorithm	Parameter
Naïve Bayes	None known in sklearn
Decision Tree	Tree Depth
KNN	K neighbors
Logistic Regression	Penalty (l1, l2) and C

3. Take a look at the data and get a feel for the types of feature values you observe.
We include in Appendix B a series of plots generated from Weka from a load of test.data in csv format into Weka. Please refer to Appendix B.
4. Preprocess the data. Make sure to explain how you preprocessed the data and why.
Data was preprocessed after load with a general pre-processing routine. After general pre-processing, the data was further pre-processed with specific actions prior to submission to each of the learning algorithms. General preprocessing consisted of:
 - Removed rows with missing values in columns "workclass", "occupation" and "native_country"
 - Dropped column "education" since it's redundant (it matches column "education-num")
 - Converted all categorical columns to numbers
 - Converted all columns to floatsPre-processing specific to each algorithm is described under question b)

- a) It is important to note that some of the features are missing values (indicated with a "?"). How will you handle missing feature values? Explain your solution and why you chose it. You can treat a missing feature value as its own feature value, but this might not be the best solution, so you should

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

occupation and native_country. These are all categorical features. Even though we could **have** used mode as a measure of centrality to impute missing values we opted to drop rows containing missing data to avoid adding bias to the training dataset.

- b) Some of the feature values are continuous and some are categorical. Do you want to discretize the continuous features? If so, why? Do you want to discretize the feature values for all of the classifiers? It is okay to preprocess the data differently for each of the classifiers, just make sure to explain what you did and why.

The following table summarizes additional pre-processing for each specific algorithm:

Algorithm	Pre-Processing (What)	Reasoning (Why)
Naïve Bayes	Discretized continuous features: age, fnlwgt, capital_gain, capital_loss, hours_per_week Replaced original feature values with upper bound of percentile for the feature range.	Initially thought discretization of continuous variables was required since $P(X=x Y=y) = 0$ if x is a number in the Real set. Found that not to be the case as sklearn performs better for Naïve Bayes with continuous feature values.
Decision Tree	Tried to discretize continuous features using the same pre-processing function initially used for Naïve Bays. In the end used original continuous values since sklearn Decision Tree Classifier algorithm performed better that way.	Discretized continuous features under the assumption that calculation of entropy would require discrete values. However, sklearn algorithm performed better with continuous feature values.
KNN	Same as Naïve Bayes. Also tried using continuous features, which yielded poorer results.	Discretized continuous features to increase feature grouping hoping to improve clustering of data points.
Logistic Regression	Standardized continuous features and converted categorical features into dummy variables.	Since logistic regression learns a liner model we must prepare features accordingly. Standardization will give equal weight to features like age (0-100 range) and capital_gain (0-99999). Using dummy variables to represent categorical features is the recommended practice in regression models.

5. If you are using Weka, then convert the CSV into the appropriate ARFF format. If you are using Python, then I recommend looking into Numpy's loadtxt() function.

I only used Weka to plot features in the train data set (see Appendix B).

6. If you are using Python/SKlearn, then implement the learning algorithm.

Learning algorithm was implemented with Python and sklearn.

Code is include in Appendix C and also uploaded with PA#1 assignment report in Canvas (file predict_adult.py)

Please refer to code in Appendix C, for example, function “learn_logit”

8. Evaluate the learned model on the test set.

- a) If you used a validation set to tune your hyper-parameters, then make sure to use the model with the best performance on the test set.

Each model is tuned with 10 kfold validation and tested on the train dataset validation partition (20% of training data) as explained under the answer for Question 2. Tuning parameters are also described under Question 2 above. Program output is included in Appendix A.

Note that instead of running only the best model on the test dataset (file “adult.test”) we ran all models against it, always using the best hyper parameter setting identified during training and tuning. This was done for illustrative purposes.

As explained in the answer to Question 2, prior to running the learned model against the test dataset we learn the model against the entire training set with the proper hyper parameter setting.

Training results are listed in the program output in Appendix A under the title “TRAIN VALIDATION PARTITION”.

Training results are listed in the program output in Appendix A under the title “TEST DATASET”.

9. Report the performance of the model.

- a) Make sure to use an appropriate performance measure (ie. accuracy vs. ROC vs. F1-score). Explain why you chose that performance measure.

Since the proportion of adults with salary >\$50k vs. salary <=\$50 is slightly unbalanced (approximately 29% proportion with salary >\$50k) we chose to use F1-score rather than accuracy to evaluate performance. We tuned all algorithms to capture the hyper parameter yielding the best F1-score.

- b) Explain the effect different hyper-parameters have on the performance of the algorithm. For example, what effect does the depth of the tree have on the performance of the decision tree?

For Naïve Bayes

Algorithm	Parameter	Effect
Naïve Bayes	None known in sklearn	N/A. Could not find tuning parameter for Naïve Bayes in sklearn.
Decision Tree	Tree Depth	Higher depth will increase overfitting; lower depth will increase bias.
KNN	K neighbors	Lower K will increase overfitting; higher K will increase bias.
Logistic Regression	Penalty (l1, l2) and C	Both parameters impact regularization of the regression function and reduce overfitting.

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

- Briefly compare the performance of the different learning algorithms on the test set and make some hypothesizes about why you might have observed these differences.

The table below summarizes our results. We attribute performance differences among the classifiers to undetermined characteristics of the dataset features and their relationship to each other and the target label.

Best performance in each metric in the training set is highlighted in light blue.

Best performance in each metric in the test set is highlighted in light green.

Bold font indicates an improvement or no performance change for the metric on the test dataset (vs. train dataset).

Metric	Naïve Bayes		Decision Tree (depth=9)		KNN (K=13)		Logistic Regression (l1, C=1)	
	Train	Test	Train	Test	Train	Test	Train	Test
Accuracy	0.788	0.796	0.854	0.842	0.811	0.788	0.848	0.848
Error	0.212	0.204	0.146	0.158	0.189	0.212	0.152	0.152
Precision	0.603	0.591	0.675	0.623	0.612	0.500	0.667	0.661
Recall	0.648	0.600	0.607	0.533	0.599	0.432	0.614	0.604
F1-Score	0.564	0.582	0.760	0.750	0.625	0.595	0.730	0.731
ROC AUC	0.741	0.730	0.772	0.738	0.740	0.668	0.770	0.766

Logistic regression clearly offered the best performance on the test dataset. Decision tree offered the best performance on the training dataset. Most algorithms show overfitting with the exception of Naïve Bayes (accuracy and F1-Score improved) and, to a certain extent Logistic Regression (accuracy unchanged and F1-Score improved). The decision tree performance decreased noticeably on the test dataset demonstrating the overfitting tendency of this algorithm. Interestingly, we rerun our program forcing the maximum depth of the decision tree to 8 and got better results in the train test partition and the test dataset. The results are replicated below. They seem to suggest that it is worth decreasing the max tree depth hyper parameter further after tuning. We attempted a similar strategy for KNN, where we increased K in an attempt to reduce overfitting after tuning, but the results were disappointing.

Performance scores highlighted in dark green indicate scores superior to Logistic Regression scores. In a production setting we would select between Logistic Regression and Decision Tree (max depth=8) based on cost analysis of specific misclassification scenarios (for example, FP vs. FN cost).

Metric	Decision Tree (depth=9)		Decision Tree (depth=8)	
	Train	Test	Train	Test
Accuracy	0.854	0.842	0.848	0.851
Error	0.146	0.158	0.152	0.149
Precision	0.675	0.623	0.635	0.641
Recall	0.607	0.533	0.533	0.542
F1-Score	0.760	0.750	0.786	0.784
ROC AUC	0.772	0.738	0.742	0.747

On a final note, we point out that our results are comparable to the results reported on the web site where the datasets were published (<https://archive.ics.uci.edu/ml/machine-learning-databases/adult/adult.names>). The best result reported on the site shows accuracy=0.8595 for “FSS Naïve Bayes”. That is a mere 1% improvement over our best accuracy result with the Decision Tree of max depth 8.

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

References:

- [1] McKinney, Wes. *Python for data analysis: Data wrangling with Pandas, NumPy, and IPython*. " O'Reilly Media, Inc.", 2012.
- [2] James, Gareth, et al. *An introduction to statistical learning*. Vol. 6. New York: springer, 2013.
- [3] Raschka, Sebastian. *Python machine learning*. Packt Publishing Ltd, 2015.
- [4] Friedman, Jerome, Trevor Hastie, and Robert Tibshirani. *The elements of statistical learning*. Vol. 1. Springer, Berlin: Springer series in statistics, 2001.
- [5] Downey, Allen. *Think Python*. " O'Reilly Media, Inc.", 2012.

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017
Carlos Sathler (cssathler@gmail.com)

APPENDIX A

RESULTS FOR NAIVE BAYES CLASSIFIER

TRAIN VALIDATION PARTITION

Accuracy.....: 0.788
Precision Score.: 0.603
Recall Score....: 0.648
F1 Score.....: 0.564
ROC AUC.....: 0.741

Confusion Matrix:

		Prediction	
		0	1
Truth	0	3780	752
	1	528	973

TEST DATASET

Accuracy.....: 0.796
Precision Score.: 0.591
Recall Score....: 0.600
F1 Score.....: 0.582
ROC AUC.....: 0.730

Confusion Matrix:

		Prediction	
		0	1
Truth	0	9767	1593
	1	1479	2221

RESULTS FOR DECISION TREE CLASSIFIER

Best tree depth = 9

TRAIN VALIDATION PARTITION

Accuracy.....: 0.855
Precision Score.: 0.676
Recall Score....: 0.608
F1 Score.....: 0.760
ROC AUC.....: 0.772

Confusion Matrix:

		Prediction	
		0	1
	0	1211	288
	1		

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

+---+-----+

TEST DATASET

Accuracy.....: 0.841
Precision Score.: 0.624
Recall Score....: 0.535
F1 Score.....: 0.748
ROC AUC.....: 0.738

Confusion Matrix:

		Prediction	
		0	1
Truth	0	10692	668
	1	1720	1980

RESULTS FOR KNN CLASSIFIER

Best neighbors value = 13

TRAIN VALIDATION PARTITION

Accuracy.....: 0.811
Precision Score.: 0.612
Recall Score....: 0.599
F1 Score.....: 0.625
ROC AUC.....: 0.740

Confusion Matrix:

		Prediction	
		0	1
Truth	0	3993	539
	1	602	899

TEST DATASET

Accuracy.....: 0.788
Precision Score.: 0.500
Recall Score....: 0.432
F1 Score.....: 0.595
ROC AUC.....: 0.668

Confusion Matrix:

		Prediction	
		0	1
Truth	0	10273	1087
	1	2103	1597

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

RESULTS FOR LOGISTIC REGRESSION CLASSIFIER

Best penalty= 11
Best C value= 1.0

TRAIN VALIDATION PARTITION

Accuracy.....: 0.848
Precision Score.: 0.667
Recall Score.....: 0.614
F1 Score.....: 0.730
ROC AUC.....: 0.770

Confusion Matrix:

		Prediction	
		0	1
Truth	0	4191	341
	1	579	922

TEST DATASET

Accuracy.....: 0.848
Precision Score.: 0.661
Recall Score.....: 0.604
F1 Score.....: 0.731
ROC AUC.....: 0.766

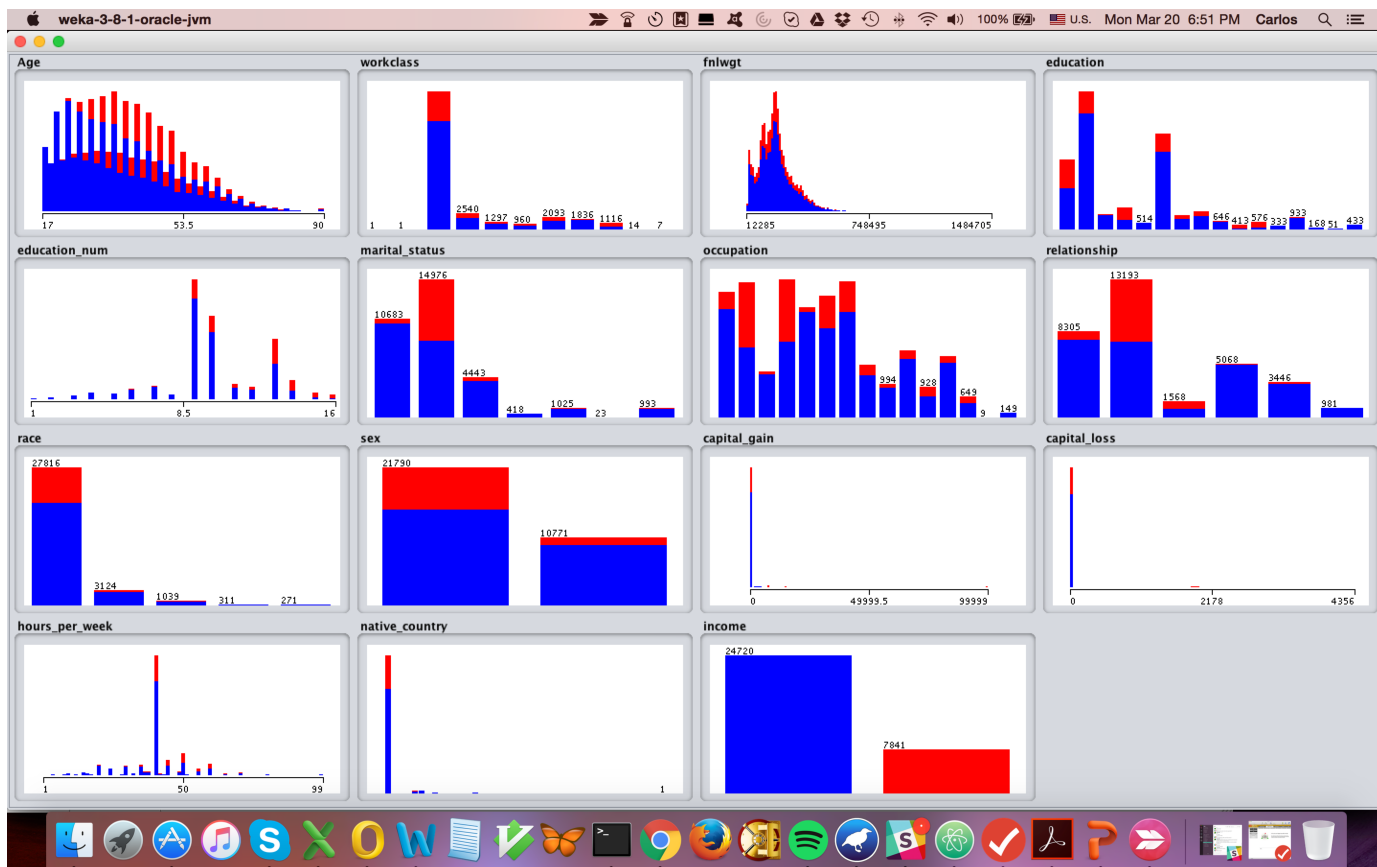
Confusion Matrix:

		Prediction	
		0	1
Truth	0	10536	824
	1	1465	2235

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

APPENDIX B



APPENDIX C

```

1 import pandas as pd
2 import numpy as np
3 from sklearn.model_selection import train_test_split
4 from sklearn.model_selection import KFold
5 from sklearn.neighbors import KNeighborsClassifier
6 from sklearn import tree
7 from sklearn.metrics import precision_score, recall_score, f1_score, confusion_matrix
8 from sklearn.metrics import roc_curve, roc_auc_score
9 from sklearn.naive_bayes import GaussianNB
10 from sklearn.linear_model import LogisticRegression
11 from sklearn.preprocessing import scale
12
13
14 #-----
15 # show_detailed_results
16 # Show detailed score results for algo, along with confusion matrix
17 # in: accuracy score, true target values, predictions
18 # out:
19 #
20 def show_detailed_results(msg, ac, y_test, y_pred):
21
22     f1 = f1_score(y_test, y_pred)
23     rc = recall_score(y_test, y_pred)
24     pr = precision_score(y_test, y_pred)
25     cm = confusion_matrix(y_test, y_pred)
26     roc = roc_auc_score(y_test, y_pred)
27
28     print('\n          ' + msg)
29     print("          Accuracy.....: %.3f" % ac)
30     print("          Precision Score.: %.3f" % f1)
31     print("          Recall Score....: %.3f" % rc)
32     print("          F1 Score.....: %.3f" % pr)
33     print("          ROC AUC.....: %.3f" % roc)
34     print("\n          Confusion Matrix:")
35
36     print("          +-----+")
37     print("          |          Prediction          |")
38     print("          +-----+")
39     print("          |          0          |          1          |")
40     print("          +---+-----+")
41     print("          | 0 |          %5d          |          %5d          |" % (cm[0][0], cm[0][1]))
42     print("          Truth + |-----+")
43     print("          | 1 |          %5d          |          %5d          |" % (cm[1][0], cm[1][1]))
44     print("          +---+-----+")
45
46
47 #-----
48 # gets dictionary and swaps keys with values
49 # in : my_dict_in
50 # out: my_dict_out
51 #
52 def invert_dict(my_dict_in):
53     my_dict_out = {}
54     for x, y in my_dict_in.iteritems(): my_dict_out[y] = float(x)
55     return my_dict_out
56
57 #-----
58 # loads dataset
59 # in: none
60 # out: dataset as dataframe
61 #
62 def load_data(fname):
63     # define column names for data import
64     df_col_names = ['age', 'workclass', 'fnlwt', 'education', 'education_num', 'marital_status',
65                     'occupation', 'relationship', 'race', 'sex', 'capital_gain', 'capital_loss',
66                     'hours_per_week', 'native_country', 'income']
67     df_col_dtype = {'age': np.float64, 'fnlwt': np.float64, 'education_num': np.float64,
68                     'capital_gain': np.float64, 'capital_loss': np.float64, 'hours_per_week':
np.float64}
69     # read data
70     df = pd.read_csv(fname, index_col=None, header=None, na_values='?',
71                     names=df_col_names, dtype=df_col_dtype)

```

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

```

75
76 #-----
77 # general precessing
78 # Performs the following preprocessing to benefit all algorithms
79 # 1 - converts all columns in dataset to float
80 #     all categorical and ordinal features are converted
81 # 2 - remove rows with nan's
82 # 3 - remove redundant features
83 #
84 # in:  dataset as dataframe
85 # out: tuple with three elements
86 #     1 - source dataframe w/o nan's and redundant features
87 #     2 - source dataframe used to create np.arrays - all features as floats
88 #     3 - dictionary with mapping between original feature values and new number values
89 #
90 def general_preprocess(df):
91
92     # found 2399 rows with NaN in the following columns: workclass, occupation, native_country
93     # cannot think of a way to impute values, so will dropna
94     df.dropna(inplace=True)
95
96     # drop education since we have education_num
97     df.drop('education', axis=1, inplace=True)
98
99     df_source = df.copy()
100
101     col_map = {}
102
103     # create dictionary objects for each column containing nominal or ordinal values
104     workclass_dict = pd.Series(df.workclass.unique()).to_dict()
105     marital_status_dict = pd.Series(df.marital_status.unique()).to_dict()
106     occupation_dict = pd.Series(df.occupation.unique()).to_dict()
107     relationship_dict = pd.Series(df.relationship.unique()).to_dict()
108     race_dict = pd.Series(df.race.unique()).to_dict()
109     sex_dict = pd.Series(df.sex.unique()).to_dict()
110     native_country_dict = pd.Series(df.native_country.unique()).to_dict()
111     income_dict = pd.Series(df.income.unique()).to_dict()
112
113     # create inverted dictionary objects to update dataframe columns
114     workclass_dict_inv = invert_dict( workclass_dict )
115     marital_status_dict_inv = invert_dict( marital_status_dict )
116     occupation_dict_inv = invert_dict( occupation_dict )
117     relationship_dict_inv = invert_dict( relationship_dict )
118     race_dict_inv = invert_dict( race_dict )
119     sex_dict_inv = invert_dict( sex_dict )
120     native_country_dict_inv = invert_dict( native_country_dict )
121     income_dict_inv = invert_dict( income_dict )
122
123     col_map['workclass'] = [ workclass_dict, workclass_dict_inv ]
124     col_map['marital_status'] = [ marital_status_dict, marital_status_dict_inv ]
125     col_map['occupation'] = [ occupation_dict, occupation_dict_inv ]
126     col_map['relationship'] = [ relationship_dict, relationship_dict_inv ]
127     col_map['race'] = [ race_dict, race_dict_inv ]
128     col_map['sex'] = [ sex_dict, sex_dict_inv ]
129     col_map['native_country'] = [ native_country_dict, native_country_dict_inv ]
130     col_map['income'] = [ income_dict, income_dict_inv ]
131
132     # convert each nominal column value to a number using inverted dictionaries
133     df.workclass.replace( workclass_dict_inv , inplace=True )
134     df.marital_status.replace( marital_status_dict_inv , inplace=True )
135     df.occupation.replace( occupation_dict_inv , inplace=True )
136     df.relationship.replace( relationship_dict_inv , inplace=True )
137     df.race.replace( race_dict_inv , inplace=True )
138     df.sex.replace( sex_dict_inv , inplace=True )
139     df.native_country.replace( native_country_dict_inv , inplace=True )
140     df.income.replace( income_dict_inv , inplace=True )
141
142     return (df_source, df, col_map)
143
144 #-----
145 # Get percentile using global percentile list
146 # Used by discretize function
147 # in: value from dataframe continuous feature
148 # out: return percentile corresponding to value from dataframe column using global percentiles list
149 #

```

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

```

153     elif (x <= percentiles[1]): return 2.0
154     elif (x <= percentiles[2]): return 3.0
155     elif (x <= percentiles[3]): return 4.0
156     elif (x <= percentiles[4]): return 5.0
157     elif (x <= percentiles[5]): return 6.0
158     elif (x <= percentiles[6]): return 7.0
159     elif (x <= percentiles[7]): return 8.0
160     elif (x <= percentiles[8]): return 9.0
161     else: return 10.0
162
163 #-----
164 # Discretize
165 # Replace continous feature with discrete value from 1-10 corresponding to it's percentile
166 # in: dataframe and name of the column to be discretize
167 #
168 def discretize(df, col_name):
169
170     # create list with percentiles
171     global percentiles
172     percentiles = []
173     x = []
174     for i in range(1,11,1): x.append(float(i)/10)
175     percentiles = df[col_name].quantile(x).tolist()
176
177     # discretize column
178     df[col_name] = df[col_name].apply(lambda x: getpercentile(x))
179
180     return df
181
182 #-----
183 # logit_preprocess
184 # Pre-processing for logistic regression algo
185 # will do the following:
186 # - standardize continuous variables
187 # - discretize income variable
188 # - create dummy variable for each categorical feature
189 # in: dataframe not ready for logit algo
190 # out: dataframe ready for logit algo
191 #
192 def logit_preprocess(df):
193
194     # standardize continous variables
195     df.age = scale(df.age)
196     df.fnlwgt = scale(df.fnlwgt)
197     df.education_num = scale(df.education_num)
198     df.capital_gain = scale(df.capital_gain)
199     df.capital_loss = scale(df.capital_loss)
200     df.hours_per_week = scale(df.hours_per_week)
201
202     # discretize income variable
203     col_map = {}
204     income_dict = pd.Series(df.income.unique()).to_dict()
205     income_dict_inv = invert_dict( income_dict )
206     col_map['income'] = [ income_dict, income_dict_inv ]
207     df.income.replace(income_dict_inv, inplace=True )
208
209     # create dummy variables for categorical features
210     df_dummies = pd.get_dummies(df)
211
212     return df_dummies
213
214 #-----
215 # nbayes_preprocess
216 # Naive Bayes specific preprocessing
217 # in: dataframe not ready for Naive Bayes algo
218 # out: dataframe ready for Naive Bayes algo
219 #
220 #
221 def nbayes_preprocess(df):
222
223     # Discretize all continous columns
224     df = discretize( df, 'age' )
225     df = discretize( df, 'fnlwgt' )
226     df = discretize( df, 'capital_gain' )
227     df = discretize( df, 'capital_loss' )

```

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

```

231
232
233 #-----
234 # Learn naive_bayes
235 # in: dataframe ready for naive bayes algo
236 # out:
237 #
238 def learn_naive_bayes(df, df_TEST):
239
240     # partition dataset
241     # will tune using stratified kfold and get final score on test set
242     X_train, X_test, y_train, y_test = train_test_split(
243         df.iloc[:, :13].values, df.iloc[:, 13:14].values, test_size=0.2, random_state=0)
244
245     # create naive bayes classifier object
246     clf = GaussianNB()
247
248     # will train model using kfold validation and compute average score
249     scores = []
250     kf = KFold(n_splits=10)
251
252     # split "Generate indices to split data into training and test set" test = validation
253     for train_idx, valid_idx in kf.split(X_train, y_train):
254
255         # load train and validation partition for this iteration of score computing
256         X_train_part = X_train[np.ravel(train_idx)]
257         y_train_part = y_train[np.ravel(train_idx)]
258         X_valid_part = X_train[np.ravel(valid_idx)]
259         y_valid_part = y_train[np.ravel(valid_idx)]
260
261         # learn/fit model for this fold
262         clf = clf.fit(X_train_part, y_train_part)
263
264         # calculate f1 score
265         y_pred = clf.predict(X_valid_part)
266         scores.append(f1_score(y_valid_part, y_pred) )
267
268     # display results
269     #
270     print('\nRESULTS FOR NAIVE BAYES CLASSIFIER')
271     print('-----')
272
273     # first results on training dataset - validation partition
274     #
275     ac = clf.score(X_test, y_test)
276     y_pred = clf.predict(X_test)
277     show_detailed_results('TRAIN VALIDATION PARTITION', ac, y_test, y_pred)
278
279     # then results on test dataset
280     #
281     X_train = df.copy()
282     X_train.drop('income', axis=1, inplace=True)
283     y_train = df.income
284     clf = clf.fit(X_train.values, y_train.values)
285     #
286     X_TEST = df_TEST.copy()
287     X_TEST.drop('income', axis=1, inplace=True)
288     y_TEST = df_TEST.income
289     ac = clf.score(X_TEST, y_TEST)
290     y_pred_TEST = clf.predict(X_TEST)
291     show_detailed_results('TEST DATASET', ac, y_TEST, y_pred_TEST)
292
293
294
295 #-----
296 # Learn decision tree
297 # in: dataframe ready for decision tree algo
298 # out:
299 #
300 def learn_decision_tree(df, df_TEST):
301
302     # partition dataset
303     # will tune using stratified kfold and get final score on test set
304     X_train, X_test, y_train, y_test = train_test_split(
305         df.iloc[:, :13].values, df.iloc[:, 13:14].values, test_size=0.2, random_state=0)

```

INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

```

309
310 for parm_value in np.arange(15)+1:
311
312     # create nvaive bayes classifier object with max_depth parameter
313     clf = tree.DecisionTreeClassifier(max_depth=parm_value)
314
315     # will train model using kfold validation and compute average score
316     scores = []
317     kf = KFold(n_splits=10)
318
319     # split "Generate indices to split data into training and test set" test = validation
320     for train_idx, valid_idx in kf.split(X_train, y_train):
321
322         # load train and validation partition for this iteration of score computing
323         X_train_part = X_train[np.ravel(train_idx)]
324         y_train_part = y_train[np.ravel(train_idx)]
325         X_valid_part = X_train[np.ravel(valid_idx)]
326         y_valid_part = y_train[np.ravel(valid_idx)]
327
328         # learn/fit model for this fold
329         clf = clf.fit(X_train_part, np.ravel(y_train_part))
330
331         # calculate f1 score to select best hyperparameter
332         y_pred = clf.predict( X_valid_part)
333         scores.append( f1_score(y_valid_part, y_pred) )
334
335     # track best score and corresponding hyperparameter value
336     if best_score < np.mean(scores):
337         best_score = np.mean(scores)
338         best_parm = parm_value
339
340     # repeat learning on full train partition with best parm value
341     clf = tree.DecisionTreeClassifier(max_depth=best_parm)
342     clf = clf.fit( X_train, np.ravel(y_train))
343
344     # calculate test scores
345     #
346     print('\nRESULTS FOR DECISION TREE CLASSIFIER')
347     print('-----')
348     print("Best tree depth = %i" % best_parm)
349
350     # first results on training dataset - validation partition
351     #
352     ac = clf.score( X_test, y_test )
353     y_pred = clf.predict(X_test)
354     show_detailed_results('TRAIN VALIDATION PARTITION', ac, y_test, y_pred)
355
356     # then results on test dataset
357     #
358     X_train = df.copy()
359     X_train.drop('income', axis=1, inplace=True)
360     y_train = df.income
361     clf = clf.fit( X_train.values, y_train.values )
362     #
363     X_TEST = df_TEST.copy()
364     X_TEST.drop('income', axis=1, inplace=True)
365     y_TEST = df_TEST.income
366     ac = clf.score( X_TEST, y_TEST)
367     y_pred_TEST = clf.predict(X_TEST)
368     show_detailed_results('TEST DATASET', ac, y_TEST, y_pred_TEST)
369
370
371 #-----
372 # Learn knn algorithm
373 # in: dataframe ready for knn algo
374 # out:
375 #
376 def learn_knn(df, df_TEST):
377
378     # partition dataset
379     # will tune using stratified kfold and get final score on test set
380     X_train, X_test, y_train, y_test = train_test_split(
381         df.iloc[:, :13].values, df.iloc[:, 13:14].values, test_size=0.2, random_state=0)
382
383     best_parm = 0

```


INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

```

465
466 l = ['l1'] * 5 + ['l2'] * 5
467 C = [1.0, 10.0, 100.0, 1000.0, 1000000.0] * 2
468
469 for l_, C_ in zip(l,C):
470
471     # create nvaive bayes classifier object with max_depth parameter
472     clf = LogisticRegression(penalty=l_, C=C_)
473
474     # will train model using kfold validation and compute average score
475     scores = []
476     kf = KFold(n_splits=10)
477
478     # split "Generate indices to split data into training and test set" test = validation
479     for train_idx, valid_idx in kf.split(X_train, y_train):
480
481         # load train and validation partition for this iteration of score computing
482         X_train_part = X_train[np.ravel(train_idx)]
483         y_train_part = y_train[np.ravel(train_idx)]
484         X_valid_part = X_train[np.ravel(valid_idx)]
485         y_valid_part = y_train[np.ravel(valid_idx)]
486
487         # learn/fit model for this fold
488         clf = clf.fit(X_train_part, np.ravel(y_train_part))
489
490         # calculate f1 score to select best hyperparameter
491         y_pred = clf.predict( X_valid_part)
492         scores.append( f1_score(y_valid_part, y_pred) )
493
494     # track best score and corresponding hyperparameter value
495     if best_score < np.mean(scores):
496         best_score = np.mean(scores)
497         best_parm = [ l_, C_ ]
498
499     # repeat learning on full train partition with best parm value
500     clf = LogisticRegression(penalty=best_parm[0], C=best_parm[1])
501     clf = clf.fit( X_train, np.ravel(y_train) )
502
503     # calculate test scores
504     #
505     print('\nRESULTS FOR LOGISTIC REGRESSION CLASSIFIER')
506     print('-----')
507     print("Best penalty= " + best_parm[0])
508     print("Best C value= " + str(best_parm[1]))
509
510     # first results on training dataset - validation partition
511     #
512     ac = clf.score( X_test, y_test )
513     y_pred = clf.predict(X_test)
514     show_detailed_results('TRAIN VALIDATION PARTITION', ac, y_test, y_pred)
515
516     # then results on test dataset
517     #
518     X_train = df.copy()
519     X_train.drop('income', axis=1, inplace=True)
520     # removing column which doesn't exist in TEST dataset
521     X_train.drop('native_country_Holand-Netherlands', axis=1, inplace=True)
522     y_train = df.income
523     clf = clf.fit( X_train.values, y_train.values )
524     #
525     X_TEST = df_TEST.copy()
526     X_TEST.drop('income', axis=1, inplace=True)
527     y_TEST = df_TEST.income
528     ac = clf.score( X_TEST, y_TEST)
529     y_pred_TEST = clf.predict(X_TEST)
530     show_detailed_results('TEST DATASET', ac, y_TEST, y_pred_TEST)
531
532
533
534 #####
535 #
536 if __name__ == '__main__':
537
538     # load data
539     df = load_data('adult.data')

```


INFO-I526 APPLIED MACHINE LEARNING – SPRING 2017

Carlos Sathler (cssathler@gmail.com)

```
543 # called "general" because all algorithms will have this preprocessing in common
544 # - get df_source back without nan rows and without redundant features
545 # - get df_floats with all categorical and nominal features converted to float
546 # - get dictionary with mapping info for each converted field
547 #
548 (df, df_floats, df_col_map) = general_preprocess( df.copy() )
549 (df_TEST, df_floats_TEST, df_col_map_TEST) = general_preprocess( df_TEST.copy() )
550
551 # performs pre-processing for naive bayes algorithm
552 # - get df_nbayes with continous features properly discretized
553 df_nbayes = nbayes_preprocess( df_floats.copy() )
554 df_nbayes_TEST = nbayes_preprocess( df_floats_TEST.copy() )
555 learn_naive_bayes( df_nbayes, df_nbayes_TEST )
556
557 # performs pre-processing for decision tree algorithm
558 # same as naive bayes
559 df_dtree = df_nbayes.copy()
560 df_dtree_TEST = df_nbayes_TEST.copy()
561 #learn_decision_tree( df_dtree, df_dtree_TEST )
562 learn_decision_tree( df_floats, df_floats_TEST ) # performed better with continous features
563
564 # performs pre-processing for knn algorithm
565 # will use discretized functions to increase clustering of datapoints
566 learn_knn( df_nbayes, df_nbayes_TEST ) # performed better with discretized
features
567 # also trying using continous feature values
568 df_knn = df_floats.copy()
569 df_knn_TEST = df_floats_TEST.copy()
570 #learn_knn( df_knn, df_knn_TEST )
571
572 # performs pre-processing for logistic regression algorithm
573 df_logit = logit_preprocess( df.copy() )
574 df_logit_TEST = logit_preprocess( df_TEST.copy() )
575 learn_logit( df_logit, df_logit_TEST )
```