

Ahmed Awadallah

Imputation and Multi Classification

Introduction

Hospitals have a direct interest in the future health of their patients while simultaneously also collecting and recording a lot of data on patients as part of their standard procedure. This provides a unique opportunity to create models to further benefit the patients. For this specific dataset the goal is to predict whether a patient will or will not be readmitted in the near future.

Methods

Our data comes from the DS 7333 Case Study 2. Our objective was to impute the data and then perform multiclass classification using logistic regression.

The dataset was shuffled to remove any chance of an inherited ordering impacting the model.

The features admission_type_id, discharge_disposition_id, and admission_source_id were coded with numerical values. A supplemental dataset supplied mappings from the code to descriptive states relating to the feature. All of these features had their codes mapped to their descriptive counterparts.

The dataset was explored for any possible na like value. Upon an intensive search na values were represented by “?” or “unknown/invalid”. The features that contained possible na variants were weight, race, diag_1, diag_2, diag_3, gender, payer_code, medical_specialty, admission_type_id, discharge_disposition_id, and admission_source_id. Upon inspection of how many NAs each of these features had the following percentages per feature.

- race : 2.234 %
- gender : 0.003 %
- weight : 96.858 %
- admission_type_id : 5.199 %
- discharge_disposition_id : 3.627 %
- admission_source_id : 6.663 %
- payer_code : 39.557 %
- medical_specialty : 49.082 %
- diag_1 : 0.021 %
- diag_2 : 0.352 %
- diag_3 : 1.398 %

Based on this data the following decisions about how to deal with these features were made.

The weight variable would be completely removed as too much of it was NA.

Race, gender, admission_type_id, discharge_disposition_id, admission_source_id, diag_1, diag_2, and diag_3 would all have their NA variants converted to the more suitable type None and then imputed using the highest frequency category within their respective feature. This method was chosen since there was only a small percentage of the row as NA.

Finally, payer_code and medical_specialty would keep their value of “?” as a special category indicating something was different about these rows. This was decided as too much was NA to consider imputing but not an overwhelming amount to have to completely remove the feature. Additionally, I’ve never used this method so I was curious about this approach.

The target feature, readmitted, was encoded so that it could be used in a logistic regression model.

Continuous variables were scaled using standardization to avoid the model biasing towards features with larger numbers.

Categorical features were one hot encoded so that they could be fed into a logistic regression model.

An optimal hyperparameter C was chosen using a custom made search function through a log scale. This hyperparameter was optimized using macro f1 score. Macro f1 since I wanted all the classes to be equally weighted for overall performance.

Using the optimal C a final model was trained and analyzed.

Results

Overall the performance of the model is not terrible but not great. Our final model achieved a macro f1 score of 0.4 with the following classification report.

Classification Report						
			precision	recall	f1-score	support
		NO	0.63	0.85	0.72	54864
		<30	0.49	0.05	0.10	11357
		>30	0.53	0.39	0.45	35545
	accuracy				0.60	101766
	macro avg		0.55	0.43	0.42	101766
	weighted avg		0.58	0.60	0.56	101766

Figure 1: Classification report for all three classes

The accuracy was above 60% which is atleast above randomly guessing the most popular class (the NO class). As a reminder the target for this model was hospital readmittance. "NO" means the patient was not readmitted, "<30" means the patient was readmitted in less than 30 days, and ">30" means the patient was readmitted in over 30 days.

For the "NO" class performance was decent. "NO" also happens to be the most populated class. The additional data may factor into why the model did a better job predicting it. Additionally "NO" is fundamentally different from the "<30" and ">30" classes.

For the "<30" class the performance overall is poor and it is the least populated class. Its poor performance is likely a mix of low data, lack of directly correlated features, and lack of a conceptual difference between "<30" and ">30".

Lastly the ">30" class had ok performance. This category leaves a lot to be desired but the model is at least able to make some sensical predictions.

Below are the ROC curves for all three models and their AUC metrics. The ROC curve gives further evidence that these models are unfortunately unable to achieve a respectable True Positive Rate without sacrificing too many False Positives. This is further reflected in the AUC metrics being far from the maximal value of 1. An interesting observation to make is "<30" ROC curve is actually better than ">30" which implies its model is better as well. But we know from the classification report that the ">30" performed better.

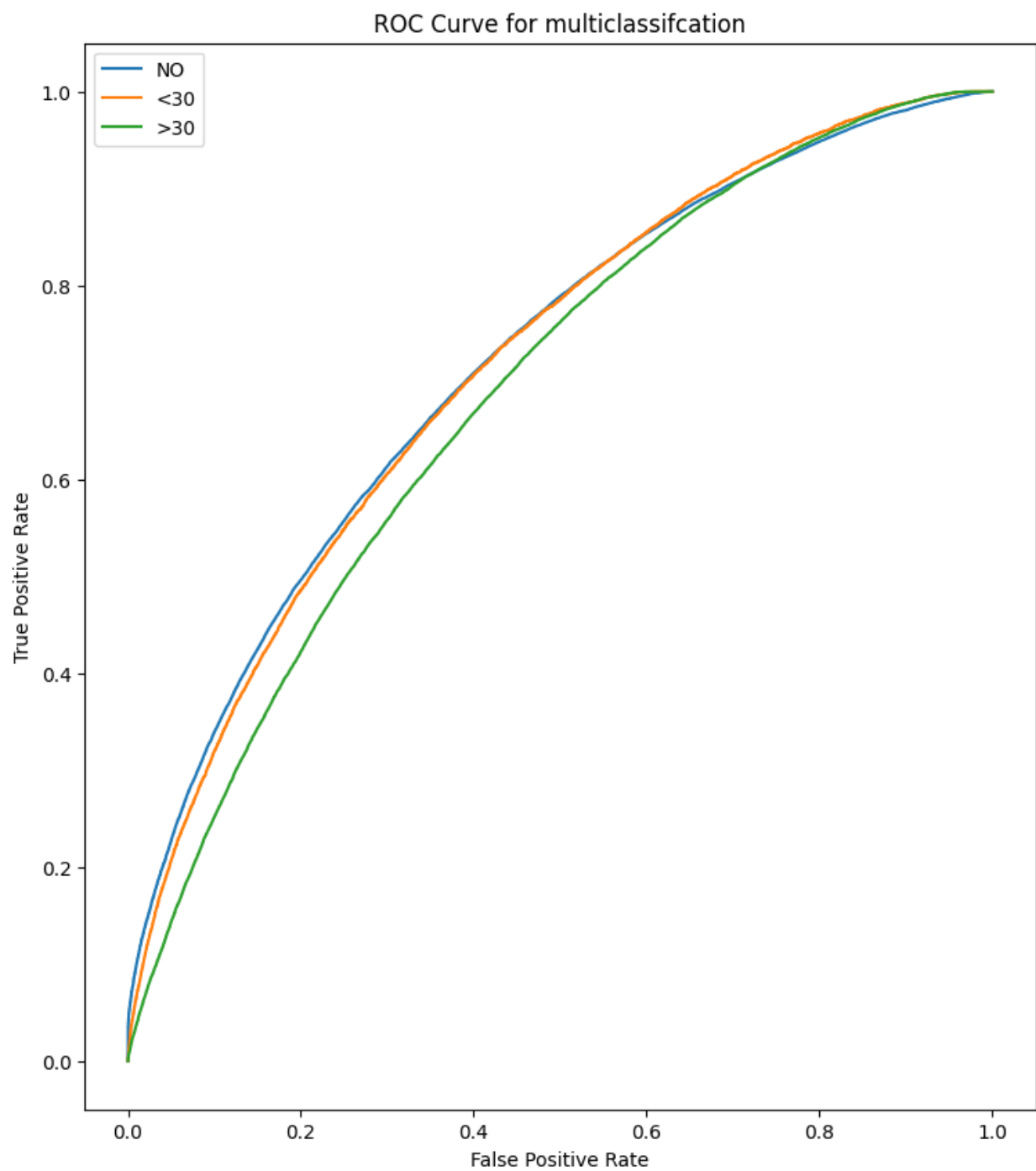
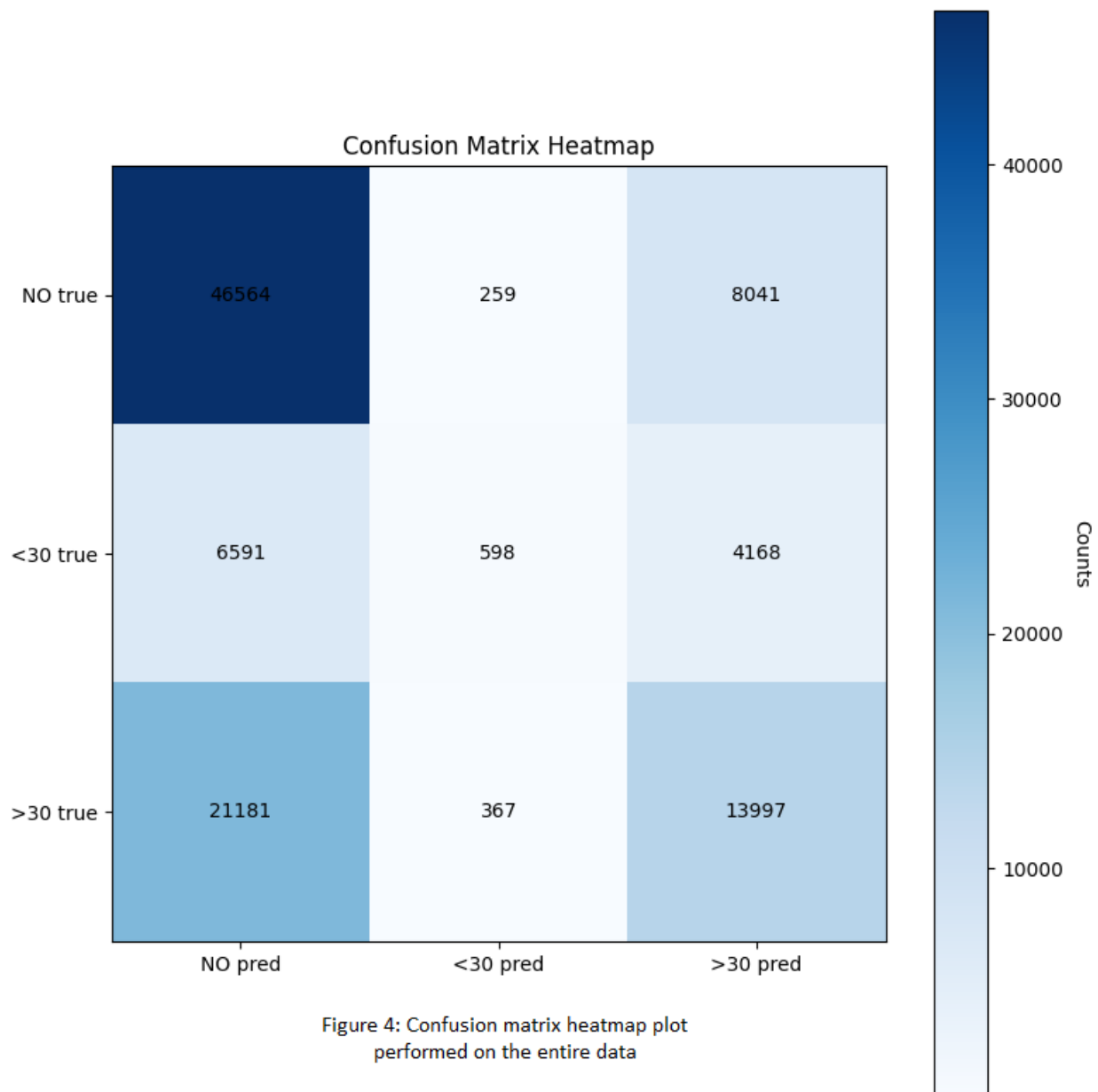


Figure 2: ROC Curve for each classes one versus rest logistic regression model

```
ROC/AUC Metrics
~~~~~
NO AUC: 0.7199
<30 AUC: 0.7170
>30 AUC: 0.6878
```

Figure 3: AUC for each class

In addition to the above a confusion matrix for our models decision boundary is displayed below.



The “<30” class barely got any predictions (1224) despite having a support of 11357. This was known from the classification report but this further shows that this class was hard to predict..

There is some evidence to suggest that a better model could be made if we changed the target variable of the model to be “NO” and “readmitted”. This would be instead of splitting the readmitted category into two time checkpoints. The reason for this suggestion is that 4168 of

our true “<30” were classified as “>30” as well as there was barely any “<30” predictions compared to the overall size of the dataset.

Lastly in the “NO” pred there were more “>30” and “<30” predictions than in the “<30” pred and “>30” pred columns. This displays that there is some difficulty distinguishing between “NO” and the other classes. Perhaps more useful features need to be researched for a future model.

The following is a figure of the top 8 most important weights for each class.

```
Most Relevant Coefs Report
~~~~~
Class: NO | discharge_disposition_id_Expired: 3.9772
Class: NO | discharge_disposition_id_Hospice / medical facility: 1.4593
Class: NO | diag_2_513: -1.1361
Class: NO | diag_3_250.91: -1.1323
Class: NO | diag_1_526: -1.0922
Class: NO | admission_type_id_Trauma Center: 1.0744
Class: NO | diag_3_V60: -1.0294
Class: NO | diag_3_481: -1.0168
-----
Class: <30 | discharge_disposition_id_Expired: -2.0111
Class: <30 | diag_1_271: 1.3030
Class: <30 | diag_3_156: 1.2925
Class: <30 | diag_2_136: 1.2416
Class: <30 | diag_1_643: 1.1953
Class: <30 | diag_1_356: 1.1660
Class: <30 | diag_1_82: 1.1329
Class: <30 | diag_1_358: 1.1327
-----
Class: >30 | discharge_disposition_id_Expired: -1.9661
Class: >30 | discharge_disposition_id_Hospice / medical facility: -1.3181
Class: >30 | diag_3_481: 1.3135
Class: >30 | diag_3_150: -1.2852
Class: >30 | diag_1_583: 1.1994
Class: >30 | diag_2_813: -1.1275
Class: >30 | diag_2_136: -1.0879
Class: >30 | diag_1_199: -1.0714
-----
```

Figure 5: Top 8 Weights absolute weights for each class.
Negative weight values restored to maintain proper interpretation.

From a glance it appears the various diag codes do have useful information about patient readmittance. A deeper study and analysis of the listed diag code is warranted and may provide additional knowledge about the domain.

For the named weights there are only three unique types. Discharge Disposition Expired, Discharge Disposition Hospice, and Discharge Disposition Trauma Center.

Discharge Disposition Expired refers to patients who are discharged from the hospital because they died within the hospital. To see a high association of this feature with "NO"(3.97) and a negative association with other two classes(-2.01, -1.966) is not surprising. A dead patient isn't going to be readmitted to a hospital. It's reassuring to see the model picked this feature up properly but it's also a feature that is completely trivial to understand without a model.

Discharge Disposition Hospice refers to patients who were discharged into a hospice. A hospice is a place where terminally ill patients are taken care of until they die. Similar to discharge disposition expired we see similar weight associations for the "NO"(positive) and ">30" classes(negative) although it is noticeably missing from the "<30" class. This may indicate hospice patients may in fact bounce between the hospice and being readmitted but only in short time spans.

Discharge Disposition Trauma Center refers to patients who have been discharged into a trauma center. Trauma centers deal with patients who have suffered from trauma type injuries such as falls, car crashes, or firearms. For this feature we only see a positive association with the "NO" class. This association is possibly because traumatic injuries generally lead to two future conditions. Either the patient is able to be patched up and heals their injuries over time or the injuries are too traumatic and they pass away. Both cases would not lead to a hospital readmittance.

Finally, it's interesting to note that "<30" and ">30" largely picked different weights. While the model did perform poorly guessing the "<30" category, the difference in weights provides some evidence that it may be possible. If predicting readmittance of <30 days is important further research into better correlated features may be warranted.

A figure displaying the weights associated with race is shown below.

```

Finding the Weights of the Race Category
~~~~~
Class: NO | race_Asian: 0.1828
Class: NO | race_Caucasian: -0.0290
Class: NO | race_Hispanic: 0.0794
Class: NO | race_Other: 0.1315
-----
Class: <30 | race_Asian: 0.0381
Class: <30 | race_Caucasian: 0.0196
Class: <30 | race_Hispanic: 0.0091
Class: <30 | race_Other: -0.0547
-----
Class: >30 | race_Asian: -0.2209
Class: >30 | race_Caucasian: 0.0094
Class: >30 | race_Hispanic: -0.0885
Class: >30 | race_Other: -0.0768
-----

```

Figure 3: Weights for each race in each model

Race can be a controversial feature to discuss but we felt it's important to not let cultural stigma prevent us from exploring what reality is telling us. If race is in fact an important factor for predicting patient readmittance there would be a direct benefit with recording this information.

The weights for each model tell us that some races do have some predictive power in hospital readmittance. For example asian has a positive association with not being readmitted while caucasian has near 0 relationship with all three classes. This discrepancy is curious and further research is probably warranted. A hypothesis test on the race should be conducted to see if the weights are in fact non zero as well as if the additional predictive power these features give is significant for a model. Only when both these conditions have been confirmed can we say race has a definitive and useful impact.

Conclusion

The goal of predicting hospital readmittance was achieved but there is a lot of improvement that must be made if the model is to be considered exceptional. Regardless, interesting information was learned about the features embedded within the model and their exploration may provide further insight into patient care.

Code

Appendix A

```
1  import pandas as pd
2  import numpy as np
3  import math
4  import sys
5  from sklearn.impute import SimpleImputer
6  from sklearn.preprocessing import StandardScaler
7  from sklearn.linear_model import LogisticRegression
8  from sklearn.model_selection import cross_val_score
9  import time
10 from sklearn.metrics import classification_report
11 from sklearn.metrics import confusion_matrix
12 import matplotlib.pyplot as plt
13 from sklearn import metrics as mt
14
15 # Your case study is to build a classifier using logistic regression to
16 # predict hospital readmittance. There is missing data that must be imputed.
17 # Once again, discuss variable importances as part of your submission.
18
19
20 INVALID_INPUT = 1
21
22
23 def try_to_read_input_csvs(inputs):
24     try:
25         read_input_csvs(inputs)
26     except Exception:
27         print("This script expects the two csv files associated with Case "
28               "Study1. You either input nothing, 1 file path, or one or more "
29               "of your inputs contains an invalid path.")
30         sys.exit(INVALID_INPUT)
31
32
33 def read_input_csvs(inputs):
34     global df1, df2
35     df1_path = inputs[1]
36     df2_path = inputs[2]
37     df1 = pd.read_csv(df1_path)
38     df2 = pd.read_csv(df2_path)
```

```

41 def get_mapping(range, id_map_df):
42     map_dict = {}
43     for i in range:
44         id = int(id_map_df.loc[i][0])
45         descrip = id_map_df.loc[i][1]
46         if type(descrip) == float and math.isnan(descrip):
47             descrip = None
48         map_dict[id] = descrip
49     return map_dict

```

```

52 v def logreg_cv_compare(x, y, print_progress=False, n_jobs=1):
53     best_C = 0
54     best_f1 = 0
55     first = True
56     iter_count = 20
57     report_freq = int(iter_count/20)
58     i = 0
59     t1 = time.time()
60 v for param_C in np.logspace(-6, 0, iter_count):
61
62 v     model = LogisticRegression(multi_class='multinomial',
63                               solver='lbfgs',
64                               C=param_C,
65                               n_jobs=6,
66                               max_iter=1000)
67
68 v     f1_cv = cross_val_score(model, x, y, cv=5,
69                             scoring='f1_macro')
70
71     f1 = sum(f1_cv)/len(f1_cv)
72
73 v     if first:
74         first = False
75         best_C = param_C
76         best_f1 = f1
77 v     else:
78 v         if (f1 > best_f1):
79             best_C = param_C
80             best_f1 = f1
81
82 v     if ((i+1) % report_freq == 0) and print_progress:
83 v         print(f"Iteration {i+1}. "
84               f"Percent done = {(i+1)/iter_count*100:.4f}%")
85
86     i += 1
87     t2 = time.time()
88     elapsed_time = t2-t1
89
90     report_str = f"Best f1={best_f1:.4f}, "
91     report_str += f"C={best_C}, time={elapsed_time:.4f}s\n"
92
93     return best_C, report_str

```

```

96 def get_top_n_coef_multiclass_logreg(feature_list, coef_list, n=5):
97     feature_scores_abs = {}
98     top_n_per_class = []
99     feature_scores = {}
100    class_count = coef_list.shape[0]
101    for class_id in range(class_count):
102        for i in range(len(feature_list)):
103            feature = feature_list[i]
104            coef = coef_list[class_id, i]
105            feature_scores_abs[feature] = abs(coef)
106            feature_scores[feature] = coef
107
108            feature_scores_abs = dict(sorted(feature_scores_abs.items(),
109                                           key=lambda item: item[1], reverse=True))
110            top_n_features = list(feature_scores_abs.keys())[0:n]
111
112            top_feature_scores = {}
113            for feature in top_n_features:
114                top_feature_scores[feature] = feature_scores[feature]
115
116            top_n_per_class.append(top_feature_scores)
117
118    return top_n_per_class
119
120
121 def print_top_coefs_multiclass_logreg(top_n_coefs, id_targ_map):
122     str_report = ""
123     for class_id in range(len(top_n_coefs)):
124         feature_top_n = top_n_coefs[class_id]
125         class_name = id_targ_map[class_id]
126         for feature in feature_top_n.keys():
127             coef = feature_top_n[feature]
128             str_report += f"Class: {class_name} | {feature}: {coef:.4f}\n"
129         str_report += "-----\n"
130    return str_report

```

```

133  def plot_heatmap(data,
134                  xlab,
135                  ylab,
136                  size=9,
137                  title="untitled heatmap",
138                  cbar_label="untitled",
139                  save_name="untitled.png"):
140
141      fig, ax = plt.subplots(figsize=(size, size))
142      im = ax.imshow(data, cmap=plt.cm.Blues)
143
144      for i in range(len(ylab)):
145          for j in range(len(xlab)):
146              ax.text(j, i, data[i, j],
147                     ha="center", va="center", color="black")
148
149      # Show all ticks and label them with the respective list entries
150      ax.set_xticks(np.arange(len(xlab)), labels=xlab)
151      ax.set_yticks(np.arange(len(ylab)), labels=ylab)
152
153      cbar = ax.figure.colorbar(im, ax=ax)
154      cbar.ax.set_ylabel(cbar_label, rotation=-90, va="bottom")
155
156      ax.set_title(title)
157      fig.tight_layout()
158      plt.savefig(save_name, bbox_inches='tight')
159      plt.clf()
160
161
162  if __name__ == '__main__':
163      log_file = open("log.txt", "w")
164      log_file.write("Case Study 2 Report\n\n")
165
166      inputs = sys.argv
167      try_to_read_input_csvs(inputs)

```

```

169 # Find which dataframe contains the cts variables that need to be scaled
170 if (df1.shape == (101766, 50)) and (df2.shape == (67, 2)):
171     raw_data = df1
172     id_map_df = df2
173 elif (df2.shape == (101766, 50)) and (df1.shape == (67, 2)):
174     raw_data = df2
175     id_map_df = df1
176 else:
177     raise Exception("One or both of dataframes has unexpected shape")
178
179 # Shuffle raw data
180 raw_data = raw_data.sample(frac=1)
181
182 # Use metadata to provide useful descriptions for some features
183 admission_type_map = get_mapping(range(0, 8), id_map_df)
184 discharge_map = get_mapping(range(10, 40), id_map_df)
185 admission_source_map = get_mapping(range(42, 67), id_map_df)
186
187 raw_data["admission_type_id"] = raw_data[
188     "admission_type_id"].map(admission_type_map)
189 raw_data["discharge_disposition_id"] = raw_data[
190     "discharge_disposition_id"].map(discharge_map)
191 raw_data["admission_source_id"] = raw_data[
192     "admission_source_id"].map(admission_source_map)
193
194 # Data Cleaning Decisions
195 # Remove weight as a feature as too much of it is missing
196 clean_df = raw_data.drop(["weight"], axis=1)
197 # Remove the following features as they should have no relationship with
198 # target
199 features_to_remove = ["patient_nbr", "encounter_id"]
200 clean_df = clean_df.drop(features_to_remove, axis=1)
201 # Impute desired features
202 clean_df["gender"] = clean_df["gender"].replace("Unknown/Invalid", "?")
203 features_to_impute = ["race", "diag_1", "diag_2", "diag_3", "gender",
204     "admission_type_id", "admission_source_id",
205     "discharge_disposition_id"]
206 # All these categories use "?" to signify na data so replace with None
207 for feature in features_to_impute:
208     clean_df[feature] = clean_df[feature].replace("?", None)
209 imputer = SimpleImputer(missing_values=None, strategy="most_frequent")
210 for feature in features_to_impute:
211     imputed_feature = imputer.fit_transform(clean_df[[feature]]).ravel()
212     clean_df[feature] = imputed_feature

```

```

214 # Encode Target
215 readmitted_to_id_map = {"NO": 0, "<30": 1, ">30": 2}
216 id_to_readmitted_map = {0: "NO", 1: "<30", 2: ">30"}
217 clean_df["readmitted"] = clean_df["readmitted"].map(readmitted_to_id_map)
218
219 # Specifying cts v categorical features
220 v cts_features = ["time_in_hospital",
221                  "num_lab_procedures",
222                  "num_procedures",
223                  "num_medications",
224                  "number_outpatient",
225                  "number_emergency",
226                  "number_inpatient",
227                  "number_diagnoses"]
228
229 categ_features = []
230 v for feature in list(clean_df):
231 v     if feature not in cts_features:
232         categ_features.append(feature)
233
234 # Scaling Cts Data
235 scaler = StandardScaler()
236 clean_df[cts_features] = scaler.fit_transform(clean_df[cts_features])
237
238 # Split data into train data and target data
239 y = clean_df["readmitted"]
240 x = clean_df.drop("readmitted", axis=1)
241
242 # One hot encode
243 x = pd.get_dummies(x, drop_first=True)
244
245 # Subset data for demonstration purposes
246 subset_count = 500
247 v print("Data subsetting is ACTIVE, only", subset_count,
248       "samples being used to train")
249 print("This is just to demonstrate the code behavior quickly")
250 x = x.head(subset_count)
251 y = y.head(subset_count)

```

```

253 # Search optimal regularization C
254 print("Finding optimal model hyperparameters...")
255 print("")
256 optimal_c, report_str = logreg_cv_compare(
257     x,
258     y,
259     print_progress=True,
260     n_jobs=4)
261 log_file.write("Optimization Report\n")
262 log_file.write("~~~~~\n")
263 log_file.write(report_str)
264 log_file.write("\n")
265
266 # Final Model
267 final_model = LogisticRegression(multi_class='multinomial',
268     solver='lbfgs',
269     C=optimal_c,
270     max_iter=1000,
271     n_jobs=6)
272 final_model.fit(x, y)
273 y_pred = final_model.predict(x)
274
275 # Final Model Analysis
276
277 # Top Coefs Per Target Feature Class
278 class_top_n = get_top_n_coef_multiclass_logreg(
279     list(x),
280     final_model.coef_,
281     n=8)
282
283 log_file.write("Most Relevant Coefs Report\n")
284 log_file.write("~~~~~\n")
285 log_file.write(
286     print_top_coefs_multiclass_logreg(
287         class_top_n,
288         id_to_readmitted_map
289     )
290 )
291 log_file.write("\n")

```

```

293 # Print Race Coef
294 log_file.write("Finding the Weights of the Race Category\n")
295 log_file.write("~~~~~\n")
296 races = ['race_Asian',
297          'race_Caucasian',
298          'race_Hispanic',
299          'race_Other',
300          'race_AfricanAmerican']
301 class_race_score = []
302 class_count = final_model.coef_.shape[0]
303 for class_id in range(class_count):
304     race_feature_score = {}
305
306     for race in races:
307         if race in list(x):
308             race_ind = list(x).index(race)
309             coef = final_model.coef_[class_id, race_ind]
310
311             race_feature_score[race] = coef
312
313     class_race_score.append(race_feature_score)
314 log_file.write(
315     print_top_coefs_multiclass_logreg(
316         class_race_score,
317         id_to_readmitted_map
318     )
319 )
320 log_file.write("\n")
321
322 # Classif Report
323 log_file.write("Classification Report\n")
324 log_file.write("~~~~~\n")
325 log_file.write(
326     classification_report(
327         y,
328         y_pred,
329         target_names=id_to_readmitted_map.values()
330     )
331 )
332 log_file.write("\n")
333
334 # Conf Matrix Plot
335 conf_mat = confusion_matrix(y, y_pred)

```



```

337 base_labels = list(id_to_readmitted_map.values())
338 label_true = []
339 label_pred = []
340 for i in range(len(base_labels)):
341     true_label = base_labels[i] + " true"
342     pred_label = base_labels[i] + " pred"
343     label_pred.append(pred_label)
344     label_true.append(true_label)
345
346 plot_heatmap(conf_mat,
347              label_pred,
348              label_true,
349              size=8,
350              title="Confusion Matrix Heatmap",
351              cbar_label="Counts",
352              save_name="conf_heatmap.png")
353
354 # ROC Curve and AUC
355 log_file.write("ROC/AUC Metrics\n")
356 log_file.write("~~~~~\n")
357 y_pred_prob = final_model.predict_proba(x)
358 for i in range(y_pred_prob.shape[1]):
359     y_label = id_to_readmitted_map[i]
360     preds = y_pred_prob[:, i]
361     fpr, tpr, thresholds = mt.roc_curve(y, preds, pos_label=i)
362     plt.plot(fpr, tpr, label=y_label)
363     title_str = 'ROC Curve for multiclassifcation'
364     plt.title(title_str)
365     plt.xlabel('False Positive Rate')
366     plt.ylabel('True Positive Rate')
367     log_file.write(f"{y_label} AUC: {mt.auc(fpr, tpr):.4f}\n")
368
369 plt.legend()
370 plt.savefig("roc_plot.png", bbox_inches='tight')
371 plt.clf()
372 log_file.write("\n")
373 log_file.close()

```