


```
In [22]: log_reg_sag = LogisticRegression(solver='sag', max_iter=10, penalty='none')
log_reg_sag.fit(df_train, df_sick_train)

test_predictions_1 = log_reg_sag.predict(df_test)

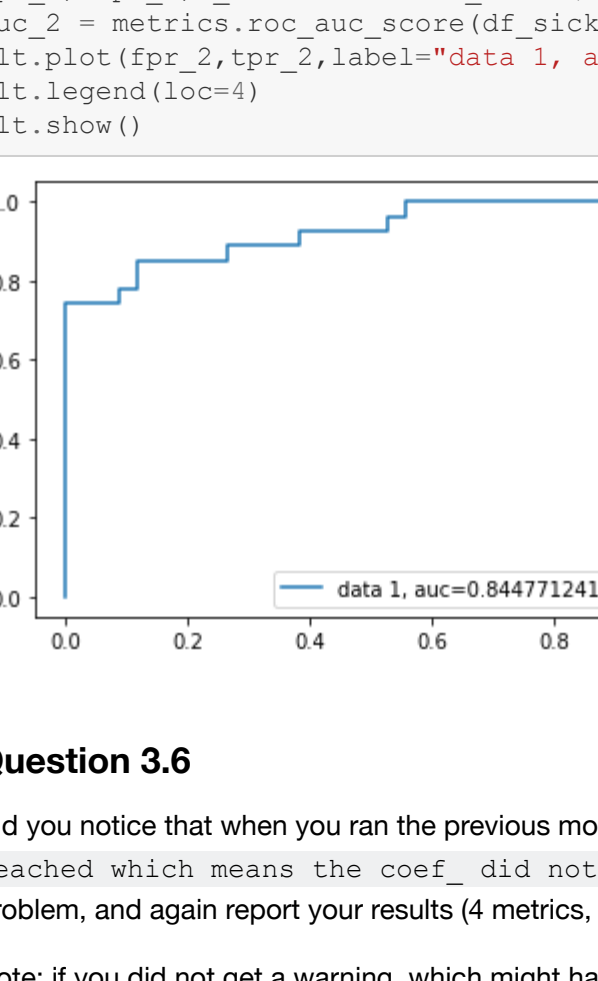
print("Accuracy: %.3f" % metrics.accuracy_score(df_sick_test, test_predictions_1))
print("Precision: %.3f" % metrics.precision_score(df_sick_test, test_predictions_1))
print("Recall: %.3f" % metrics.recall_score(df_sick_test, test_predictions_1))
print("F1 Score: %.3f" % metrics.f1_score(df_sick_test, test_predictions_1))

draw_confusion_matrix(df_sick_test, test_predictions_1, title_name="Confusion Matrix")

Accuracy: 0.852
Precision: 0.875
Recall: 0.778
F1 Score: 0.824

/opt/homebrew/lib/python3.9/site-packages/sklearn/linear_model/_sag.py:352: ConvergenceWarning: The max_iter was reached which means the coef_ did not converge
  warnings.warn()
```

Confusion Matrix



True label

Predicted label

```
In [23]: y_pred_proba_2 = log_reg_sag.predict_proba(df_test)[::,1]
fpr_2, tpr_2, _ = metrics.roc_curve(df_sick_test, y_pred_proba_2)
auc_2 = metrics.roc_auc_score(df_sick_test, test_predictions_2_fix)
plt.plot(fpr_2, tpr_2, label='data 1, auc='+str(auc_2))
plt.legend(loc=4)
plt.show()
```

ROC curve for data 1, auc=0.8447712418300652

Question 3.6

Did you notice that when you ran the previous model you got the following warning: `ConvergenceWarning: The max_iter was reached which means the coef_ did not converge.` Check the documentation and see if you can implement a fix for this problem, and again report your results (4 metrics, confusion matrix, ROC curve).

Note: if you did not get a warning, which might happen to those running this notebook in VSCode, please try running the following code, as described [here](#):

```
import warnings
warnings.simplefilter(action='default')
```

```
In [24]: log_reg_sag_fix = LogisticRegression(solver='sag', max_iter=1000, penalty='none')
log_reg_sag_fix.fit(df_train, df_sick_train)

test_predictions_2_fix = log_reg_sag_fix.predict(df_test)

print("Accuracy: %.3f" % metrics.accuracy_score(df_sick_test, test_predictions_2_fix))
print("Precision: %.3f" % metrics.precision_score(df_sick_test, test_predictions_2_fix))
print("Recall: %.3f" % metrics.recall_score(df_sick_test, test_predictions_2_fix))
print("F1 Score: %.3f" % metrics.f1_score(df_sick_test, test_predictions_2_fix))

draw_confusion_matrix(df_sick_test, test_predictions_2_fix, title_name="Confusion Matrix")

Accuracy: 0.852
Precision: 0.875
Recall: 0.778
F1 Score: 0.824

Confusion Matrix
```



True label

Predicted label

```
In [25]: y_pred_proba_2 = log_reg_sag_fix.predict_proba(df_test)[::,1]
fpr_2, tpr_2, _ = metrics.roc_curve(df_sick_test, y_pred_proba_2_fix)
auc_2_fix = metrics.roc_auc_score(df_sick_test, test_predictions_2_fix)
plt.plot(fpr_2, tpr_2, label='data 1, auc='+str(auc_2_fix))
plt.legend(loc=4)
plt.show()
```

ROC curve for data 1, auc=0.8447712418300652

Question 3.7

Explain what you changed and why this fixed the `ConvergenceWarning` problem. Are there any downsides of your fix? How might you have harmed the outcome instead? What other parameters you set may be playing a factor in affecting the results?

I increased `max_iter` from 10 to 1000, which allowed the algorithm to run more iterations, so it could converge. One potential downside of increasing maximum iterations is that slow convergence could be indicative of several problems, and increasing max iterations ignores these potential problems. Ignoring possible problems could result in an inaccurate model. The penalty we set and solver we chose could also be affecting our results.

Question 3.8

Rerun your logistic classifier, but modify the `penalty='l2'`, `solver='newton-cg'` and again report the results (4 metrics, confusion matrix, ROC curve)

```
In [26]: log_reg_new = LogisticRegression(solver='newton-cg', max_iter=100, penalty='l2')
log_reg_new.fit(df_train, df_sick_train)

test_predictions_3 = log_reg_new.predict(df_test)

print("Accuracy: %.3f" % metrics.accuracy_score(df_sick_test, test_predictions_3))
print("Precision: %.3f" % metrics.precision_score(df_sick_test, test_predictions_3))
print("Recall: %.3f" % metrics.recall_score(df_sick_test, test_predictions_3))
print("F1 Score: %.3f" % metrics.f1_score(df_sick_test, test_predictions_3))

draw_confusion_matrix(df_sick_test, test_predictions_3, title_name="Confusion Matrix")

Accuracy: 0.852
Precision: 0.875
Recall: 0.778
F1 Score: 0.824

Confusion Matrix
```



True label

Predicted label

```
In [27]: y_pred_proba_3 = log_reg_new.predict_proba(df_test)[::,1]
fpr_3, tpr_3, _ = metrics.roc_curve(df_sick_test, y_pred_proba_3)
auc_3 = metrics.roc_auc_score(df_sick_test, test_predictions_3)
plt.plot(fpr_3, tpr_3, label='data 1, auc='+str(auc_3))
plt.legend(loc=4)
plt.show()
```

ROC curve for data 1, auc=0.8447712418300652

Question 3.9

Explain how the 2 solvers work and how they differ from each other.

sag stands for Stochastic Average Gradient, and it's a variation of gradient descent. It takes random samples of previous gradient descent values, which makes it fast for big datasets. On the other hand, newton-cg uses an exact Hessian matrix, meaning it has to compute up to second-order derivatives, making it slow for large datasets.

Question 3.10

We also played around with different penalty terms (none, L1 etc.) Describe what the purpose of a penalty term is and the difference between L1 and L2 penalties.

The purpose of having a penalty term is for regularization, to tune the model to prevent overfitting. The L1 penalty is lasso regression and L2 penalty is ridge regression. Ridge regression adds a penalty to the cost function equivalent to square of the magnitude of the coefficients, putting a constraint on the coefficients. If the coefficients are large, the optimization function is penalized. Lasso regression does something similar, but instead of taking the square of the coefficients, it takes the magnitudes, so it can lead to zero coefficients. This helps with feature elimination.

Question 3.11 Support Vector Machine (SVM)

A support vector machine (SVM) is a discriminative classifier formally defined by a separating hyperplane. In other words, given labeled training data (supervised learning), the algorithm outputs an optimal hyperplane which categorizes new examples. In 2-D space this hyperplane is a line dividing a plane in two parts where in each class lay in either side.

Implement an SVM classifier on your pipelined data (recommend using `scikit-learn`) For this implementation you can simply use the default settings, but set `probability=True`.

```
In [28]: clf_1 = SVC(probability=True)
clf_1.fit(df_train, df_sick_train)
svm_pred_1 = clf_1.predict(df_test)

Confusion Matrix
```



True label

Predicted label

```
In [29]: print("Accuracy: %.3f" % metrics.accuracy_score(df_sick_test, svm_pred_1))
print("Precision: %.3f" % metrics.precision_score(df_sick_test, svm_pred_1))
print("Recall: %.3f" % metrics.recall_score(df_sick_test, svm_pred_1))
print("F1 Score: %.3f" % metrics.f1_score(df_sick_test, svm_pred_1))

draw_confusion_matrix(df_sick_test, svm_pred_1, title_name="Confusion Matrix")

Accuracy: 0.885
Precision: 0.917
Recall: 0.915
F1 Score: 0.863

Confusion Matrix
```



True label

Predicted label

```
In [30]: y_pred_proba_4 = clf_1.predict_proba(df_test)[::,1]
fpr_4, tpr_4, _ = metrics.roc_curve(df_sick_test, y_pred_proba_4)
auc_4 = metrics.roc_auc_score(df_sick_test, svm_pred_1)
plt.plot(fpr_4, tpr_4, label='data 1, auc='+str(auc_4))
plt.legend(loc=4)
plt.show()
```

ROC curve for data 1, auc=0.8779956427015251

Question 3.13

Rerun your SVM, but now modify your model parameter kernel to be `linear`. Again report your accuracy, precision, recall, F1 scores, and confusion matrix and plot the new ROC curve.


```
In [31]: clf_2 = SVC(probability=True, kernel='linear')
clf_2.fit(df_train, df_sick_train)
svm_pred_2 = clf_2.predict(df_test)

print("Accuracy: %.3f" % metrics.accuracy_score(df_sick_test, svm_pred_2))
print("Precision: %.3f" % metrics.precision_score(df_sick_test, svm_pred_2))
print("Recall: %.3f" % metrics.recall_score(df_sick_test, svm_pred_2))
print("F1 Score: %.3f" % metrics.f1_score(df_sick_test, svm_pred_2))

draw_confusion_matrix(df_sick_test, svm_pred_2, title_name="Confusion Matrix")

Accuracy: 0.852
Precision: 0.875
Recall: 0.778
F1 Score: 0.824

Confusion Matrix
```



True label

Predicted label

```
In [32]: y_pred_proba_5 = clf_2.predict_proba(df_test)[::,1]
fpr_5, tpr_5, _ = metrics.roc_curve(df_sick_test, y_pred_proba_5)
auc_5 = metrics.roc_auc_score(df_sick_test, svm_pred_2)
plt.plot(fpr_5, tpr_5, label='data 1, auc='+str(auc_5))
plt.legend(loc=4)
plt.show()
```

ROC curve for data 1, auc=0.8447712418300652

Question 3.14

Explain the what the new results you've achieved mean. Read the documentation to understand what you've changed about your model and explain why changing the kernel parameter might impact the results in the manner you've observed.

When we set kernel to 'linear', we can see that for our four metrics all decreased, which indicates that a linear kernel may not be the best fit for this dataset. By changing the kernel from the default 'rbf' to 'linear', we try to separate our data linearly, to see if it's better to use the simpler linear classifier. If our data fits well with this linear classifier, we would likely see our four metrics stay the same, or change slightly. If the linear classifier doesn't make sense for our data, then we would likely see our metrics reflect this.

Question 3.15

Both logistic regression and linear SVM are trying to classify data points using a linear decision boundary. How do they differ in how they try to find this boundary?

Logistic regression uses a logistic function to try and find the relationship between variables, while SVM creates a decision boundary to separate the classes, trying to find the best separator between the classes. SVM is less vulnerable to overfitting.

Question 3.16

We also learned about linear regression in class. Why is linear regression not a suitable model for this classification task?

We chose logistic rather than linear regression because we are not trying to predict a quantitative variable, but rather a categorical one. With classification problems, we're trying to place each observation into a category, thus we use logistic regression.

Statistical Classification Methods

Now we'll explore a statistical classification method, the naive Bayes classifier.

Naive Bayes methods are a set of supervised learning algorithms based on applying Bayes' theorem with the "naive" assumption of conditional independence between every pair of features given the value of the class variable. Bayes' theorem states the following relationship, given class variable C_i and dependent feature vector $\mathbf{bold{x}} = [x_1, x_2, \dots, x_d]^T$.

$$P(C_i | \mathbf{bold{x}}) = \frac{P(C_i) \prod \mathbf{bold{P}}(C_i)}{P(\mathbf{bold{x}})}$$

Note for our purposes, there are 2 possible classes (sick or not sick), so i ranges from 1 to 2.

Question 3.17

Implement a Naive Bayes Classifier on the pipelined data. Use the `GaussianNB` model. For this model, simply use the default parameters. Report out the number of mislabeled points that result (i.e. both the false positives and false negatives), along with the accuracy, precision, recall, F1 score and confusion matrix. Also, plot an ROC curve.

```
In [33]: gnb = GaussianNB()
gnb_class = gnb.fit(df_train, df_sick_train)
gnb_pred = gnb.classify(df_test)

print("Number of mislabeled points out of a total", df_test.shape[0], "points:", (df_sick_test != gnb_pred).sum())

print("Accuracy: %.3f" % metrics.accuracy_score(df_sick_test, gnb_pred))
print("Precision: %.3f" % metrics.precision_score(df_sick_test, gnb_pred))
print("Recall: %.3f" % metrics.recall_score(df_sick_test, gnb_pred))
print("F1 Score: %.3f" % metrics.f1_score(df_sick_test, gnb_pred))

draw_confusion_matrix(df_sick_test, gnb_pred, title_name="Confusion Matrix")

Number of mislabeled points out of a total 61 points: 11
Accuracy: 0.820
Precision: 0.808
Recall: 0.778
F1 Score: 0.792

Confusion Matrix
```



True label

Predicted label

```
In [34]: y_pred_gnb = gnb.class.predict_proba(df_test)[::,1]
fpr_gnb, tpr_gnb, _ = metrics.roc_curve(df_sick_test, y_pred_gnb)
auc_gnb = metrics.roc_auc_score(df_sick_test, gnb_pred)
plt.plot(fpr_gnb, tpr_gnb, label='data 1, auc='+str(auc_gnb))
plt.legend(loc=4)
plt.show()
```

ROC curve for data 1, auc=0.81535947124183

Question 3.18

Discuss the observed results. What assumptions about our data are there and why might those be inaccurate?

We observe that we produced high accuracy and precision scores, but our recall is low. As mentioned earlier, recall is essentially a model's ability to predict positives from actual positives, and is good to look at in scenarios where false negatives have a high cost. Since we are using a naive bayes classifier, we make the assumption that our input variables are independent. However, this may not necessarily true. Earlier, we plotted the correlation matrix, and observed high correlations between some variables, which indicates that they might not be independent of each other. This could potentially be problematic, as this isn't representative of real-life scenarios, and could produce incorrect models and predictions.

Part 2: Cross Validation and Model Selection

We've sampled a number of different classification techniques, leveraging nearest neighbors, linear classifiers, and statistical classifiers. You've also tweaked with a few parameters for those models to optimize performance. Based on these experiments you should have settled on a particular model that performs most optimally on this dataset. Before our work is done though, we want to ensure that our results are not the result of the random sampling of our data we did with the train-test split. To check this, we will conduct a K-fold cross validation on our top 2 performing models, assess their cumulative performance across folds (report accuracy, precision, recall, and F1 score), and determine the best model for our particular data.

Question 4.1

Select your top 2 performing models and run a 10-Fold cross validation on both. Report your best performing model. Use the `model_selection.KFold` class from `scikit-learn`.

```
In [35]: kf = KFold(n_splits=10, random_state=1, shuffle=True)
knn = KNeighborsClassifier()
print("KNN:")
sc = cross_validate(knn, df_test, df_sick_test, scoring=("accuracy", "precision", "recall", "f1"), cv=kf, n_jobs=-1)
print("Accuracy:" + str(sum(sc['test_accuracy'])/10))
print("Precision:" + str(sum(sc['test_precision'])/10))
print("Recall:" + str(sum(sc['test_recall'])/10))
print("F1 Score:" + str(sum(sc['test_f1'])/10))

print("\n")

clf_1 = SVC(probability=True)
print("SVM:")
sc = cross_validate(clf_1, df_test, df_sick_test, scoring=("accuracy", "precision", "recall", "f1"), cv=kf, n_jobs=-1)
print("Accuracy:" + str(sum(sc['test_accuracy'])/10))
print("Precision:" + str(sum(sc['test_precision'])/10))
print("Recall:" + str(sum(sc['test_recall'])/10))
print("F1 Score:" + str(sum(sc['test_f1'])/10))

KNN:
Accuracy:0.8238095238095238
Precision:0.8083333333333332
Recall:0.7633333333333334
F1 Score:0.7679365079365079

SVM:
Accuracy:0.8214285714285714
Precision:0.775
Recall:0.7633333333333334
F1 Score:0.7479365079365079
```

Question 4.2

Discuss your results and why they differ slightly from what you got for the 2 models above.

These four metrics differ from the first time we ran the model because we performed cross validation, where we assessed the performance of our model across all validation sets, not just on a single test set of data.

Question 4.3

Out of these 2 models, based on their scores for the 4 metrics, which one would you pick for this specific case of predicting if someone has heart disease or not?

Since the cost of a false negative is high, we should be looking at recall, but both models have the same recall values. In this case, we look at the other scores, and see that the KNN classifier has better scores, thus that is the model I would choose.