


```
[30]: X_val_RESAMP.reset_index(drop=True,inplace=True)
      y_val_RESAMP.reset_index(drop=True,inplace=True)
      X_val_RESAMP_OR.reset_index(drop=True,inplace=True)
      y_val_RESAMP_OR.reset_index(drop=True,inplace=True)

      sample_index1 = sample(X_train_RESAMP.index.tolist(),10000)
      sample_index2 = sample(X_val_RESAMP.index.tolist(),10000)

In [37]: knn1 = KNeighborsClassifier(n_neighbors=6)
          knn1.fit(X_train_RESAMP.iloc[sample_index1], y_train_RESAMP_OR.iloc[sample_index1])

Out[37]: KNeighborsClassifier(n_neighbors=6)

In [38]: preds_RESAMP = knn1.predict(X_val_RESAMP.iloc[sample_index2])

In [39]: from sklearn.metrics import confusion_matrix, classification_report
          print(confusion_matrix(preds_RESAMP, y_val_RESAMP_OR.iloc[sample_index2]))

[[8756  592]
 [ 153 4391]]

In [41]: knn2 = KNeighborsClassifier(n_neighbors = 6)
          knn2.fit(X_train_RESAMP_OR.iloc[sample_index2], y_val_RESAMP_OR.iloc[sample_index2])

Out[41]: KNeighborsClassifier(n_neighbors=6)

In [43]: preds_VAL = knn2.predict(X_val_RESAMP_OR.iloc[sample_index2])

In [44]: print(confusion_matrix(preds_VAL, y_val_RESAMP_OR.iloc[sample_index2]))

[[8756 592]
 [ 153 4391]]
```

Just by comparing the two confusion matrices, we can see that the numbers are better. There is a small increase in accuracy in the model trained on the validation set, but it is low enough to say that the clustering is generalisable.

We can also check the predictions with Cohen's Kappa. We get 64% agreement, which is a relatively decent score.

```
In [57]: cohen_kappa_score(preds_VAL, preds_RESAMP)

Out[57]: 0.6457374613708246
```

Finally we can do the same for the data with outliers removed, and check if there is any improvement in accuracy.

```
In [62]: knn3 = KNeighborsClassifier(n_neighbors = 6)
          knn3.fit(X_train_RESAMP_OR.iloc[sample_index1], y_train_RESAMP_OR.iloc[sample_index1])
          preds_OR = knn2.predict(X_val_RESAMP_OR.iloc[sample_index2])

In [63]: knn4 = KNeighborsClassifier(n_neighbors = 6)
          knn4.fit(X_val_RESAMP_OR.iloc[sample_index2], y_val_RESAMP_OR.iloc[sample_index2])
          preds_VAL_OR = knn4.predict(X_val_RESAMP_OR.iloc[sample_index2])

In [64]: cohen_kappa_score(preds_OR, preds_VAL_OR)

Out[64]: 0.5689034638049317
```

The Kappa score has reduced here - perhaps the outliers were useful.

Decision Trees and Random Forests

A random forest is just a collection of decision trees. The gain in accuracy is offset by a loss of interpretability.

```
In [68]: from sklearn.tree import DecisionTreeClassifier
          from sklearn.ensemble import RandomForestClassifier

In [72]: tree1 = DecisionTreeClassifier()
          tree1.fit(X_train_FULL, y_train_FULL)
          tree2 = DecisionTreeClassifier()
          tree2.fit(X_train_OR, y_train_OR)
          tree3 = DecisionTreeClassifier()
          tree3.fit(X_train_RESAMP, y_train_RESAMP)
          tree4 = DecisionTreeClassifier()
          tree4.fit(X_train_RESAMP_OR, y_train_RESAMP_OR)

Out[72]: DecisionTreeClassifier()

In [73]: pred1 = tree1.predict(X_val_FULL)
          pred2 = tree2.predict(X_val_OR)
          pred3 = tree3.predict(X_val_RESAMP)
          pred4 = tree4.predict(X_val_RESAMP_OR)

In [76]: print('pred1')
          print(classification_report(pred1,y_val_FULL))
          print('\n')
          print('pred2')
          print(classification_report(pred2,y_val_OR))
          print('\n')
          print('pred3')
          print(classification_report(pred3,y_val_RESAMP))
          print('\n')
          print('pred4')
          print(classification_report(pred4,y_val_RESAMP_OR))
          print('\n')
```

pred1	precision	recall	f1-score	support
0	0.98	0.98	0.98	22915
1	0.42	0.36	0.39	891
accuracy				23806
macro avg	0.70	0.67	0.68	23806
weighted avg	0.95	0.96	0.96	23806

pred2	precision	recall	f1-score	support
0	0.99	0.99	0.99	17167
1	0.29	0.25	0.27	159
accuracy				17326
macro avg	0.64	0.62	0.63	17326
weighted avg	0.99	0.99	0.99	17326

pred3	precision	recall	f1-score	support
0	0.95	0.96	0.95	9500
1	0.64	0.59	0.62	1258
accuracy				10758
macro avg	0.79	0.77	0.78	10758
weighted avg	0.91	0.91	0.91	10758

pred4	precision	recall	f1-score	support
0	0.96	0.97	0.97	7088
1	0.79	0.71	0.74	934
accuracy				8022
macro avg	0.87	0.84	0.86	8022
weighted avg	0.94	0.94	0.94	8022

```
In [77]: forest1 = RandomForestClassifier()
          forest1.fit(X_train_FULL, y_train_FULL)
          forest2 = RandomForestClassifier()
          forest2.fit(X_train_OR, y_train_OR)
          forest3 = RandomForestClassifier()
          forest3.fit(X_train_RESAMP, y_train_RESAMP)
          forest4 = RandomForestClassifier()
          forest4.fit(X_train_RESAMP_OR, y_train_RESAMP_OR)

Out[77]: RandomForestClassifier()

In [78]: forestpred1 = forest1.predict(X_val_FULL)
          forestpred2 = forest2.predict(X_val_OR)
          forestpred3 = forest3.predict(X_val_RESAMP)
          forestpred4 = forest4.predict(X_val_RESAMP_OR)

In [79]: print('forestpred1')
          print(classification_report(forestpred1,y_val_FULL))
          print('\n')
          print('forestpred2')
          print(classification_report(forestpred2,y_val_OR))
          print('\n')
          print('forestpred3')
          print(classification_report(forestpred3,y_val_RESAMP))
          print('\n')
          print('forestpred4')
          print(classification_report(forestpred4,y_val_RESAMP_OR))
          print('\n')
```

forestpred1	precision	recall	f1-score	support
0	1.00	0.98	0.99	23559
1	0.28	0.84	0.42	247
accuracy				23806
macro avg	0.64	0.91	0.70	23806
weighted avg	0.39	0.98	0.58	23806

forestpred2	precision	recall	f1-score	support
0	1.00	0.99	1.00	17302
1	0.18	1.00	0.30	24
accuracy				17326
macro avg	0.59	1.00	0.65	17326
weighted avg	1.00	0.99	1.00	17326

forestpred3	precision	recall	f1-score	support
0	0.98	0.96	0.97	9820
1	0.66	0.82	0.73	938
accuracy				10758
macro avg	0.82	0.89	0.85	10758
weighted avg	0.95	0.95	0.95	10758

forestpred4	precision	recall	f1-score	support
0	1.00	0.98	0.99	7308
1	0.81	0.95	0.88	714
accuracy				8022
macro avg	0.90	0.97	0.93	8022
weighted avg	0.98	0.98	0.98	8022

Explanation in document.

Support Vector Classifier

If we have p variables, we can think of this as a p -dimensional space.

In 3 dimensions, a plane is 2 dimensional and flat in the 3 dimensional space. It separates points from one side of the plane, from points on the other.

This idea generalises to n dimensions. A plane in $(n-1)$ dimensions will be flat in n -dimensional space. This plane is known as a hyperplane.

One issue with SVC is that it requires a flat division of the points from different classes. This is not always the case in datasets.

```
In [81]: from sklearn.svm import SVC

In [82]: svc1 = SVC()
          svc1.fit(X_train_FULL, y_train_FULL)
          svc2 = SVC()
          svc2.fit(X_train_OR, y_train_OR)
          svc3 = SVC()
          svc3.fit(X_train_RESAMP, y_train_RESAMP)
          svc4 = SVC()
          svc4.fit(X_train_RESAMP_OR, y_train_RESAMP_OR)

Out[82]: SVC()

In [83]: svc_pred1 = svc1.predict(X_val_FULL)
          svc_pred2 = svc2.predict(X_val_OR)
          svc_pred3 = svc3.predict(X_val_RESAMP)
          svc_pred4 = svc4.predict(X_val_RESAMP_OR)

In [84]: print('svc_pred1')
          print(classification_report(svc_pred1,y_val_FULL))
          print('\n')
          print('svc_pred2')
          print(classification_report(svc_pred2,y_val_OR))
          print('\n')
          print('svc_pred3')
          print(classification_report(svc_pred3,y_val_RESAMP))
          print('\n')
          print('svc_pred4')
          print(classification_report(svc_pred4,y_val_RESAMP_OR))
          print('\n')
```

svc_pred1	precision	recall	f1-score	support
0	1.00	0.97	0.98	23776
1	0.02	0.53	0.04	30
accuracy				23806
macro avg	0.51	0.75	0.61	23806
weighted avg	1.00	0.97	0.98	23806

svc_pred2	precision	recall	f1-score	support
0	1.00	0.99	1.00	17326
1	0.00	0.00	0.00	0
accuracy				17326
macro avg	0.50	0.50	0.50	17326
weighted avg	1.00	0.99	1.00	17326

svc_pred3	precision	recall	f1-score	support
0	0.97	0.94	0.95	9933
1	0.46	0.66	0.54	825
accuracy				10758
macro avg	0.72	0.80	0.75	10758
weighted avg	0.93	0.92	0.92	10758

svc_pred4	precision	recall	f1-score	support
0	0.99	0.93	0.96	7653
1	0.33	0.74	0.45	369
accuracy				8022
macro avg	0.66	0.83	0.70	8022
weighted avg	0.96	0.92	0.93	8022

C:\Users\goggar\anaconda3\lib\site-packages\sklearn\metrics\classification.py:1245: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to control this behavior.
C:\Users\goggar\anaconda3\lib\site-packages\sklearn\metrics\classification.py:1245: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to control this behavior.
C:\Users\goggar\anaconda3\lib\site-packages\sklearn\metrics\classification.py:1245: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to control this behavior.
C:\Users\goggar\anaconda3\lib\site-packages\sklearn\metrics\classification.py:1245: UndefinedMetricWarning: Recall and F-score are ill-defined and being set to 0.0 in labels with no true samples. Use 'zero_division' parameter to control this behavior.

Explanation in document.

Linear & Quadratic Discriminant Analysis

Unlike Support Vector Classifiers, Discriminant Analysis assumes that the data follow a multivariate normal distribution. This is not always true, but LDA and QDA have a good track-record of predictions, even for data which is not actually normally distributed.

Both methods seek to find the boundary between classes, and both use Bayes' Theorem. LDA assumes that the covariance structure is the same for both classes, whereas QDA allows different classes to have different covariance structures. Because of how this works out mathematically, we get a linear boundary for LDA and a curvy boundary for QDA.

```
In [58]: from sklearn.discriminant_analysis import LinearDiscriminantAnalysis, QuadraticDiscriminantAnalysis

In [87]: lda1 = LinearDiscriminantAnalysis()
          lda1.fit(X_train_FULL, y_train_FULL)
          lda2 = LinearDiscriminantAnalysis()
          lda2.fit(X_train_OR, y_train_OR)
          lda3 = LinearDiscriminantAnalysis()
          lda3.fit(X_train_RESAMP, y_train_RESAMP)
          lda4 = LinearDiscriminantAnalysis()
          lda4.fit(X_train_RESAMP_OR, y_train_RESAMP_OR)

Out[87]: LinearDiscriminantAnalysis()

In [88]: lda_pred1 = lda1.predict(X_val_FULL)
          lda_pred2 = lda2.predict(X_val_OR)
          lda_pred3 = lda3.predict(X_val_RESAMP)
          lda_pred4 = lda4.predict(X_val_RESAMP_OR)

In [89]: print('lda_pred1')
          print(classification_report(lda_pred1,y_val_FULL))
          print('\n')
          print('lda_pred2')
          print(classification_report(lda_pred2,y_val_OR))
          print('\n')
          print('lda_pred3')
          print(classification_report(lda_pred3,y_val_RESAMP))
          print('\n')
          print('lda_pred4')
          print(classification_report(lda_pred4,y_val_RESAMP_OR))
          print('\n')
```

lda_pred1	precision	recall	f1-score	support
0	0.97	0.98	0.98	22768
1	0.48	0.35	0.40	1038
accuracy				23806
macro avg	0.73	0.66	0.69	23806
weighted avg	0.95	0.96	0.95	23806

lda_pred2	precision	recall	f1-score	support
0	0.98	1.00	0.99	16907
1	0.46	0.15	0.22	419
accuracy				17326
macro avg	0.72	0.57	0.61	17326
weighted avg	0.97	0.98	0.97	17326

lda_pred3	precision	recall	f1-score	support
0	0.98	0.95	0.96	9888
1	0.54	0.73	0.62	870
accuracy				10758
macro avg	0.76	0.84	0.79	10758
weighted avg	0.94	0.93	0.93	10758

lda_pred4	precision	recall	f1-score	support
0	1.00	0.95	0.98	7520
1	0.59	0.96	0.73	502
accuracy				8022
macro avg	0.79	0.96	0.85	8022
weighted avg	0.97	0.95	0.96	8022

These results are not good, and show that perhaps the assumption that the covariance matrix for both classes is the same, is a bad assumption.

```
In [90]: qda1 = QuadraticDiscriminantAnalysis()
          qda1.fit(X_train_FULL, y_train_FULL)
          qda2 = QuadraticDiscriminantAnalysis()
          qda2.fit(X_train_OR, y_train_OR)
          qda3 = QuadraticDiscriminantAnalysis()
          qda3.fit(X_train_RESAMP, y_train_RESAMP)
          qda4 = QuadraticDiscriminantAnalysis()
          qda4.fit(X_train_RESAMP_OR, y_train_RESAMP_OR)

Out[90]: QuadraticDiscriminantAnalysis()

In [91]: qda_pred1 = qda1.predict(X_val_FULL)
          qda_pred2 = qda2.predict(X_val_OR)
          qda_pred3 = qda3.predict(X_val_RESAMP)
          qda_pred4 = qda4.predict(X_val_RESAMP_OR)

In [92]: print('qda_pred1')
          print(classification_report(qda_pred1,y_val_FULL))
          print('\n')
          print('qda_pred2')
          print(classification_report(qda_pred2,y_val_OR))
          print('\n')
          print('qda_pred3')
          print(classification_report(qda_pred3,y_val_RESAMP))
          print('\n')
          print('qda_pred4')
          print(classification_report(qda_pred4,y_val_RESAMP_OR))
          print('\n')
```

weighted avg	0.82	0.27	0.24	10758
--------------	------	------	------	-------

qda_pred4	precision	recall	f1-score	support
0	0.35	0.98	0.52	2551
1	0.95	0.15	0.25	5471
accuracy			0.41	8022
macro avg	0.65	0.57	0.38	8022
weighted avg	0.76	0.41	0.34	8022

The collinearity has caused a lot of problems with the predictions. Two options would be:

- removing variables which we think are collinear (did this before).
- use variables from Lasso method.


```

In [65]: qda1 = QuadraticDiscriminantAnalysis()
qda1.fit(X_train_FULL[['VOL_DS','VOL_DEP','BETRAG_EINGANG_3M']], y_train_FULL)

qda2 = QuadraticDiscriminantAnalysis()
qda2.fit(X_train_OR[['VOL_DS','VOL_DEP','BETRAG_EINGANG_3M']], y_train_OR)

qda3 = QuadraticDiscriminantAnalysis()
qda3.fit(X_train_RESAMP[['VOL_DS','VOL_DEP','BETRAG_EINGANG_3M']], y_train_RESAMP)

qda4 = QuadraticDiscriminantAnalysis()
qda4.fit(X_train_RESAMP_OR[['VOL_DS','VOL_DEP','BETRAG_EINGANG_3M']], y_train_RESAMP_OR)
  
```



```

Out[65]: QuadraticDiscriminantAnalysis()
  
```

The collinearity has caused a lot of problems with the predictions. Two options would be:

- removing variables which we think are collinear (did this before).
- use variables from Lasso method.

```
In [65]: qda1 = QuadraticDiscriminantAnalysis()
          qda1.fit(X_train_FULL[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']], y_train_FULL)
          qda2 = QuadraticDiscriminantAnalysis()
          qda2.fit(X_train_OR[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']], y_train_OR)
          qda3 = QuadraticDiscriminantAnalysis()
          qda3.fit(X_train_RESAMP[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']], y_train_RESAMP)
          qda4 = QuadraticDiscriminantAnalysis()
          qda4.fit(X_train_RESAMP_OR[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']], y_train_RESAMP_OR)

Out[65]: QuadraticDiscriminantAnalysis()

In [66]: qda_pred5 = qda1.predict(X_val_FULL[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']])
          qda_pred6 = qda2.predict(X_val_OR[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']])
          qda_pred7 = qda3.predict(X_val_RESAMP[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']])
          qda_pred8 = qda4.predict(X_val_RESAMP_OR[['VOL_DS', 'ANZ_DEP', 'BETRAG_EINGANG_3M']])

In [68]: print('qda_pred5')
          print(classification_report(qda_pred5,y_val_FULL))
          print('\n')
          print('qda_pred6')
          print(classification_report(qda_pred6,y_val_OR))
          print('\n')
          print('qda_pred7')
          print(classification_report(qda_pred7,y_val_RESAMP))
          print('\n')
          print('qda_pred8')
          print(classification_report(qda_pred8,y_val_RESAMP_OR))
          print('\n')
```

qda_pred5	precision	recall	f1-score	support
0	0.96	0.98	0.97	22775
1	0.29	0.21	0.24	1031
accuracy				23806
macro avg	0.63	0.59	0.61	23806
weighted avg	0.94	0.94	0.94	23806

qda_pred6	precision	recall	f1-score	support
0	0.94	1.00	0.97	16171
1	0.56	0.06	0.12	1155
accuracy				17326
macro avg	0.75	0.53	0.63	17326
weighted avg	0.91	0.53	0.51	17326

qda_pred7	precision	recall	f1-score	support
0	0.96	0.92	0.94	9935
1	0.38	0.48	0.40	823
accuracy				10758
macro avg	0.65	0.70	0.67	10758
weighted avg	0.91	0.89	0.90	10758

qda_pred8	precision	recall	f1-score	support
0	0.93	0.93	0.93	7198
1	0.42	0.43	0.43	824
accuracy				8022
macro avg	0.68	0.68	0.68	8022
weighted avg	0.88	0.88	0.88	8022

Results are not as good as when we manually selected the variables to use.

```
In [60]: qda5 = QuadraticDiscriminantAnalysis()
          qda5.fit(X_train_FULL[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])
          qda6 = QuadraticDiscriminantAnalysis()
          qda6.fit(X_train_OR[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])
          qda7 = QuadraticDiscriminantAnalysis()
          qda7.fit(X_train_RESAMP[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])
          qda8 = QuadraticDiscriminantAnalysis()
          qda8.fit(X_train_RESAMP_OR[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])

Out[60]: QuadraticDiscriminantAnalysis()

In [61]: qda_pred5 = qda5.predict(X_val_FULL[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])
          qda_pred6 = qda6.predict(X_val_OR[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])
          qda_pred7 = qda7.predict(X_val_RESAMP[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])
          qda_pred8 = qda8.predict(X_val_RESAMP_OR[['VOL_DS', 'ANZ_DEP', 'VOL_DEP', 'KSEIT', 'ALTER', 'BETRAG_EINGANG_3M', 'LOGIN_WEB_3M', 'LOGIN_MX_3M']])

In [63]: print('qda_pred5')
          print(classification_report(qda_pred5,y_val_FULL))
          print('\n')
          print('qda_pred6')
          print(classification_report(qda_pred6,y_val_OR))
          print('\n')
          print('qda_pred7')
          print(classification_report(qda_pred7,y_val_RESAMP))
          print('\n')
          print('qda_pred8')
          print(classification_report(qda_pred8,y_val_RESAMP_OR))
          print('\n')
```

```
print(classification_report(qda_pred5,y_val_FUL5))
print('\n')
print('qda_pred6')
print(classification_report(qda_pred6,y_val_OR))
print('\n')
print('qda_pred7')
print(classification_report(qda_pred7,y_val_RESAMP))
print('\n')
print('qda_pred8')
print(classification_report(qda_pred8,y_val_RESAMP_OR))
print('\n')
```

qda_pred5	precision	recall	f1-score	support
0	0.96	0.98	0.97	22775
1	0.29	0.21	0.24	1031
accuracy			0.94	23806
macro avg	0.63	0.59	0.61	23806
weighted avg	0.94	0.94	0.94	23806

qda_pred6	precision	recall	f1-score	support
0	0.94	1.00	0.97	16171
1	0.56	0.06	0.12	1155
accuracy			0.93	17326
macro avg	0.75	0.53	0.64	17326
weighted avg	0.91	0.53	0.51	17326

qda_pred7	precision	recall	f1-score	support
0	0.96	0.92	0.94	9935
1	0.38	0.48	0.40	823
accuracy			0.92	10758
macro avg	0.65	0.70	0.67	10758
weighted avg	0.91	0.89	0.90	10758

qda_pred8	precision	recall	f1-score	support
0				