

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5 import plotly.express as px
6
7 from sklearn.preprocessing import LabelEncoder
8 from sklearn.model_selection import cross_val_score
9
10 # Supervised learning
11 from sklearn.neural_network import MLPClassifier
12 from sklearn.neighbors import KNeighborsClassifier
13 from sklearn.svm import SVC
14 from sklearn.gaussian_process import GaussianProcessClassifier
15 from sklearn.gaussian_process.kernels import RBF
16 from sklearn.tree import DecisionTreeClassifier
17 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
18 from sklearn.naive_bayes import GaussianNB
19 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
20 from sklearn.model_selection import KFold, train_test_split
21 import statsmodels.formula.api as smf
22 import statsmodels.api as sm
23
24 # Unsupervised learning
25 from sklearn.decomposition import PCA, FastICA
26 from sklearn.manifold import TSNE
27 from sklearn.cluster import KMeans
28 from sklearn.preprocessing import StandardScaler, Normalizer
29 from scipy.optimize import linear_sum_assignment
```

▼ Project

Madeleine Ekblom, Matias Jääskeläinen, Jakub Kubečka

Time used: many hours

Read in the data

```
1 dummy = pd.read_csv("/content/sample_data/dummy.csv")
2 npf_train = pd.read_csv("/content/sample_data/npf_train.csv")
3 npf_test_hidden = pd.read_csv("/content/sample_data/npf_test_hidden.csv")
```

See the columns

```
1 npf_train.head(2)
```



	id	date	event	partlybad	HYY_META.CO2168.mean	HYY_META.CO2168.std	HYY_
0	1	2000-01-23	nonevent	False	373.496585	0.189497	

1	2	2000-01-25	nonevent	False	381.752738	1.701439	
---	---	------------	----------	-------	------------	----------	--

```
1 npf_test_hidden.head(2)
```

	id	date	event	partlybad	HYY_META.CO2168.mean	HYY_META.CO2168.std	HYY_
0	725	NaN	NaN	False	372.575187	10.051405	
1	726	NaN	NaN	False	382.408306	0.752684	

2 rows x 104 columns

Split the data into features and labels

```
1 #del x
2 X = npf_train.iloc[:,4::2].copy()
3 #X = npf_train.iloc[:,4:]
4
5 #del test
6 test = npf_test_hidden.iloc[:,4::2].copy()
7 #test = npf_test_hidden.iloc[:,4:]
8
9 #del y
10 Y = npf_train.iloc[:,2]
```

```
1 X.head(1)
```

	HYY_META.CO2168.mean	HYY_META.CO2336.mean	HYY_META.CO242.mean	HYY_META.C
0	373.496585	373.382593	373.961481	

```
1 Y[:2]
```

```
0 nonevent
1 nonevent
Name: event, dtype: object
```

```
1 test.head(1)
```

	HYY_META.CO2168.mean	HYY_META.CO2336.mean	HYY_META.CO242.mean	HYY_META.C
0	372.575187	372.324439	376.135989	

Make Y (labels) categorical

```

1  labeler = LabelEncoder()
2  Y = labeler.fit_transform(Y.astype('str'))
3  list(labeler.classes_), set(Y)

↳ (['II', 'Ia', 'Ib', 'nonevent'], {0, 1, 2, 3})

```

```

1  Y[:2]

↳ array([3, 3])

```

Scale X (features)

```

1  scaler = StandardScaler()
2  X = scaler.fit(X).transform(X)
3
4  test = scaler.transform(test)

1  X.shape, test.shape

↳ ((724, 50), (724, 50))

```

Reduce the number of dimensions

First test how many PCA-components should be used

```

1  pca_test = PCA(n_components=30)
2  pca_test.fit(X)

↳ PCA(copy=True, iterated_power='auto', n_components=30, random_state=None,
      svd_solver='auto', tol=0.0, whiten=False)

1  np.sum(pca_test.explained_variance_ratio_)

↳ 0.9999269315979487

1  plt.plot(np.arange(1,31),pca_test.explained_variance_ratio_,'o-');
2  plt.plot(np.arange(1,31),np.cumsum(pca_test.explained_variance_ratio_), 'o-');

↳

```

Transform the labels and the test set into 4 PCs:

```

1  pca = PCA(n_components=4)
2  npf_pca = pca.fit_transform(X)
3
4  test_pca = pca.transform(test)

1  npf_pca.shape, type(npf_pca), test_pca.shape, type(test_pca)

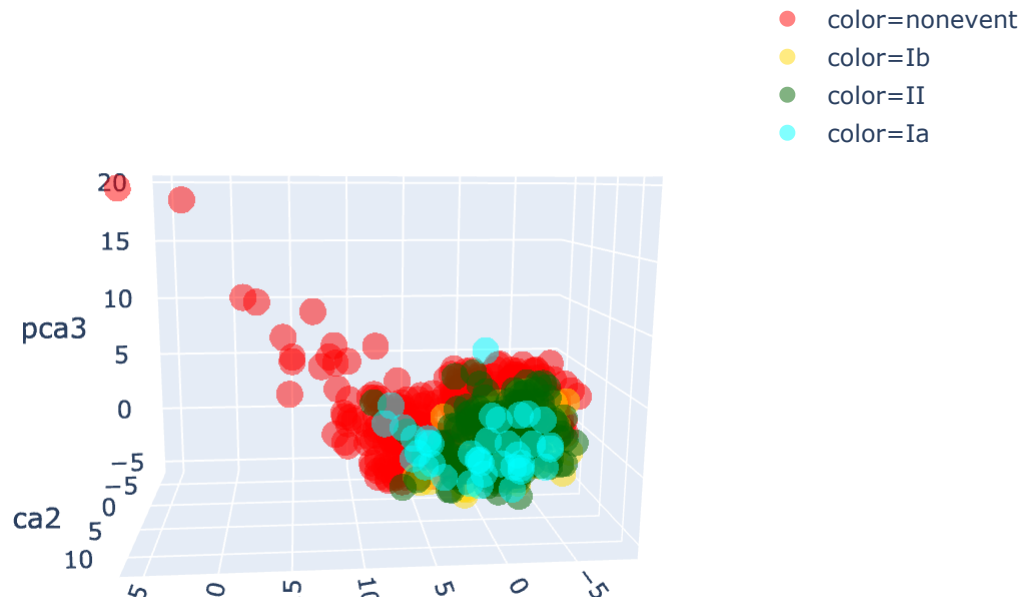
↳ ((724, 4), numpy.ndarray, (724, 4), numpy.ndarray)

1  own_cmap = [
2      "red",
3      "gold",
4      "darkgreen",
5      "cyan",
6      "dodgerblue",
7      "royalblue",
8      "blue",
9      "blueviolet",
10     "plum",
11     "deeppink",
12     "magenta",
13     "purple",
14     "saddlebrown",
15     "lightsalmon",
16     "k"]

1  npf_pca = pd.DataFrame(npf_pca, columns=["pca1", "pca2", "pca3", "pca4"])
2
3  test_pca = pd.DataFrame(test_pca, columns=["pca1", "pca2", "pca3", "pca4"])
4
5  klabels = labeler.inverse_transform(Y)
6  px.scatter_3d(data_frame=npf_pca,
7                x="pca1",
8                y="pca2",
9                z="pca3",
10               color=klabels,
11               color_discrete_sequence=own_cmap,
12               hover_name=npf_pca.index,
13               opacity=0.5,
14               width= 600,
15               height= 500)

```





Test ICA also

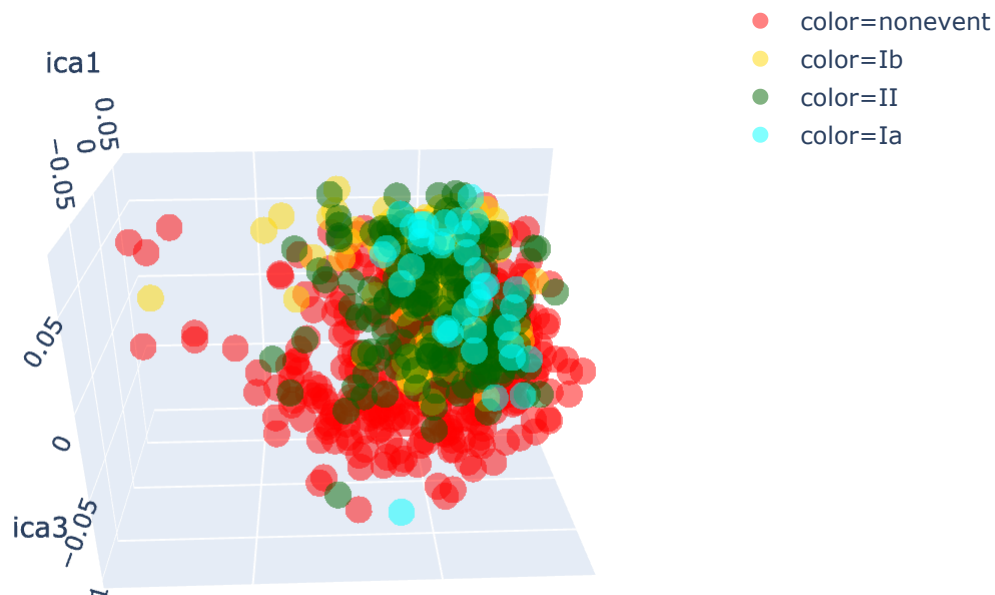
```
1 ica = FastICA(n_components=4)
2 npf_ica = ica.fit_transform(X)
3
4 test_ica = ica.transform(test)
```

```
1 npf_ica.shape, type(npf_ica), test_ica.shape, type(test_ica)
```

```
↳ ((724, 4), numpy.ndarray, (724, 4), numpy.ndarray)
```

```
1 npf_ica = pd.DataFrame(npf_ica, columns=["ica1", "ica2", "ica3", "ica4"])
2
3 test_ica = pd.DataFrame(test_ica, columns=["ica1", "ica2", "ica3", "ica4"])
4
5 px.scatter_3d(data_frame=npf_ica,
6               x="ica1",
7               y="ica2",
8               z="ica3",
9               color=klabels,
10              color_discrete_sequence=own_cmap,
11              hover_name=npf_ica.index,
12              opacity=0.5,
13              width= 600,
14              height= 500)
```

```
↳
```



Try t-SNE for visualisation

It does maybe make sense but it looks cool in the CV

```

/\  /\
((ovo))
():::()
vvv

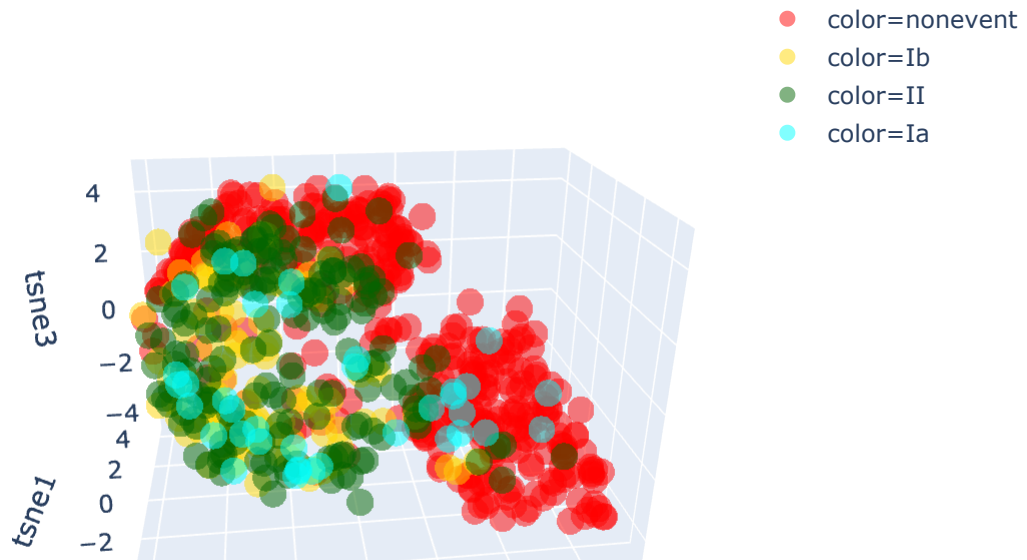
```

```

1  npf_tsne = TSNE(n_components=3,perplexity=100,random_state=1).fit_transform(X)
2
3  npf_tsne = pd.DataFrame(npf_tsne,columns=['tsne1','tsne2','tsne3'])
4
5  px.scatter_3d(data_frame=npf_tsne,
6                x="tsne1",
7                y="tsne2",
8                z="tsne3",
9                color=klabels,
10                 color_discrete_sequence=own_cmap,
11                 hover_name=npf_tsne.index,
12                 opacity=0.5,
13                 width= 600,
14                 height= 500)

```





Split the train data into train and validation

`x_train` = training features

`x_test` = validation features

`y_train` = training labels

`y_test` = validation labels

```
1  X_train,X_test,Y_train,Y_test = train_test_split(npf_ica,Y,random_state=42)
```

```
1  X_train.shape,X_test.shape,Y_train.shape,Y_test.shape
```

```
↳ ((543, 4), (181, 4), (543,), (181,))
```

▼ Try different classification models

```
1  names = ["Nearest Neighbors", "Linear SVM", "RBF SVM", "Gaussian Process",
2          "Decision Tree", "Random Forest", "Neural Net", "AdaBoost",
3          "Naive Bayes", "QDA"]
4
5  classifiers = [
6      KNeighborsClassifier(3),
7      SVC(kernel="linear", C=0.025),
8      SVC(gamma=2, C=1),
9      GaussianProcessClassifier(1.0 * RBF(1.0)),
10     DecisionTreeClassifier(max_depth=5),
11     RandomForestClassifier(max_depth=5, n_estimators=10, max_features=1),
12     MLPClassifier(alpha=1, max_iter=1000),
13     AdaBoostClassifier(),
14     GaussianNB(),
15     QuadraticDiscriminantAnalysis()]
```

Use cross-validation to find the best model

```

1  # iterate over classifiers
2  for name, clf in zip(names, classifiers):
3      clf.fit(X_train, Y_train)
4      score = clf.score(X_test, Y_test)
5      print("%s:  %0.2f  %" % (name, score*100))

```

```

☞ Nearest Neighbors:  64.64 %
   Linear SVM:  55.25 %
   RBF SVM:  59.12 %
   Gaussian Process:  65.75 %
   Decision Tree:  61.33 %
   Random Forest:  64.09 %
   Neural Net:  61.33 %
   AdaBoost:  58.56 %
   Naive Bayes:  60.77 %
   QDA:  63.54 %

```

Use 10-fold cross-validation to find the best model

```

1  #cross validation
2  # DECLARATION #
3  kf = KFold(n_splits=10)
4  # CODE #
5  scores=[]
6  # iterate over k-folds
7  for train, val in kf.split(npf_pca):
8      scores_iter=[]
9      # iterate over classifiers
10     for name, clf in zip(names, classifiers):
11         clf.fit(npf_pca.iloc[train], Y[train])
12         score = clf.score(npf_pca.iloc[val], Y[val])
13         scores_iter.append(score)
14         #print("{}:  {}".format(name, score))
15     scores.append(scores_iter)
16 #print(scores)
17 scores=np.mean(scores,axis=0)
18 #print(scores)
19 for name, score in zip(names, scores):
20     print("%s:  %0.2f  %" % (name, score*100))

```

```

☞ Nearest Neighbors:  62.87 %
   Linear SVM:  61.08 %
   RBF SVM:  57.52 %
   Gaussian Process:  68.41 %
   Decision Tree:  61.92 %
   Random Forest:  64.82 %
   Neural Net:  67.58 %
   AdaBoost:  59.70 %
   Naive Bayes:  64.81 %
   QDA:  65.21 %

```


Train the model

We pick Gaussian processes

```
1 model = GaussianProcessClassifier(1.0 * RBF(1.0))
2 # Wanted to test also the neural net
3 #model = MLPClassifier(alpha=1, max_iter=1000)
```

Fit the model

```
1 model.fit(npf_ica,Y)
```

```
➞ GaussianProcessClassifier(copy_X_train=True, kernel=1**2 * RBF(length_scale=1)
    max_iter_predict=100, multi_class='one_vs_rest',
    n_jobs=None, n_restarts_optimizer=0,
    optimizer='fmin_l_bfgs_b', random_state=None,
    warm_start=False)
```

Predict values for npf_test_hidden.csv

Using ICA transformed data

```
1 prediction = model.predict(test_ica)
2 #prediction = model.predict(npf_ica)      # this is a test
3 prediction = labeler.inverse_transform(prediction)
4
5 probs = model.predict_proba(test_ica)
6 #probs = model.predict_proba(npf_ica)    # this is a test
7 #probs = np.max(probs, axis=1)
```

```
1 # sum probabilities of event = p(1a) + p(1b) + p(2), p(nonevent) is 3rd column
2 probs = np.sum(probs[:, :3], axis=1)
```

```
1 results_df = pd.DataFrame(zip(prediction, probs), columns=["event", "p_event"])
```

```
1 results_df.head(10)
```

➞

	event	p_event
0	nonevent	0.467028

Print out a csv

```
1 results_df.to_csv("/content/drive/My Drive/Colab Notebooks/answers.csv",index=
3 nonevent 0.334076
1
5 nonevent 0.153963
6 nonevent 0.166439
7 II 0.946102
8 nonevent 0.058983
9 II 0.710914
```