

ML Classification Template

Feature Scaling are applied to improve the overall result even though some algorithms DO NOT need it.

```
In [1]: 1 import numpy as np
2 import matplotlib.pyplot as plt
3 import pandas as pd
```

```
In [2]: 1 df = pd.read_csv("Data.csv")
2 df.head()
```

Out[2]:

	Sample code number	Clump Thickness	Uniformity of Cell Size	Uniformity of Cell Shape	Marginal Adhesion	Single Epithelial Cell Size	Bare Nuclei	Bland Chromatin	Normal Nucleoli	Mitoses	Class
0	1000025	5	1	1	1	2	1	3	1	1	2
1	1002945	5	4	4	5	7	10	3	2	1	2
2	1015425	3	1	1	1	2	2	3	1	1	2
3	1016277	6	8	8	1	3	4	3	7	1	2
4	1017023	4	1	1	3	2	1	3	1	1	2

```
In [3]: 1 df.isnull().sum()
```

```
Out[3]: Sample code number      0
Clump Thickness               0
Uniformity of Cell Size       0
Uniformity of Cell Shape       0
Marginal Adhesion              0
Single Epithelial Cell Size    0
Bare Nuclei                     0
Bland Chromatin                  0
Normal Nucleoli                  0
Mitoses                         0
Class                           0
dtype: int64
```

```
In [4]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

Decision Trees

```
In [5]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

```
In [6]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
```

```
9 from sklearn.tree import DecisionTreeClassifier
10 classifier = DecisionTreeClassifier(criterion = 'entropy', random_state = 0)
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))
```

	precision	recall	f1-score	support
2	0.97	0.96	0.97	107
4	0.94	0.95	0.95	64
accuracy			0.96	171
macro avg	0.96	0.96	0.96	171
weighted avg	0.96	0.96	0.96	171
[103 4]				
[3 61]				
0.9590643274853801				

K Nearest Neighbors

```
In [7]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]
```

```
In [8]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
```

```

5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from sklearn.neighbors import KNeighborsClassifier
10 classifier = KNeighborsClassifier(n_neighbors = 5, metric = 'minkowski', p = 2)
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

      precision    recall  f1-score   support

           2       0.95     0.96     0.96      107
           4       0.94     0.92     0.93      64

    accuracy                           0.95      171
   macro avg       0.95     0.94     0.94      171
weighted avg       0.95     0.95     0.95      171

[[103  4]
 [ 5 59]]
0.9473684210526315

```

Support Vector Machine (kernel='linear')

```

In [9]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [10]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from sklearn.svm import SVC
10 classifier = SVC(kernel = 'linear', random_state = 0)
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

      precision    recall  f1-score   support

           2       0.95     0.95     0.95      107
           4       0.92     0.92     0.92      64

    accuracy                           0.94      171
   macro avg       0.94     0.94     0.94      171
weighted avg       0.94     0.94     0.94      171

[[102  5]
 [ 5 59]]
0.9415204678362573

```

Support Vector Machine (Kernel="rbf")

```

In [11]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [12]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from sklearn.svm import SVC
10 classifier = SVC(kernel = 'rbf', random_state = 0)
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

      precision    recall  f1-score   support

           2       0.97     0.95     0.96      107
           4       0.92     0.95     0.94      64

    accuracy                           0.95      171
   macro avg       0.95     0.95     0.95      171
weighted avg       0.95     0.95     0.95      171

[[102  5]
 [ 5 59]]

```

```
l > 0.9532163742690059
```

Logistic Regression

```
In [13]: 1 x = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [14]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from sklearn.linear_model import LogisticRegression
10 classifier = LogisticRegression(random_state = 0)
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

      precision    recall  f1-score   support

          2       0.95     0.96     0.96      107
          4       0.94     0.92     0.93       64

   accuracy                           0.95      171
macro avg       0.95     0.94     0.94      171
weighted avg    0.95     0.95     0.95      171

[[103  4]
 [ 5 59]]
0.9473684210526315
```

Naive Bayes (GaussianNB)

```
In [15]: 1 x = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [16]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from sklearn.naive_bayes import GaussianNB
10 classifier = GaussianNB()
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

      precision    recall  f1-score   support

          2       0.98     0.93     0.95      107
          4       0.89     0.97     0.93       64

   accuracy                           0.94      171
macro avg       0.93     0.95     0.94      171
weighted avg    0.94     0.94     0.94      171

[[99  8]
 [ 2 62]]
0.9415204678362573
```

Naive Bayes (MultinomialNB)

```
In [17]: 1 x = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [18]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from sklearn.naive_bayes import MultinomialNB
10 classifier = MultinomialNB()
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
```

```

15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

-----
ValueError                                Traceback (most recent call last)
<ipython-input-18-6c27a474db04> in <module>
      9 from sklearn.naive_bayes import MultinomialNB
     10 classifier = MultinomialNB()
--> 11 classifier.fit(X_train, y_train)
     12
     13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report

~/anaconda3/lib/site-packages\sklearn\naive_bayes.py in fit(self, X, y, sample_weight)
  636
  637         self._init_counters(n_effective_classes, n_features)
--> 638         self._count(X, Y)
  639         alpha = self._check_alpha()
  640         self._update_feature_log_prob(alpha)

~/anaconda3/lib/site-packages\sklearn\naive_bayes.py in _count(self, X, Y)
  769     def _count(self, X, Y):
  770         """Count and smooth feature occurrences."""
--> 771         check_non_negative(X, "MultinomialNB (input X)")
  772         self.feature_count_ += safe_sparse_dot(Y.T, X)
  773         self.class_count_ += Y.sum(axis=0)

~/anaconda3/lib/site-packages\sklearn\utils\validation.py in check_non_negative(X, whom)
 1123
 1124     if X_min < 0:
-> 1125         raise ValueError("Negative values in data passed to %s" % whom)
 1126
 1127

ValueError: Negative values in data passed to MultinomialNB (input X)

```

▼ Stochastic Gradient Descent

```

In [19]: 1 x = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [20]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from sklearn.linear_model import SGDClassifier
10 # classifier = SGDClassifier(loss="hinge", penalty="L2", max_iter=5)
11 classifier = SGDClassifier()
12 classifier.fit(X_train, y_train)
13
14 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
15 y_pred = classifier.predict(X_test)
16 print(classification_report(y_test, y_pred))
17 print(confusion_matrix(y_test, y_pred))
18 print(accuracy_score(y_test, y_pred))

          precision    recall  f1-score   support

           2       0.95      0.97      0.96      107
           4       0.95      0.91      0.93       64

      accuracy                           0.95      171
     macro avg       0.95      0.94      0.94      171
weighted avg       0.95      0.95      0.95      171

[[104   3]
 [ 6  58]]
0.9473684210526315

```

▼ XGBoost

```

In [21]: 1 x = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [22]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from xgboost import XGBClassifier
10 classifier = XGBClassifier()
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))


```

```
[12:32:31] WARNING: C:/Users/Administrator/workspace/xgboost-win64_release_1.4.0/src/learner.cc:1095: Starting in XGBoost 1.3.0, the default evaluation metric used with the objective 'binary:logistic' was changed from 'error' to 'logloss'. Explicitly set eval_metric if you'd like to restore the old behavior.
      precision    recall   f1-score   support
      2          0.95     0.96     0.96      107
      4          0.94     0.92     0.93       64

      accuracy          0.95      171
macro avg          0.95     0.94     0.94      171
weighted avg       0.95     0.95     0.95      171

[[103  4]
 [ 5 59]]
0.9473684210526315
```

C:\Users\Perry\anaconda3\lib\site-packages\xgboost\sklearn.py:1146: UserWarning: The use of label encoder in XGBClassifier is deprecated and will be removed in a future release. To remove this warning, do the following: 1) Pass option use_label_encoder=False when constructing XGBClassifier object; and 2) Encode your labels (y) as integers starting with 0, i.e. 0, 1, 2, ..., [num_class - 1].
 warnings.warn(label_encoder_deprecation_msg, UserWarning)

▼ Catboost (Tune by itself and no need to parameter tuning, good results when having multiple categorical variables)

```
In [23]: 1 # pip install catboost

In [24]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [25]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from catboost import CatBoostClassifier
10 classifier = CatBoostClassifier()
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

992: learn: 0.0097515      total: 1.84s  remaining: 12.9ms
993: learn: 0.0097357      total: 1.84s  remaining: 11.1ms
994: learn: 0.0097215      total: 1.84s  remaining: 9.24ms
995: learn: 0.0096939      total: 1.84s  remaining: 7.39ms
996: learn: 0.0096689      total: 1.84s  remaining: 5.54ms
997: learn: 0.0096557      total: 1.84s  remaining: 3.69ms
998: learn: 0.0096378      total: 1.84s  remaining: 1.85ms
999: learn: 0.0096212      total: 1.85s  remaining: 0us
      precision    recall   f1-score   support
      2          0.95     0.96     0.96      107
      4          0.94     0.92     0.93       64

      accuracy          0.95      171
macro avg          0.95     0.94     0.94      171
weighted avg       0.95     0.95     0.95      171

[[103  4]
 [ 5 59]]
0.9473684210526315
```

▼ LightGBM (Light Gradient Boosting Machine)

```
In [26]: 1 # pip install Lightgbm

In [27]: 1 X = df.iloc[:, :-1]
2 y = df.iloc[:, -1]

In [28]: 1 from sklearn.model_selection import train_test_split
2 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)
3
4 from sklearn.preprocessing import StandardScaler
5 sc = StandardScaler()
6 X_train = sc.fit_transform(X_train)
7 X_test = sc.transform(X_test)
8
9 from lightgbm import LGBMClassifier
10 classifier = LGBMClassifier()
11 classifier.fit(X_train, y_train)
12
13 from sklearn.metrics import confusion_matrix, accuracy_score, classification_report
14 y_pred = classifier.predict(X_test)
15 print(classification_report(y_test, y_pred))
16 print(confusion_matrix(y_test, y_pred))
17 print(accuracy_score(y_test, y_pred))

precision    recall   f1-score   support
```

	2	0.96	0.96	0.96	107
	4	0.94	0.94	0.94	64
accuracy				0.95	171
macro avg		0.95	0.95	0.95	171
weighted avg		0.95	0.95	0.95	171

```
[[103  4]
 [ 4 60]]
0.9532163742690059
```

```
In [29]: █ 1  #!/usr/bin/python
2  # -*- coding: utf-8 -*-
3
4  """
5  =====
6  Classifier comparison
7  =====
8
9  A comparison of several classifiers in scikit-learn on synthetic datasets.
10 The point of this example is to illustrate the nature of decision boundaries
11 of different classifiers.
12 This should be taken with a grain of salt, as the intuition conveyed by
13 these examples does not necessarily carry over to real datasets.
14
15 Particularly in high-dimensional spaces, data can more easily be separated
16 linearly and the simplicity of classifiers such as naive Bayes and linear SVMs
17 might lead to better generalization than is achieved by other classifiers.
18
19 The plots show training points in solid colors and testing points
20 semi-transparent. The lower right shows the classification accuracy on the test
21 set.
22 """
23 print(__doc__)
24
25
26 # Code source: Gaël Varoquaux
27 #             Andreas Müller
28 # Modified for documentation by Jaques Grobler
29 # License: BSD 3 clause
30
31 import numpy as np
32 import matplotlib.pyplot as plt
33 from matplotlib.colors import ListedColormap
34 from sklearn.model_selection import train_test_split
35 from sklearn.preprocessing import StandardScaler
36 from sklearn.datasets import make_moons, make_circles, make_classification
37 from sklearn.neural_network import MLPClassifier
38 from sklearn.neighbors import KNeighborsClassifier
39 from sklearn.svm import SVC
40 from sklearn.gaussian_process import GaussianProcessClassifier
41 from sklearn.gaussian_process.kernels import RBF
42 from sklearn.tree import DecisionTreeClassifier
43 from sklearn.ensemble import RandomForestClassifier, AdaBoostClassifier
44 from sklearn.naive_bayes import GaussianNB
45 from sklearn.discriminant_analysis import QuadraticDiscriminantAnalysis
46 from tqdm.notebook import tqdm
47
48 h = .02 # step size in the mesh
49
50 names = ["Nearest Neighbors", "Linear SVM", "RBF SVM", "Gaussian Process",
51          "Decision Tree", "Random Forest", "Neural Net", "AdaBoost",
52          "Naive Bayes", "QDA"]
53
54 classifiers = [
55     KNeighborsClassifier(3),
56     SVC(kernel="linear", C=0.025),
57     SVC(gamma=2, C=1),
58     GaussianProcessClassifier(1.0 * RBF(1.0)),
59     DecisionTreeClassifier(max_depth=5),
60     RandomForestClassifier(max_depth=5, n_estimators=10, max_features=1),
61     MLPClassifier(alpha=1, max_iter=1000),
62     AdaBoostClassifier(),
63     GaussianNB(),
64     QuadraticDiscriminantAnalysis()]
65
66 X, y = make_classification(n_features=2, n_redundant=0, n_informative=2,
67                            random_state=1, n_clusters_per_class=1)
68 rng = np.random.RandomState(2)
69 X += 2 * rng.uniform(size=X.shape)
70 linearly_separable = (X, y)
71
72 datasets = [make_moons(noise=0.3, random_state=0),
73              make_circles(noise=0.2, factor=0.5, random_state=1),
74              linearly_separable
75          ]
76
77 figure = plt.figure(figsize=(27, 9))
78 i = 1
79 # iterate over datasets
80 for ds_cnt, ds in tqdm(enumerate(datasets)):
81     # preprocess dataset, split into training and test part
82     X, y = ds
83     X = StandardScaler().fit_transform(X)
84     X_train, X_test, y_train, y_test = \
85         train_test_split(X, y, test_size=.4, random_state=42)
86
87     x_min, x_max = X[:, 0].min() - .5, X[:, 0].max() + .5
88     y_min, y_max = X[:, 1].min() - .5, X[:, 1].max() + .5
89     xx, yy = np.meshgrid(np.arange(x_min, x_max, h),
```

```

98         np.arange(y_min, y_max, n))
99
100    # just plot the dataset first
101    cm = plt.cm.RdBu
102    cm_bright = ListedColormap(['#FF0000', '#0000FF'])
103    ax = plt.subplot(len(datasets), len(classifiers) + 1, i)
104    if ds_cnt == 0:
105        ax.set_title("Input data")
106    # Plot the training points
107    ax.scatter(X_train[:, 0], X_train[:, 1], c=y_train, cmap=cm_bright,
108                edgecolors='k')
109    # Plot the testing points
110    ax.scatter(X_test[:, 0], X_test[:, 1], c=y_test, cmap=cm_bright, alpha=0.6,
111                edgecolors='k')
112    ax.set_xlim(xx.min(), xx.max())
113    ax.set_ylim(yy.min(), yy.max())
114    ax.set_xticks(())
115    ax.set_yticks(())
116    i += 1
117
118    # iterate over classifiers
119    for name, clf in zip(names, classifiers):
120        ax = plt.subplot(len(datasets), len(classifiers) + 1, i)
121        clf.fit(X_train, y_train)
122        score = clf.score(X_test, y_test)
123
124        # Plot the decision boundary. For that, we will assign a color to each
125        # point in the mesh [x_min, x_max]x[y_min, y_max].
126        if hasattr(clf, "decision_function"):
127            Z = clf.decision_function(np.c_[xx.ravel(), yy.ravel()])
128        else:
129            Z = clf.predict_proba(np.c_[xx.ravel(), yy.ravel()])[:, 1]
130
131        # Put the result into a color plot
132        Z = Z.reshape(xx.shape)
133        ax.contourf(xx, yy, Z, cmap=cm, alpha=.8)
134
135        # Plot the training points
136        ax.scatter(X_train[:, 0], X_train[:, 1], c=y_train, cmap=cm_bright,
137                    edgecolors='k')
138        # Plot the testing points
139        ax.scatter(X_test[:, 0], X_test[:, 1], c=y_test, cmap=cm_bright,
140                    edgecolors='k', alpha=0.6)
141
142        ax.set_xlim(xx.min(), xx.max())
143        ax.set_ylim(yy.min(), yy.max())
144        ax.set_xticks(())
145        ax.set_yticks(())
146        if ds_cnt == 0:
147            ax.set_title(name)
148            ax.text(xx.max() - .3, yy.min() + .3, ('%.2f' % score).lstrip('0'),
149                    size=15, horizontalalignment='right')
150        i += 1
151
152    plt.tight_layout()
153    # plt.show()
154
```

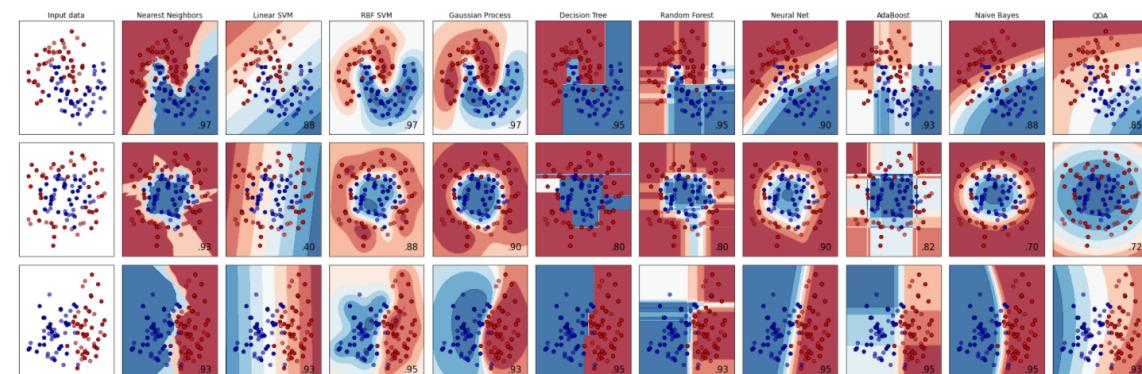
=====
Classifier comparison
=====

A comparison of a several classifiers in scikit-learn on synthetic datasets.
The point of this example is to illustrate the nature of decision boundaries
of different classifiers.
This should be taken with a grain of salt, as the intuition conveyed by
these examples does not necessarily carry over to real datasets.

Particularly in high-dimensional spaces, data can more easily be separated
linearly and the simplicity of classifiers such as naive Bayes and linear SVMs
might lead to better generalization than is achieved by other classifiers.

The plots show training points in solid colors and testing points
semi-transparent. The lower right shows the classification accuracy on the test
set.

3/? [00:05<00:00, 1.78s/it]



In [301]: 1 %%time

```
2 import time
3 from tqdm import tqdm, trange
4
5 for i in tqdm(range(3)):
6     time.sleep(1)
7
8 for i in trange(3):
9     time.sleep(1)

100%|██████████| 3/3 [00:03<00:00,  1.01s/it]
100%|██████████| 3/3 [00:03<00:00,  1.01s/it]

Wall time: 6.08 s
```