

scikit-learn & Bokeh

David Dobrovolný

scikit-learn

- <https://scikit-learn.org>
- built on NumPy, SciPy and matplotlib
- open source
- on-going development (6 updates in 2020, as of May 4th 2021, 2 updates in 2021)
- used by e.g. JPMorgan, Spotify, Télécom ParisTech, Booking.com

scikit-learn

- classification
- regression
- clustering
- dimensionality reduction
- model selection
- preprocessing

Naive Bayes

```
gnb = GaussianNB()  
model = gnb.fit(X_train, y_train)  
model.predict(X_test)
```

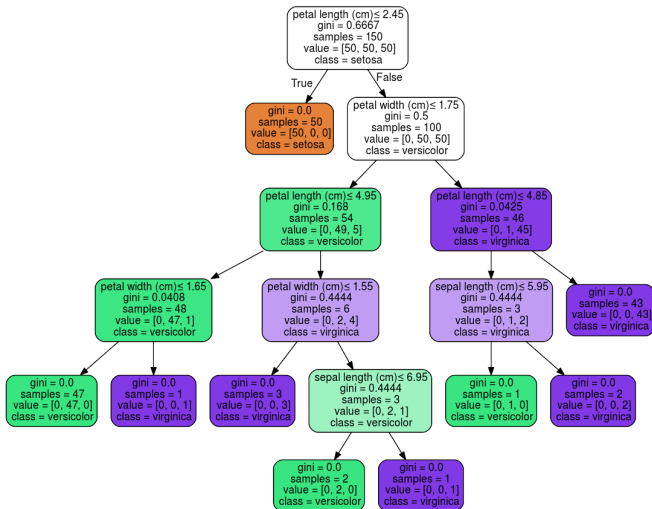
- Gaussian
- Multinomial
- Bernoulli
- ...

Decision Trees

```
clf = tree.DecisionTreeClassifier()  
clf = clf.fit(X, Y)  
clf.predict(test)  
clf.predict_proba(test)
```

- ID3, C4.5, C5.0 and CART
- classification and regression
- Maximum depth, minimum leaf samples, impurity decrease, ...
- Can export to Graphviz format (next slide)

Decision Trees



Support Vector Machines

```
clf = svm.SVC()
```

```
clf.fit(X, y)
```

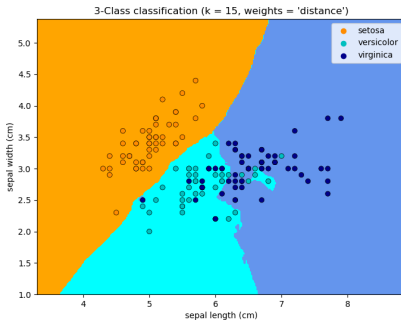
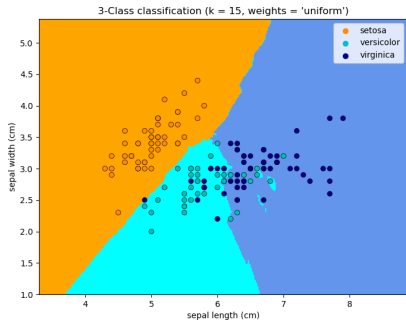
```
clf.predict(test)
```

- classification and regression
- Different kernel functions
 - Linear: $\langle x, x' \rangle$
 - Polynomial: $(\gamma \langle x, x' \rangle + r)^d$
 - RBF: $\exp(-\gamma \|x - x'\|^2)$
 - Sigmoid: $\tanh(\gamma \langle x, x' \rangle + r)$
 - custom

Nearest Neighbours

- unsupervised
 - BallTree
 - KDTree
 - brute force
 - auto (algorithm determines the best approach)
- classification
 - uniform weights
 - weights based on distance (next slide)
- regression

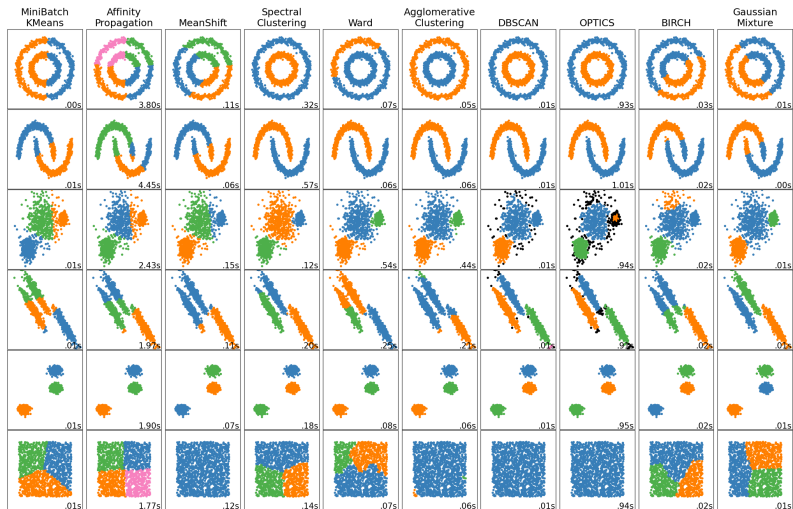
Nearest Neighbours



Other

- Clustering (next slide)
- Ensemble methods
 - Random Forests
 - Adaboost
- Semi-supervised learning
 - Self Training
- Neural Networks

Clustering



Bokeh

- <https://bokeh.org>
- Interactive web browser visualizations.
 - Server App - allows more interactive manipulation
 - Notebook
 - Standalone - limited interactivity, produces html file
 - Examples: <https://docs.bokeh.org/en/latest/docs/gallery.html>