

Theoretical and practical metagenomic approaches to viral discovery

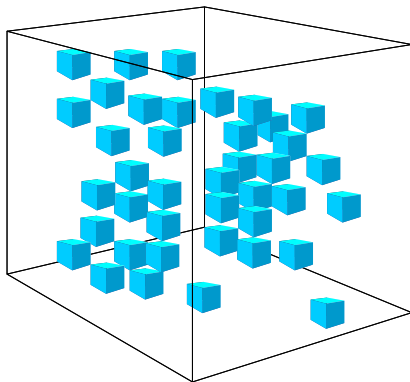
Practical Session: Dimension Reduction and Clustering

Kevin Lamkiewicz, Manja Marz

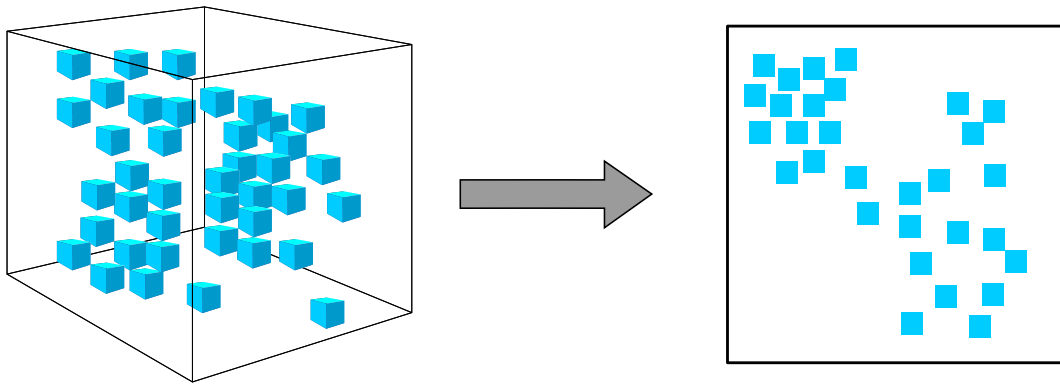
23.10.2019

European Virus Bioinformatics Center

Dimension Reduction



Dimension Reduction



PCA with scikit-learn

A QUICK EXAMPLE

```
1 from sklearn import datasets
2 from sklearn.decomposition import PCA
3 from sklearn.preprocessing import StandardScaler
4 import numpy as np
5 import pandas as pd
6 import matplotlib.pyplot as plt
7
8 iris = datasets.load_iris()
9 # combine everything into a pandas DataFrame.
10 # numpy.c_ concatenates the given vector into one (note the [] instead of ())
11 df = pd.DataFrame(data= np.c_[iris['data'], iris['target_names'][iris['target']]],
12                   columns= iris['feature_names'] + ['target'])
```

NORMALIZE YOUR DATA

```
1 features = ['sepal length (cm)', 'sepal width (cm)', 'petal length (cm)', 'petal width (cm)']
2
3 # Here we separate features from the target column
4 x = df.loc[:, features].values
5 y = df.loc[:, ['target']].values
6
7 # Standardizing the only the features
8 x = StandardScaler().fit_transform(x)
```

APPLYING PCA

```
1  pca = PCA(n_components=2)
2  principalComponents = pca.fit_transform(x)
3  principalDf = pd.DataFrame(data = principalComponents
4                             , columns = ['principal component 1', 'principal component 2'])
5  finalDf = pd.concat([principalDf, df[['target']], axis = 1)
```

Virus Classification with PCA

OTHER APPLICATION OF PCAs

Dimension Reduction in Machine Learning

Normally, you would like to use PCA (or any other method for dimension reduction) to **speed up** your machine learning algorithm. Instead of learning many instances with over 700 features (e.g. our handwritten letter example from earlier), we can **reduce the number of features**, by only taking the most **important combinations** of features into account.

EXERCISE / HOMEWORK

Have a look at the `viral_data.fasta`.

There are some known viruses, but also some unknown viruses inside. Try to use k-mer frequencies as features and cluster the sequences based on that.

Can you roughly classify the unknown viruses?

You do not have to do this in Python. If you are more familiar with, for example, R, feel free to use this as well.

I am just a little bit Python addicted. :)

COFFEE BREAK

