

# Theoretical and practical metagenomic approaches to viral discovery

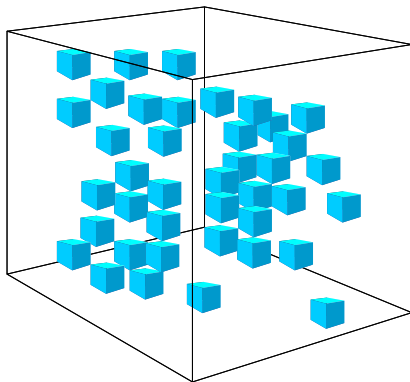
Practical Session: Dimension Reduction and Clustering

Kevin Lamkiewicz, Manja Marz

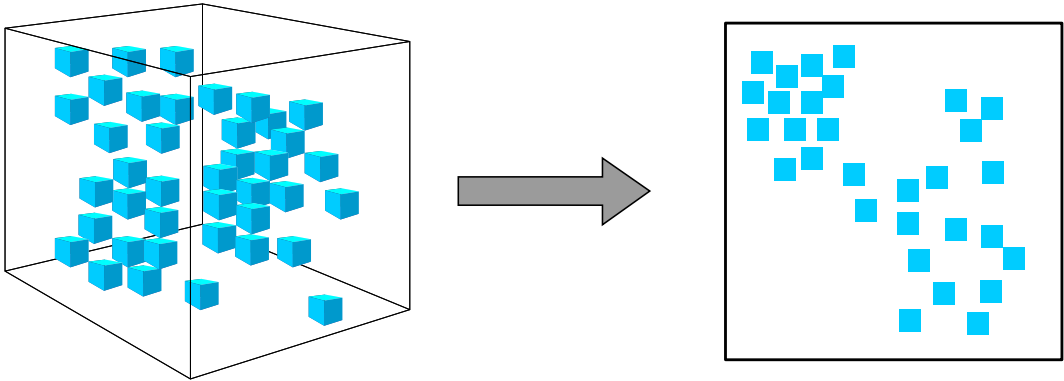
23.10.2019

European Virus Bioinformatics Center

# Dimension Reduction



# Dimension Reduction



## PCA with scikit-learn

## A QUICK EXAMPLE

```
1 from sklearn.decomposition import PCA
2 from sklearn.preprocessing import StandardScaler
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6
7
8 data = []
9 target = []
10 target_names = list(class2id.keys())
11 with open('virus.csv', 'r') as inputStream:
12     ...
```

## A QUICK EXAMPLE

```
1 from sklearn.decomposition import PCA
2 from sklearn.preprocessing import StandardScaler
3 import numpy as np
4 import pandas as pd
5 import matplotlib.pyplot as plt
6
7
8 data = []
9 target = []
10 target_names = list(class2id.keys())
11 with open('virus.csv', 'r') as inputStream:
12     ...
13 df = pd.DataFrame(data= np.c_[data, [id2class[x] for x in target]],
14                   columns= ['feature1', 'feature2', 'feature3', 'feature4'] + ['target'])
```

# NORMALIZE YOUR DATA

```
1 features = ['feature1','feature2','feature3','feature4']
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)
```



# 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.

# Virus Classification with PCA

# 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

