

### Download Chapter pdf

### Download Chapter notebook (ipynb)

### Mandatory Lesson Feedback Survey

- 
- 
- 
- Work with high-dimensional datasets
- Apply dimensionality reduction to extract few, informative features
- Use PCA for dimensionality reduction
- Estimate optimal number of features

### Prereq

- Lesson Clustering Introduction
- Lesson Clustering Images

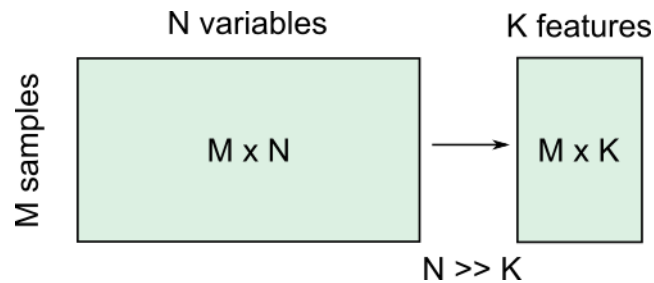
### Code Preparation

```
1 from numpy import reshape, append, mean, pi, linspace, var, array,
   cumsum, arange
2 from numpy.random import seed, multivariate_normal, binomial
3 from numpy.linalg import norm
4
5 from math import cos, sin, atan2
6 from matplotlib.pyplot import subplots, xlabel, ylabel, scatter, title
7 from matplotlib.pyplot import figure, plot, xticks, xlabel, ylabel,
   suptitle, show
8
9 import glob
10 from pathlib import Path
11 from pandas import read_csv, DataFrame
12 import csv
```

### Introduction

Clinical and scientific datasets often contain measurements of many physiological variables, expression of many genes or metabolic markers. These datasets, where each individual observation/sample is described by many variables are referred to as **high-dimensional**. However, some of these features can be highly correlated or redundant, and we sometimes want a more condensed description of the

data. This process of describing the data by a reduced number of new combined features that reduces redundancy, is called **dimensionality reduction**.



**Figure 1:** Goal of decomposition

There are several reasons for performing dimensionality reduction, both for data exploration and for preprocessing:

- Data compression, to optimize memory usage/loading speed and visualisation of “important features”.
- To find the individual variable or combinations of variables that give a condensed but accurate description of the data (sometimes also called feature extraction).
- It can also be used for “denoising” - removing features that only account for small variability.
- It is also useful to perform this step before other supervised (regression, classification) or unsupervised learning procedures (like clustering), by reducing data to “main features” to improve robustness of learning, especially with a limited number of independent observations.

### **Example: Predicting outcome from large-scale genetic or molecular studies**

Clinical research can involve doing an exhaustive measurement of all tentative diagnostic markers, for example levels of various cytokines and hormones, or searching for genetic markers for disease X. There is good reason to expect that these share regulatory/signalling mechanisms, and that their interaction determines disease outcome. Building a regression model including all measurements is problematic because of multicollinearity and the demand for a much larger number of observations.

By performing dimensionality reduction, we first find features that capture the major patterns of covariation of these factors in the sample population. Then we will use these compact features, rather than individual measurements, to train our classifier or regression model, to study outcomes. With fewer features, we train models with less data.

### Loading example data: Childhood Onset Rheumatic Disease gene expression profile

Gene expression data (for selected markers) is available for 83 individuals, including some patients with varying forms of early-onset rheumatic disease (available at EMBL-Bioinformatics database, also included in folder). This is a typical example of high dimensional datasets, with transcript levels for >50,000 genes measured!

Here, the data are read from individual files in the 'Dataset' folder and combined into a Pandas dataframe called 'geneData'.

```
1 # Load all data
2 nGenes = 1000 # keep only these genes for now
3
4 filectr=0
5
6 datapath = Path("data")    #/relative/path/to/folder/with/datasets
7
8 for txt in datapath.glob("*.txt"):
9     b=read_csv(txt, sep='\t', usecols=[1])
10    b=reshape(b.values, (1,54675)) # each file is 1 observation sample
    i.e. 1 row.
11    if filectr==0:
12        coln = read_csv(txt, sep='\t', usecols=[0])
13        a = b
14    else:
15        a=append( a, b, axis=0)
16    filectr+=1
17
18 # convert 2D array to dataframe
19 geneData = DataFrame( data=a[:,range(nGenes)]- mean(a[:,range(nGenes)]),
    axis=0), columns = coln.values[range(nGenes)])
```

The dataframe has 80 samples and 1000 features.

Apart from the large number of variables (compared to number of samples), we also have strongly correlated expression of various markers. We can see that by plotting the correlation matrix of the dataframe.

```
1 fig, ax = subplots(figsize=(8,8))
2 im = ax.matshow(geneData.corr())
3 cbar = fig.colorbar(im, shrink=0.82)
4 xlabel('Locus #', fontsize=14);
5 ylabel('Locus #', fontsize=14);
6
7 show()
```

The little yellow blocks indicate that various features are strongly correlation (values around 1). Further analyses (e.g. checking if we find clusters mapping to healthy individuals and patients, respectively, or

whether patients with different patterns of marker expression can be identified for future diagnostic purposes) can be thus be improved if we perform a dimensionality reduction, i.e. decrease the number of features while maintaining relevant information.

### Principal Component Analysis

Principal component analysis (PCA) is one of the most commonly used dimensionality reduction methods. Instead of describing the data using the original measurements, it finds a new set of ranked features (the “basis”). The new “features” are linear combinations of the original variables, such that they are all *orthogonal* to each other, and the “importance” of a feature (called the **principal component**) is generally given by how much of the variability in the dataset it captures (the so-called **explained variance**). Thus we first transform the data from the original variables into a new set of features or principal components (PCs).

To get a lower dimensional description of the data, we retain only the top features (PCs) such that the largest variations or ‘patterns’ in the dataset are preserved.

To apply **PCA** from `scikit-learn`, we need to specify the number of components (`n_components`) we want to reduce the data to. `n_components` must be less than the *rank* of the dataset, i.e. less than the smaller of the number of samples and number of variables measured.

### Simple example with 2D generated dataset

To first explore what dimensionality reduction with PCA looks like, we will work with a simple 2-dimensional (2D) dataset, where we want to find a 1D description instead.

We start with generating some observations from a two-dimensional (multivariate) Gaussian distribution, with mean and covariance specified.

```
1 RND_SEED = 7890
2 seed(RND_SEED)
3
4 # Set up 2D multivariate Gaussians
5 means = [20, 20]
6
7 cov_mat = [[1, .85], [.85, 1.5]] # 2x2 covariance matrix must be
    symmetric and positive-semidefinite(>=0)
8
9 # Generate data
10 nSamples = 300
11
12 data = multivariate_normal(means, cov_mat, nSamples)
13
14 print(data.shape)
```

```
1 (300, 2)
```

We can check the distribution of the two features in histograms, and we can see how the two data dimensions relate to each other in the scatter plot.

```
1 # Center data
2 zdata = data - mean(data, axis=0)
3
4 fig, ax = subplots(ncols=3, nrows=1, figsize=(15, 4))
5
6 ax[0].hist(zdata[:, 0]);
7 ax[0].set_title('Dim 1', fontsize=16, fontweight='bold');
8 ax[0].set_xlabel('Values', fontsize=14);
9 ax[0].set_ylabel('Frequency', fontsize=14);
10
11 ax[1].hist(zdata[:, 1]);
12 ax[1].set_title('Dim 2', fontsize=16, fontweight='bold');
13 ax[1].set_xlabel('Values', fontsize=14);
14
15
16 ax[2].scatter(zdata[:, 0], zdata[:, 1]);
17 ax[2].set_title('Sample Data', fontsize=16, fontweight='bold')
18 ax[2].set_xlabel('Dim 1', fontsize=14);
19 ax[2].set_ylabel('Dim 2', fontsize=14);
20
21 show()
```

### Note

*The last plot, with the samples plotted as scatter in the space of measured variables is called a “state space plot”. The main axes correspond to individual variables.*

Now if we want to represent data by a single feature, we can see that the individual variables are not ideal. If we only keep Dim 1, we lose substantial variability along Dim 2. We see that the data points are stretched along the main diagonal. This implies that the two variables co-vary to some extent. Thus, a linear combination of the two might capture this relevant pattern.

A linear combination correspond to a vector (or direction) in this 2D state space.

For example, the x-axis vector  $(1, 0) = [1 \times \text{Dim1} + 0 \times \text{Dim2}]$ , i.e. it contains Dim 1 only. More generally,  $(a, b)$  represents the linear combination  $[a \times \text{Dim1} + b \times \text{Dim2}]$ , i.e. a mixture of the two variables. Typically, we ensure that the length of these vectors is 1, to avoid artifactual rescaling of data.

### PCA as projection or rotation of basis

Instead of representing data by their values along Dim 1 and Dim 2, we now try to describe them along different directions. We do so by *projecting* 2D data onto that vector. To quantify the importance of each direction, we measure the variance of the projected samples.

Let us examine different vectors that point at various angles compared to the x-axis. Remember that the vector denoted by the angle =  $(\cos(\text{angle}), \sin(\text{angle}))$  denotes the following linear combination of the 2D sample  $d=[d1, d2]$ :

$$f(d1,d2) = \cos(\text{angle}) * d1 + \sin(\text{angle}) * d2$$

```
1  ## Centre data at origin by subtracting mean
2  zdata = data - mean(data, axis=0)
3
4  # Define a feature as a direction in the 2D space using the angle
5  angle = pi/4
6
7  # Obtain the vector
8  vec = [cos(angle), sin(angle)] # has norm=1 by definition, otherwise:
    vec = vec/norm(vec)
9
10 # Project 2D data onto this single dimension
11 dim1_data = zdata.dot(vec)
12
13
14 fig, ax = subplots(ncols=2, nrows=1, figsize=(15, 6))
15
16 # Plot data and the feature direction
17 ax[0].scatter(zdata[:, 0], zdata[:, 1]);
18 ax[0].quiver([0], [0], vec[0], vec[1], scale=3);
19
20 ax[0].set_title('Sample Data', fontsize=16, fontweight='bold')
21 ax[0].set_xlabel('Dim 1', fontsize=14); ax[0].set_xlim([-4,4]);
22 ax[0].set_ylabel('Dim 2', fontsize=14); ax[0].set_ylim([-4,4]);
23
24 # Plot histogram of new projected data
```

```
1  Text(0, 0.5, 'Dim 2')
2  (-4.0, 4.0)
```

```
1  ax[1].hist(dim1_data, bins = linspace(-3,3,9));
2  ax[1].set_title('Along vec', fontsize=16, fontweight='bold');
3  ax[1].set_xlabel('Values', fontsize=14);
4  ax[1].set_ylabel('Frequency', fontsize=14);
5
6  show()
```

### Do it Yourself

Change the angle in the code above and see how distribution of projected data changes. Use e.g.  $0$ ,  $\pi/2$ ,  $2/3\pi$ , and  $\pi$ . Try to estimate how much the variance changes.

Now let's formally calculate the variance by projecting along a set of different directions. We use angles between  $0$  and  $\pi$  (i.e. between  $0$  and  $180$  degrees). For each direction (projection) we calculate the variance and plot it against the angle.

```
1 allAngles = linspace( start=0, stop=pi, num=20)
2 allVar = []
3 for angle in allAngles:
4     vec = array([cos(angle), sin(angle)]) #Unit length vector on which
        data will be projected
5     dim1_data = data.dot(vec)
6     allVar.append( var(dim1_data))
7
8 # Plot variance captured by different features
9 scatter( allAngles, allVar );
10 title('Variance along different directions', fontsize=16, fontweight='
    bold')
11 xlabel('Direction (angle in radians)', fontsize=14);
12 ylabel('Variance of projected data', fontsize=14);
```

It can be seen there is a particular direction that retains maximal variance. This would be the optimal feature to give a 1D description of the 2D data while retaining maximum variability. PCA is the method that directly find this optimal direction. This makes it especially powerful for high dimensional datasets.

### Find 1-D description using PCA

Let us now use dimensionality reduction [PCA](#) to directly find the direction that captures maximal variance. For an  $n$ -dimensional dataset, PCA finds a set of  $n$  new features, ranked by the variance along each feature. See the [scikit-learn documentation](#) for information.

To get the best 1D representation of data, we instantiate the class with one component ( $n\_components = 1$ ) to be returned, and transform the data to the feature space of that component.

```
1 from sklearn.decomposition import PCA
2
3 # Initialize the PCA model
4 pcaResults = PCA(n_components=1, whiten=False) # specify no. of
        components, and whether to standardize
5
6 # Fit to data
7 pcaResults = pcaResults.fit(data)
```

```
8
9 # Transform the data, using our learned PCA representation
10 dim1_data = pcaResults.transform(data)
11
12 print( "Size of new dataset: ", dim1_data.shape )
```

```
1 Size of new dataset: (300, 1)
```

```
1 Size of new dataset: (300, 1)
```

From the fitted model 'pcaResults' we can now get:

- The fractional variance using `explained_variance_ratio_`. It tells us what fraction of the total variance is retained by the reduced 1D description of the data.
- The direction (angle) of this first PC as calculated from the arcus tangens `atan` of the vector. You can compare it directly to the angle in the above plot.

```
1 # How much variance was captured compared to original data?
2 print("Fractional variance captured by first PC: {:.14f}.".format(
    pcaResults.explained_variance_ratio_[0]))
```

```
1 Fractional variance captured by first PC: 0.8722.
```

```
1 vec1 = pcaResults.components_[0,:]
2 print("Direction of PC1 is at angle = {:.12f} radians".format(atan2(
    vec1[1],vec1[0])))
```

```
1 Direction of PC1 is at angle = 0.99 radians
```

### Do it Yourself

Exercise (DIY): Change the covariance of the synthetic data, and repeat the above steps. Try to infer the relationship between the shape of the scatter, the strength of correlation between the two dimensions, and the variance captured by the first PC.

### Find 2-D description using PCA

Our dataset has two features. We can thus obtain two principal components. The two components must be orthogonal to each other (form a right angle). Let us get the corresponding vectors and plot them on the scatterplot.

```
1 # Initialize the PCA model
```



```
2  pcaResults = PCA(n_components=2, whiten=False) # specify no. of
    components, and whether to standardize
3
4  # Fit to data
5  pcaResults = pcaResults.fit(data)
6
7  # Transform the data, using our learned PCA representation
8  dim1_data = pcaResults.transform(data)
9
10 print( "Size of new dataset: ", dim1_data.shape )
```

```
1 Size of new dataset: (300, 2)
```

```
1 print(pcaResults.components_)
```

```
1 [[ 0.54864625  0.8360546 ]
2  [ 0.8360546  -0.54864625]]
```

```
1 vec1 = pcaResults.components_[0,:]
2 vec2 = pcaResults.components_[1,:]
3
4
5 fig, ax = subplots(figsize=(6, 6))
6
7 # Plot data and the feature direction
8 ax.scatter(zdata[:, 0], zdata[:, 1]);
9 ax.quiver([0], [0], vec1[0], vec1[1], scale=3, color='r');
10 ax.quiver([0], [0], vec2[0], vec2[1], scale=3, color='r');
11
12 ax.set_title('Sample Data', fontsize=16, fontweight='bold')
13 ax.set_xlabel('Dim 1', fontsize=14); ax.set_xlim([-4,4]);
14 ax.set_ylabel('Dim 2', fontsize=14); ax.set_ylim([-4,4]);
15
16 show()
```

The red arrows show the new coordinate system, defined by the principal components. It also shows that the more correlated the data are, the less information we lose by reducing the data representation using PCA.

### Important

When reading about PCA you will come across the term **Singular Value Decomposition** (SVD). As mentioned in the scikit-learn documentation for PCA, when you apply [PCA](#) the algorithm in the background uses SVD to return the results. The difference between the two is technical but SVD is the more general method. Wikipedia has a mathematical introduction to PCA which also discusses the relationship to SVD.

### PCA of the Gene Expression Data

Now, we will use PCA to find a few features that capture most of the variability in the gene expression dataset. we arbitrarily start with 50 components assuming that this will include the important features.

```
1 nComp = 50 # Number of PCs to be returned
2 #trainIndx = binomial(1,0.9,size=filectr)
3
4 GenePCA = PCA(n_components=nComp, whiten=True)
5 GenePCA = GenePCA.fit(geneData) #.values[trainIndx==1,:]
```

How many features should we retain? To investigate this question, we plot the total (cumulative) variance for retaining the top  $k$  modes. The more we use, the higher the cumulative variance.

```
1 cumExpVar = cumsum(GenePCA.explained_variance_ratio_)
2 fig = figure(figsize=(5, 5))
3 im = plot( range(nComp), cumExpVar )
4 xticks(arange(0,nComp,5));
5 xlabel( 'Number of PCs', fontsize=14);
6 ylabel('Cumulative explained variance', fontsize=14);
7
8 show()
```

A common heuristic for choosing the number of components is by defining a set threshold for the total explained variance. These threshold commonly vary between 0.8-0.9, depending on the structure of the PCs, e.g. depending on whether there are a few top PCs or many small PCs, but also depending on the expected noise in data, and on the desirable accuracy of the reduced data set. While a dimensionality reduction is convenient it always loses some information.

As an example, we can check how many PCs we need to retain 99 % of the variance.

```
1 threshold = 0.99
2
3 keepPC = [pc for pc in range(nComp) if cumExpVar[pc]>=threshold][0]
4
5 print('Number of features needed to explain {:.2f} fraction of total
      variance is {:2d}'.format(threshold, keepPC) )
```

```
1 Number of features needed to explain 0.99 fraction of total variance is
  19.
```

There are many alternative methods for estimating the number of features. These include:

- Plot explained variance of individual PCs and cut off when there is a sharp drop/saturation in the values.

- Cross-validation - Find number of PCs that maximize explained variance on heldout data (using a bi-cross validation scheme).
- Visual inspection or interpretable PCs

### Visualization of reduced dataset

Now that we have selected the number of features to be used, we can see what the data along those features looks like.

```
1 newGeneData = GenePCA.transform(geneData)[: ,range(keepPC)]
2
3 fig, ax = subplots(ncols=3, nrows=2, figsize=(18, 8))
4 subtitle('Data in new Feature space', fontsize=20)
5
6 ax[0, 0].scatter(newGeneData[:,0], newGeneData[:,1], c=range(filectr),
7                  cmap='summer');
8 ax[0, 0].set_xlabel('PC1', fontsize=14); ax[0, 0].set_xlim([-2,2]);
9 ax[0, 0].set_ylabel('PC2', fontsize=14); ax[0, 0].set_ylim([-2,2]);
10
11 ax[0, 1].scatter(newGeneData[:,0], newGeneData[:,2], c=range(filectr),
12                  cmap='summer');
13 ax[0, 1].set_xlabel('PC1', fontsize=14); ax[0, 1].set_xlim([-2,2]);
14 ax[0, 1].set_ylabel('PC3', fontsize=14); ax[0, 1].set_ylim([-2,2]);
15
16 ax[0, 2].scatter(newGeneData[:,0], newGeneData[:,3], c=range(filectr),
17                  cmap='summer');
18 ax[0, 2].set_xlabel('PC1', fontsize=14); ax[0, 2].set_xlim([-2,2]);
19 ax[0, 2].set_ylabel('PC4', fontsize=14); ax[0, 2].set_ylim([-2,2]);
20
21
22 ax[1, 0].scatter(newGeneData[:,1], newGeneData[:,2], c=range(filectr),
23                  cmap='summer');
24 ax[1, 0].set_xlabel('PC2', fontsize=14); ax[1, 0].set_xlim([-2,2]);
25 ax[1, 0].set_ylabel('PC3', fontsize=14); ax[1, 0].set_ylim([-2,2]);
26
27 ax[1, 1].scatter(newGeneData[:,1], newGeneData[:,3], c=range(filectr),
28                  cmap='summer');
29 ax[1, 1].set_xlabel('PC2', fontsize=14); ax[1, 1].set_xlim([-2,2]);
30 ax[1, 1].set_ylabel('PC4', fontsize=14); ax[1, 1].set_ylim([-2,2]);
31
32 ax[1, 2].scatter(newGeneData[:,3], newGeneData[:,2], c=range(filectr),
33                  cmap='summer');
34 ax[1, 2].set_xlabel('PC4', fontsize=14); ax[1, 2].set_xlim([-2,2]);
```

```
34 ax[1, 2].set_ylabel('PC3', fontsize=14); ax[1, 2].set_ylim([-2,2]);  
35  
36 show()
```

### Beyond dimensionality reduction

After reducing dataset from hundreds or thousands of genes to fewer features (PCs), it can be used:

- To find groups of genes that seem to be co-regulated (using clustering).
- To build classifiers to predict the likelihood of juvenile or late-onset rheumatoid disease. For example, one feature can be highly variable in the sampled population and therefore has high explained variance, but it may not be predictive for the current objective of predicting disease phenotype (e.g. eye pigmentation genes).
- To find what combination of genes the top features correspond to.
- To suggest what individual genes are most variable.
- To discover co-expression patterns of multiple genes.

### Optional Exercises

1. The dataset contains two types of samples, expression in neutrophil and expression in peripheral blood mononuclear cells (PBMC). Download the sample-data table to get the category of each sample. Make a scatter plot of two (original) features and set the colour of the scatter based on the sample's origin. Do a PCA, scatter the top scoring components and colour the samples in the same way. Do the top PCs capture gene expression differences in the two types of cells?
2. The dataset similarly contains samples from individuals with differing phenotypes - healthy, as well as 2 different disease characteristics. Similar to the previous exercise, colour the scatter based on the disease phenotypes. Do the top PCs capture variability in gene expression for different phenotypes?

### When does PCA fail?

Although PCA is the first and simplest method of exploring high-dimensional datasets, there are important caveats to keep in mind:

1. *PCA is a linear method*

Nonlinear patterns of co-expression will be seemingly broken up into many features. Thus the ability to find whether genes are co-regulated may be reduced if the effects are nonlinear.

2. *PCA depends on scaling of individual measurements.*

As PCA measures variance along different directions, changing the scale of one variable (for example from cm to mm, or fold-change to actual levels) may drastically change the dominant “features” extracted.

The lesson on clustering of images with Gaussian mixture models contained instructions to scale data. You can explore the original dataset of this lesson, to investigate how scaling of variables impacts the results.

3. *Estimating number of features is heuristic, and depends on the purpose.*

The definition of what is the ‘relevant’ number of features can depend on:

- the signal-to-noise ratio in data. For example, noisier or smaller datasets make it harder to accurately estimate smaller PCs. In such cases, a more conservative threshold should be used.
- the classifier or regression performance following the dimensionality reduction.

## Resources

- Detailed tutorial on PCA
- Other methods for dimensionality reduction:
  - Matrix factorization methods
  - Manifold learning or topological methods
  - Example pipelines with scikit-learn

Dataset from:

Frank MB, Wang S, Aggarwal A, et al. **Disease-associated pathophysiologic structures in pediatric rheumatic diseases show characteristics of scale-free networks seen in physiologic systems: implications for pathogenesis and treatment.** *BMC Med Genomics*. 2009;2:9. Published 2009 Feb 23. doi:10.1186/1755-8794-2-9

Download data from: <https://www.ebi.ac.uk/arrayexpress/experiments/E-GEOD-11083/>

## Keypoints

-

- 
-