

Principal Component Analysis

 sebastianraschka.com/Articles/2015_pca_in_3_steps.html

January 27, 2015

Principal Component Analysis (PCA) is a simple yet popular and useful linear transformation technique that is used in numerous applications, such as stock market predictions, the analysis of gene expression data, and many more. In this tutorial, we will see that PCA is not just a “black box”, and we are going to unravel its internals in 3 basic steps.

This article just got a complete overhaul, the original version is still available at [principal_component_analysis_old.ipynb](#).

Sections

Introduction

The sheer size of data in the modern age is not only a challenge for computer hardware but also a main bottleneck for the performance of many machine learning algorithms. The main goal of a PCA analysis is to identify patterns in data; PCA aims to detect the correlation between variables. If a strong correlation between variables exists, the attempt to reduce the dimensionality only makes sense. In a nutshell, this is what PCA is all about: Finding the directions of maximum variance in high-dimensional data and project it onto a smaller dimensional subspace while retaining most of the information.

PCA Vs. LDA

Both Linear Discriminant Analysis (LDA) and PCA are linear transformation methods. PCA yields the directions (principal components) that maximize the variance of the data, whereas LDA also aims to find the directions that maximize the separation (or discrimination) between different classes, which can be useful in pattern classification problem (PCA “ignores” class labels).

In other words, PCA projects the entire dataset onto a different feature (sub)space, and LDA tries to determine a suitable feature (sub)space in order to distinguish between patterns that belong to different classes.

PCA and Dimensionality Reduction

Often, the desired goal is to reduce the dimensions of a d -dimensional dataset by projecting it onto a k ($k < d$)-dimensional subspace (where $k < d < d$) in order to increase the computational efficiency while retaining most of the information. An important question is “what is the size of k that represents the data ‘well’?”

Later, we will compute eigenvectors (the principal components) of a dataset and collect them in a projection matrix. Each of those eigenvectors is associated with an eigenvalue which can be interpreted as the “length” or “magnitude” of the corresponding eigenvector. If some eigenvalues have a significantly larger magnitude than others, then the reduction of the dataset via PCA onto a smaller dimensional subspace by dropping the “less informative” eigenpairs is reasonable.

A Summary of the PCA Approach

- Standardize the data.
- Obtain the Eigenvectors and Eigenvalues from the covariance matrix or correlation matrix, or perform Singular Value Decomposition.
- Sort eigenvalues in descending order and choose the k eigenvectors that correspond to the k largest eigenvalues where k is the number of dimensions of the new feature subspace ($k \leq d$).
- Construct the projection matrix W from the selected k eigenvectors.
- Transform the original dataset X via W to obtain a k -dimensional feature subspace Y .

Preparing the Iris Dataset

About Iris

For the following tutorial, we will be working with the famous “Iris” dataset that has been deposited on the UCI machine learning repository (<https://archive.ics.uci.edu/ml/datasets/Iris>).

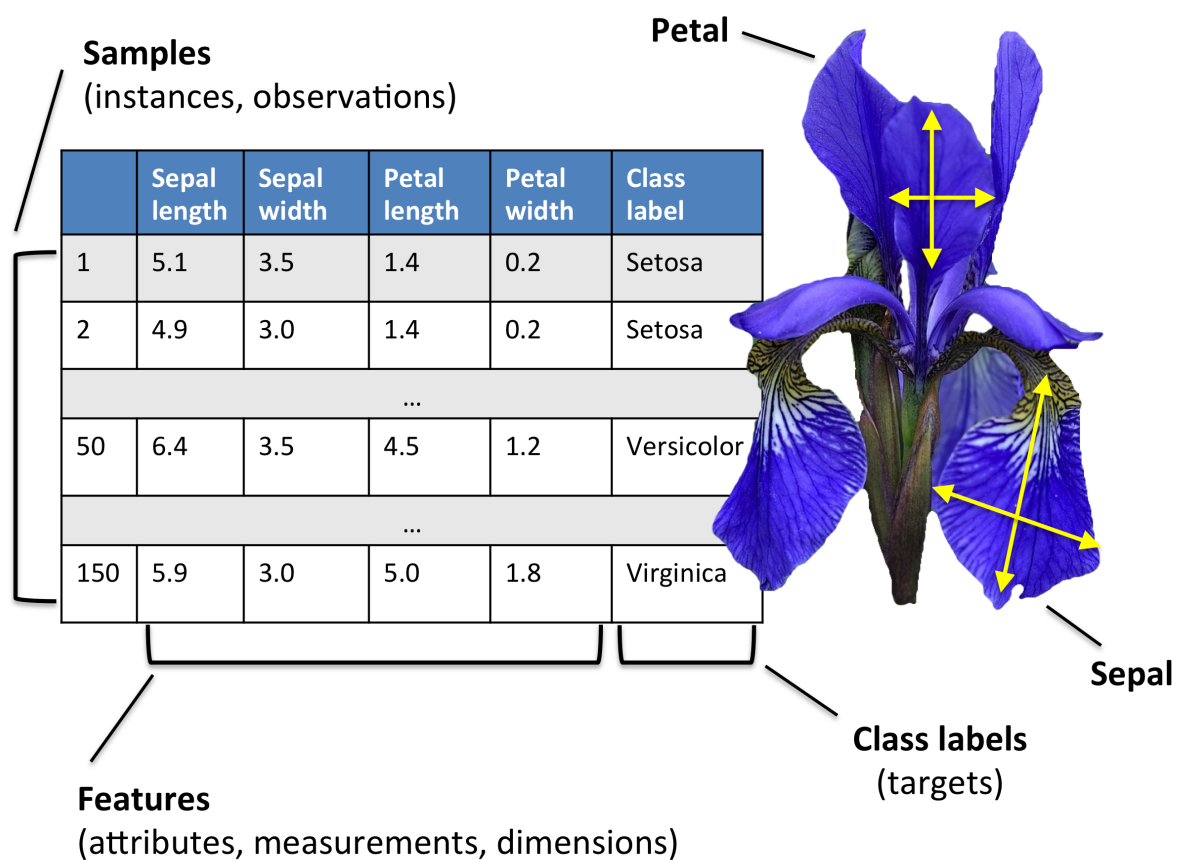
The iris dataset contains measurements for 150 iris flowers from three different species.

The three classes in the Iris dataset are:

1. *Iris-setosa* ($n=50$)
2. *Iris-versicolor* ($n=50$)
3. *Iris-virginica* ($n=50$)

And the four features of in Iris dataset are:

1. *sepal length* in cm
2. *sepal width* in cm
3. *petal length* in cm
4. *petal width* in cm



Loading the Dataset

In order to load the Iris data directly from the UCI repository, we are going to use the superb pandas library. If you haven't used pandas yet, I want encourage you to check out the pandas tutorials. If I had to name one Python library that makes working with data a wonderfully simple task, this would definitely be pandas!

```
import pandas as pd

df = pd.read_csv(
    filepath_or_buffer='https://archive.ics.uci.edu/ml/machine-learning-
databases/iris/iris.data',
    header=None,
    sep=',')

df.columns=['sepal_len', 'sepal_wid', 'petal_len', 'petal_wid', 'class']
df.dropna(how="all", inplace=True) # drops the empty line at file-end

df.tail()
```

	sepal_len	sepal_wld	petal_len	petal_wld	class
145	6.7	3.0	5.2	2.3	Iris-virginica
146	6.3	2.5	5.0	1.9	Iris-virginica
147	6.5	3.0	5.2	2.0	Iris-virginica
148	6.2	3.4	5.4	2.3	Iris-virginica
149	5.9	3.0	5.1	1.8	Iris-virginica

```
# split data table into data X and class labels y
```

```
X = df.ix[:,0:4].values
```

```
y = df.ix[:,4].values
```

Our iris dataset is now stored in form of a 150×4 matrix where the columns are the different features, and every row represents a separate flower sample. Each sample row xx can be pictured as a 4-dimensional vector

$$xT = (x_1, x_2, x_3, x_4) = (\text{sepal length}, \text{sepal width}, \text{petal length}, \text{petal width})$$

Exploratory Visualization

To get a feeling for how the 3 different flower classes are distributed along the 4 different features, let us visualize them via histograms.

```

from matplotlib import pyplot as plt
import numpy as np
import math

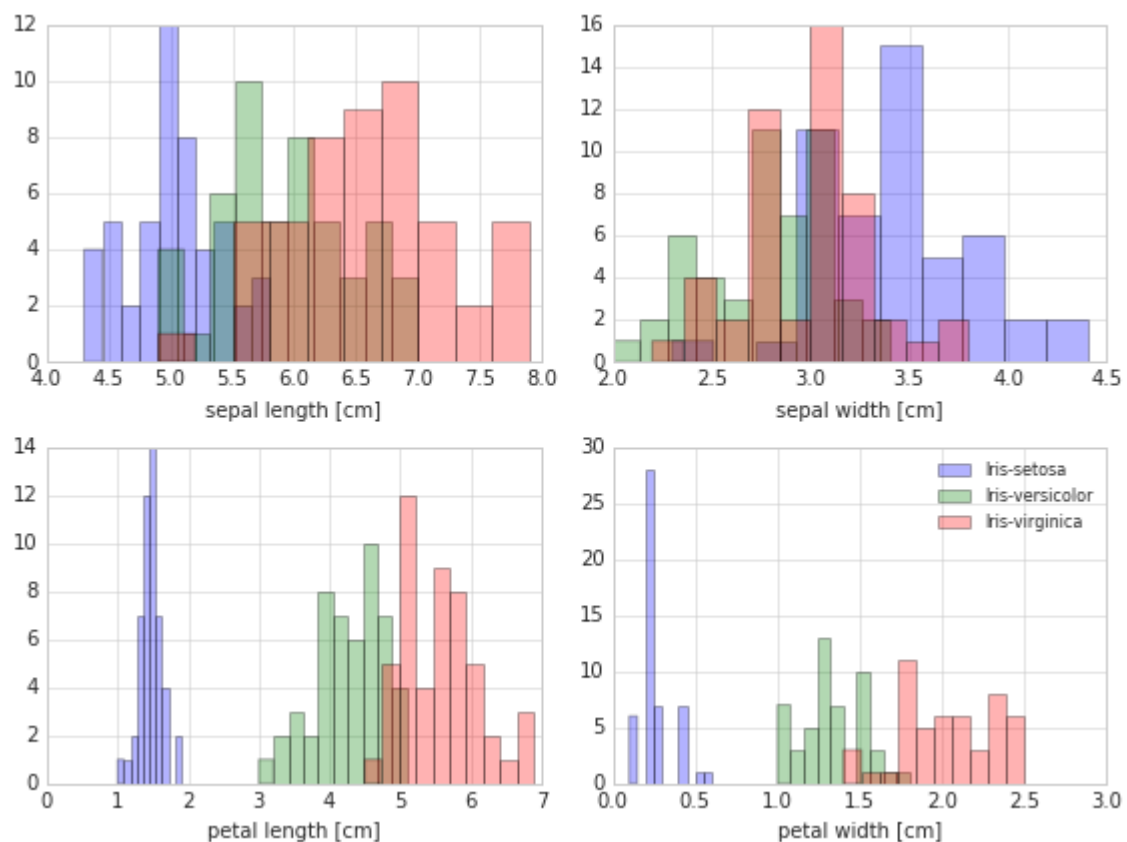
label_dict = {1: 'Iris-Setosa',
              2: 'Iris-Versicolor',
              3: 'Iris-Virginica'}

feature_dict = {0: 'sepal length [cm]',
               1: 'sepal width [cm]',
               2: 'petal length [cm]',
               3: 'petal width [cm]'}

with plt.style.context('seaborn-whitegrid'):
    plt.figure(figsize=(8, 6))
    for cnt in range(4):
        plt.subplot(2, 2, cnt+1)
        for lab in ('Iris-setosa', 'Iris-versicolor', 'Iris-virginica'):
            plt.hist(X[y==lab, cnt],
                    label=lab,
                    bins=10,
                    alpha=0.3,)
        plt.xlabel(feature_dict[cnt])
    plt.legend(loc='upper right', fancybox=True, fontsize=8)

plt.tight_layout()
plt.show()

```



Standardizing

Whether to standardize the data prior to a PCA on the covariance matrix depends on the measurement scales of the original features. Since PCA yields a feature subspace that maximizes the variance along the axes, it makes sense to standardize the data, especially, if it was measured on different scales. Although, all features in the Iris dataset were measured in centimeters, let us continue with the transformation of the data onto unit scale (mean=0 and variance=1), which is a requirement for the optimal performance of many machine learning algorithms.

```
from sklearn.preprocessing import StandardScaler
X_std = StandardScaler().fit_transform(X)
```

1 - Eigendecomposition - Computing Eigenvectors and Eigenvalues

The eigenvectors and eigenvalues of a covariance (or correlation) matrix represent the “core” of a PCA: The eigenvectors (principal components) determine the directions of the new feature space, and the eigenvalues determine their magnitude. In other words, the eigenvalues explain the variance of the data along the new feature axes.

Covariance Matrix

The classic approach to PCA is to perform the eigendecomposition on the covariance matrix Σ , which is a $d \times d \times d$ matrix where each element represents the covariance between two features. The covariance between two features is calculated as follows:

$$\sigma_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k).$$

$$\sigma_{jk} = \frac{1}{n-1} \sum_{i=1}^n (x_{ij} - \bar{x}_j)(x_{ik} - \bar{x}_k).$$

We can summarize the calculation of the covariance matrix via the following matrix equation:

$$\Sigma = \frac{1}{n-1} ((X - \bar{X})^T (X - \bar{X})) \Sigma = \frac{1}{n-1} ((X - \bar{X})^T (X - \bar{X}))$$

where \bar{x} is the mean vector $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$. $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$.

The mean vector is a d -dimensional vector where each value in this vector represents the sample mean of a feature column in the dataset.

```
import numpy as np
mean_vec = np.mean(X_std, axis=0)
cov_mat = (X_std - mean_vec).T.dot((X_std - mean_vec)) / (X_std.shape[0]-1)
print('Covariance matrix \n%s' %cov_mat)
```

```
Covariance matrix
[[ 1.00671141 -0.11010327  0.87760486  0.82344326]
 [-0.11010327  1.00671141 -0.42333835 -0.358937 ]
 [ 0.87760486 -0.42333835  1.00671141  0.96921855]
 [ 0.82344326 -0.358937    0.96921855  1.00671141]]
```

The more verbose way above was simply used for demonstration purposes, equivalently, we could have used the numpy `cov` function:

```
print('NumPy covariance matrix: \n%s' %np.cov(X_std.T))
```

NumPy covariance matrix:

```
[[ 1.00671141 -0.11010327  0.87760486  0.82344326]
 [-0.11010327  1.00671141 -0.42333835 -0.358937   ]
 [ 0.87760486 -0.42333835  1.00671141  0.96921855]
 [ 0.82344326 -0.358937    0.96921855  1.00671141]]
```

Next, we perform an eigendecomposition on the covariance matrix:

```
cov_mat = np.cov(X_std.T)
```

```
eig_vals, eig_vecs = np.linalg.eig(cov_mat)
```

```
print('Eigenvectors \n%s' %eig_vecs)
print('\nEigenvalues \n%s' %eig_vals)
```

Eigenvectors

```
[[ 0.52237162 -0.37231836 -0.72101681  0.26199559]
 [-0.26335492 -0.92555649  0.24203288 -0.12413481]
 [ 0.58125401 -0.02109478  0.14089226 -0.80115427]
 [ 0.56561105 -0.06541577  0.6338014   0.52354627]]
```

Eigenvalues

```
[ 2.93035378  0.92740362  0.14834223  0.02074601]
```

Correlation Matrix

Especially, in the field of “Finance,” the correlation matrix typically used instead of the covariance matrix. However, the eigendecomposition of the covariance matrix (if the input data was standardized) yields the same results as a eigendecomposition on the correlation matrix, since the correlation matrix can be understood as the normalized covariance matrix.

Eigendecomposition of the standardized data based on the correlation matrix:

```
cor_mat1 = np.corrcoef(X_std.T)
```

```
eig_vals, eig_vecs = np.linalg.eig(cor_mat1)
```

```
print('Eigenvectors \n%s' %eig_vecs)
print('\nEigenvalues \n%s' %eig_vals)
```

Eigenvectors

```
[ [ 0.52237162 -0.37231836 -0.72101681  0.26199559]
  [-0.26335492 -0.92555649  0.24203288 -0.12413481]
  [ 0.58125401 -0.02109478  0.14089226 -0.80115427]
  [ 0.56561105 -0.06541577  0.6338014  0.52354627]]
```

Eigenvalues

```
[ 2.91081808  0.92122093  0.14735328  0.02060771]
```

Eigendecomposition of the raw data based on the correlation matrix:

```
cor_mat2 = np.corrcoef(X.T)
```

```
eig_vals, eig_vecs = np.linalg.eig(cor_mat2)
```

```
print('Eigenvectors \n%s' %eig_vecs)
print('\nEigenvalues \n%s' %eig_vals)
```

Eigenvectors

```
[ [ 0.52237162 -0.37231836 -0.72101681  0.26199559]
  [-0.26335492 -0.92555649  0.24203288 -0.12413481]
  [ 0.58125401 -0.02109478  0.14089226 -0.80115427]
  [ 0.56561105 -0.06541577  0.6338014  0.52354627]]
```

Eigenvalues

```
[ 2.91081808  0.92122093  0.14735328  0.02060771]
```

We can clearly see that all three approaches yield the same eigenvectors and eigenvalue pairs:

- Eigendecomposition of the covariance matrix after standardizing the data.
- Eigendecomposition of the correlation matrix.
- Eigendecomposition of the correlation matrix after standardizing the data.

Singular Value Decomposition

While the eigendecomposition of the covariance or correlation matrix may be more intuitive, most PCA implementations perform a Singular Value Decomposition (SVD) to improve the computational efficiency. So, let us perform an SVD to confirm that the result are indeed the same:

```
u,s,v = np.linalg.svd(X_std.T)
```

u

```
array([[ -0.52237162, -0.37231836,  0.72101681,  0.26199559],
       [  0.26335492, -0.92555649, -0.24203288, -0.12413481],
       [ -0.58125401, -0.02109478, -0.14089226, -0.80115427],
       [ -0.56561105, -0.06541577, -0.6338014 ,  0.52354627]])
```


2 - Selecting Principal Components

Sorting Eigenpairs

The typical goal of a PCA is to reduce the dimensionality of the original feature space by projecting it onto a smaller subspace, where the eigenvectors will form the axes. However, the eigenvectors only define the directions of the new axis, since they have all the same unit length 1, which can be confirmed by the following two lines of code:

```
for ev in eig_vecs.T:
    np.testing.assert_array_almost_equal(1.0, np.linalg.norm(ev))
print('Everything ok!')
```

Everything ok!

In order to decide which eigenvector(s) can be dropped without losing too much information for the construction of a lower-dimensional subspace, we need to inspect the corresponding eigenvalues: The eigenvectors with the lowest eigenvalues bear the least information about the distribution of the data; those are the ones that can be dropped.

In order to do so, the common approach is to rank the eigenvalues from highest to lowest in order to choose the top k eigenvectors.

```
# Make a list of (eigenvalue, eigenvector) tuples
eig_pairs = [(np.abs(eig_vals[i]), eig_vecs[:,i]) for i in range(len(eig_vals))]

# Sort the (eigenvalue, eigenvector) tuples from high to low
eig_pairs.sort(key=lambda x: x[0], reverse=True)

# Visually confirm that the list is correctly sorted by decreasing eigenvalues
print('Eigenvalues in descending order:')
for i in eig_pairs:
    print(i[0])
```

```
Eigenvalues in descending order:
2.91081808375
0.921220930707
0.147353278305
0.0206077072356
```

Explained Variance

After sorting the eigenpairs, the next question is “how many principal components are we going to choose for our new feature subspace?” A useful measure is the so-called “explained variance,” which can be calculated from the eigenvalues. The explained variance tells us how much information (variance) can be attributed to each of the principal components.

```

tot = sum(eig_vals)
var_exp = [(i / tot)*100 for i in sorted(eig_vals, reverse=True)]
cum_var_exp = np.cumsum(var_exp)

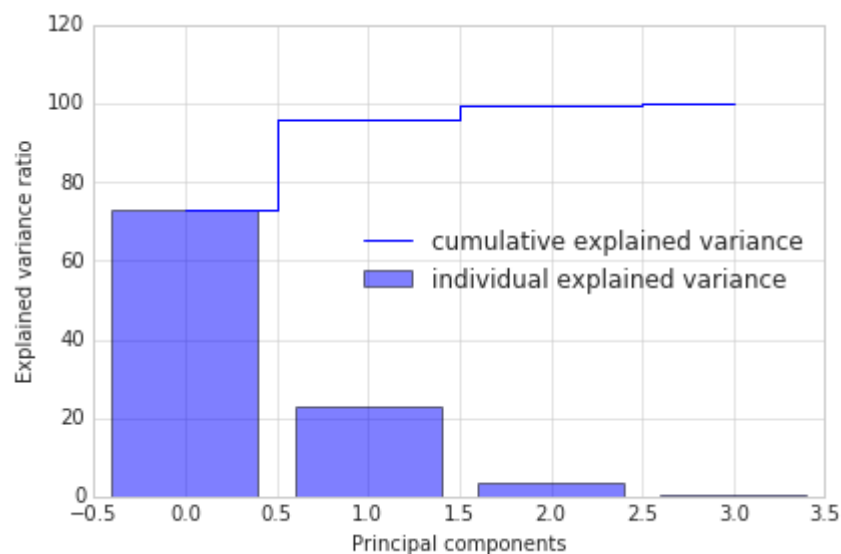
```

```

with plt.style.context('seaborn-whitegrid'):
    plt.figure(figsize=(6, 4))

    plt.bar(range(4), var_exp, alpha=0.5, align='center',
              label='individual explained variance')
    plt.step(range(4), cum_var_exp, where='mid',
             label='cumulative explained variance')
    plt.ylabel('Explained variance ratio')
    plt.xlabel('Principal components')
    plt.legend(loc='best')
    plt.tight_layout()

```



The plot above clearly shows that most of the variance (72.77% of the variance to be precise) can be explained by the first principal component alone. The second principal component still bears some information (23.03%) while the third and fourth principal components can safely be dropped without losing to much information. Together, the first two principal components contain 95.8% of the information.

Projection Matrix

It's about time to get to the really interesting part: The construction of the projection matrix that will be used to transform the Iris data onto the new feature subspace. Although, the name “projection matrix” has a nice ring to it, it is basically just a matrix of our concatenated top k eigenvectors.

Here, we are reducing the 4-dimensional feature space to a 2-dimensional feature subspace, by choosing the “top 2” eigenvectors with the highest eigenvalues to construct our $d \times k$ -dimensional eigenvector matrix W .

```

matrix_w = np.hstack((eig_pairs[0][1].reshape(4,1),
                      eig_pairs[1][1].reshape(4,1)))

print('Matrix W:\n', matrix_w)

```

Matrix W:

```
[[ 0.52237162 -0.37231836]
 [-0.26335492 -0.92555649]
 [ 0.58125401 -0.02109478]
 [ 0.56561105 -0.06541577]]
```

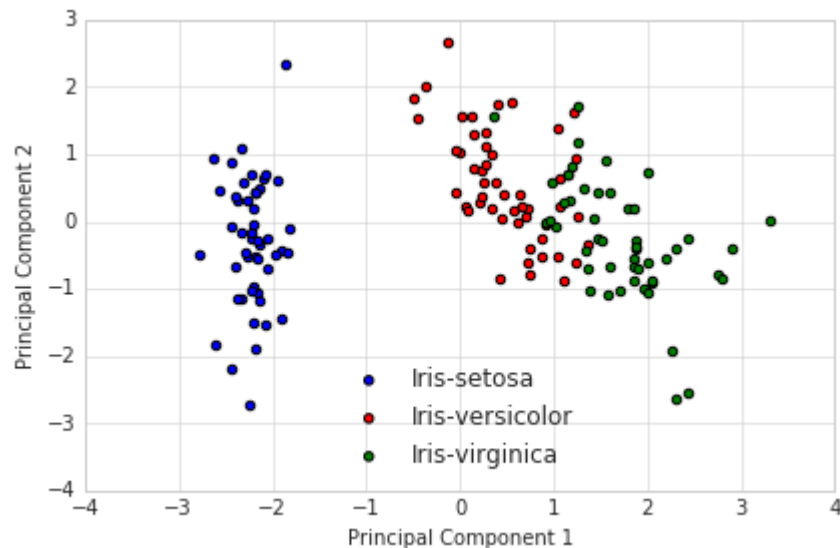
3 - Projection Onto the New Feature Space

In this last step we will use the $4 \times 24 \times 2$ -dimensional projection matrix WW to transform our samples onto the new subspace via the equation

$Y = X \times WY = X \times W$, where Y is a 150×2 matrix of our transformed samples.

```
Y = X_std.dot(matrix_w)
```

```
with plt.style.context('seaborn-whitegrid'):
    plt.figure(figsize=(6, 4))
    for lab, col in zip(('Iris-setosa', 'Iris-versicolor', 'Iris-virginica'),
                        ('blue', 'red', 'green')):
        plt.scatter(Y[y==lab, 0],
                    Y[y==lab, 1],
                    label=lab,
                    c=col)
    plt.xlabel('Principal Component 1')
    plt.ylabel('Principal Component 2')
    plt.legend(loc='lower center')
    plt.tight_layout()
    plt.show()
```



Now, what we got after applying the linear PCA transformation is a lower dimensional subspace (from 3D to 2D in this case), where the samples are “most spread” along the new feature axes.

Shortcut - PCA in scikit-learn

For educational purposes, we went a long way to apply the PCA to the Iris dataset. But luckily, there is already implementation in scikit-learn.

```

from sklearn.decomposition import PCA as sklearnPCA
sklearn_pca = sklearnPCA(n_components=2)
Y_sklearn = sklearn_pca.fit_transform(X_std)

with plt.style.context('seaborn-whitegrid'):
    plt.figure(figsize=(6, 4))
    for lab, col in zip(('Iris-setosa', 'Iris-versicolor', 'Iris-virginica'),
                        ('blue', 'red', 'green')):
        plt.scatter(Y_sklearn[y==lab, 0],
                    Y_sklearn[y==lab, 1],
                    label=lab,
                    c=col)
    plt.xlabel('Principal Component 1')
    plt.ylabel('Principal Component 2')
    plt.legend(loc='lower center')
    plt.tight_layout()
    plt.show()

```

