

# 10\_pca

February 16, 2024

## 1 Principal Component Analysis

*variationalform* <https://variationalform.github.io/>

*Just Enough: progress at pace* <https://variationalform.github.io/>

<https://github.com/variationalform>

Simon Shaw <https://www.brunel.ac.uk/people/simon-shaw>.

This work is licensed under CC BY-SA 4.0 (Attribution-ShareAlike 4.0 International)

Visit <http://creativecommons.org/licenses/by-sa/4.0/> to see the terms.

This document uses python

and also makes use of LaTeX

in Markdown

### 1.1 What this is about:

- **Principal Component Analysis**, or PCA. What it is, mathematically and in code.
- How it works, with examples.
- The connection between PCA and eigenvalues, and the SVD (**Singular Value Decomposition**).

As usual our emphasis will be on *doing* rather than *proving*: *just enough: progress at pace*

### 1.2 Assigned Reading

For this worksheet you are recommended Chapters 4 and 10 of [MML], Chapter 10 of [MLFCES], Chapter 5.3 of [IPDS],

- MML: Mathematics for Machine Learning, by Marc Peter Deisenroth, A. Aldo Faisal, and Cheng Soon Ong. Cambridge University Press. <https://mml-book.github.io>.
- MLFCES: Machine Learning: A First Course for Engineers and Scientists, by Andreas Lindholm, Niklas Wahlström, Fredrik Lindsten, Thomas B. Schön. Cambridge University Press. <http://smlbook.org>.
- IPDS: Introduction to Probability for Data Science, by Stanley H. Chan, <https://probability4datascience.com>

These can be accessed legally and without cost.

There are also these useful references for coding:

- PT: python: <https://docs.python.org/3/tutorial>
- NP: numpy: <https://numpy.org/doc/stable/user/quickstart.html>
- MPL: matplotlib: <https://matplotlib.org>

### 1.3 Context

We have seen these

- Eigenvalue decomposition
- SVD, the **Singular Value Decomposition**

Let's review them...

### 1.4 Eigen-systems of Symmetric Matrices

Given a real square symmetric  $n$ -row by  $n$ -column matrix,  $\mathbf{A} \in \mathbb{R}^{n \times n}$ , the eigenvalue problem is that of finding scalar eigenvalues  $\lambda$  and  $n$ -dimensional eigenvectors  $\mathbf{v}$  such that

$$\mathbf{A}\mathbf{v} = \lambda\mathbf{v} \implies \mathbf{A}\mathbf{V} = \mathbf{V}\mathbf{D} \implies \mathbf{A} = \sum_{k=1}^n \lambda_k \mathbf{v}_k \mathbf{v}_k^T.$$

The eigensystem is **real**.

We have the *Spectral Theorem* - see [MML, Theorem 4.15]

**Spectral Theorem (for matrices)** If  $\mathbf{A}$  is real and symmetric then its eigenvalues are all real and its eigenvector matrix  $\mathbf{V}$  can be taken as *orthogonal* so that  $\mathbf{V}^{-1} = \mathbf{V}^T$ . Hence...

$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^T$$

### 1.5 The SVD: Singular Value Decomposition

Given a real  $m$ -row by  $n$ -column matrix,  $\mathbf{B} \in \mathbb{R}^{m \times n}$

$$\mathbf{B} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \sum_{j=1}^p \sigma_j \mathbf{u}_j \mathbf{v}_j^T$$

where: for the left singular vectors:  $\mathbf{U} \in \mathbb{R}^{m \times m}$ ; for the singular values:  $\mathbf{\Sigma} \in \mathbb{R}^{m \times n}$ ; and, for the right singular vectors,  $\mathbf{V} \in \mathbb{R}^{n \times n}$ . Here  $p = \min\{m, n\}$ .

Note that  $\mathbf{\Sigma} = \text{diag}(\sigma_1, \dots, \sigma_p) + \text{zeros}$ , and we can always arrange that  $0 \leq \sigma_p \leq \dots \leq \sigma_1$ .

As  $\mathbf{B}$  is real,  $\mathbf{U}$  and  $\mathbf{V}$  are real and *orthogonal*.

If  $\sigma_r \neq 0$  and  $\sigma_p = 0$  for all  $p > r$  then  $r$  is the rank of  $\mathbf{B}$ .

## 1.6 How are these factorizations connected?

On the face of it they are very different. The first applies only to square symmetric matrices, while the second applies also to rectangular, and hence (why?) non-symmetric matrices.

But... Look at this... Given the SVD  $\mathbf{B} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$  we have,

$$\mathbf{B}^T\mathbf{B} = (\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$$

and remembering that, in general,  $(\mathbf{KL})^T = \mathbf{L}^T\mathbf{K}^T$  (this could be called *taking the transpose through*), we can write,

$$\mathbf{B}^T\mathbf{B} = \mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^T\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T = \mathbf{V}\mathbf{\Sigma}^T\mathbf{\Sigma}\mathbf{V}^T$$

because  $\mathbf{U}^T\mathbf{U} = \mathbf{I}$  (orthogonal).

Similarly, because also  $\mathbf{V}^T\mathbf{V} = \mathbf{I}$  (orthogonal),

$$\mathbf{B}\mathbf{B}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T(\mathbf{U}\mathbf{\Sigma}\mathbf{V}^T)^T = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T\mathbf{V}\mathbf{\Sigma}^T\mathbf{U}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^T\mathbf{U}^T.$$

Do you recognise these?

We have just shown that,

$$\mathbf{B}^T\mathbf{B} = \mathbf{V}\mathbf{\Sigma}^T\mathbf{\Sigma}\mathbf{V}^T \quad \text{and} \quad \mathbf{B}\mathbf{B}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^T\mathbf{U}^T.$$

Familiar? Think about  $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^T$ .

- Put  $\mathbf{A} = \mathbf{B}^T\mathbf{B}$  (symmetric) and  $\mathbf{D} = \mathbf{\Sigma}^T\mathbf{\Sigma}$ . Then,

$$\mathbf{B}^T\mathbf{B} = \mathbf{V}\mathbf{\Sigma}^T\mathbf{\Sigma}\mathbf{V}^T \quad \text{becomes} \quad \mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{V}^T.$$

- Put  $\mathbf{A} = \mathbf{B}\mathbf{B}^T$  (symmetric) and  $\mathbf{D} = \mathbf{\Sigma}\mathbf{\Sigma}^T$ . Then,

$$\mathbf{B}\mathbf{B}^T = \mathbf{U}\mathbf{\Sigma}\mathbf{\Sigma}^T\mathbf{U}^T \quad \text{becomes} \quad \mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^T.$$

- $\mathbf{V}$ , the right singular vectors in the SVD are the eigenvectors of  $\mathbf{B}^T\mathbf{B}$ .
- $\mathbf{U}$ , the left singular vectors in the SVD are the eigenvectors of  $\mathbf{B}\mathbf{B}^T$ .
- In both cases  $\mathbf{\Sigma}$  contains the positive square roots of the eigenvalues of  $\mathbf{B}^T\mathbf{B}$  and  $\mathbf{B}\mathbf{B}^T$ .
- **NOTE:**  $\mathbf{B}^T\mathbf{B}$  and  $\mathbf{B}\mathbf{B}^T$  have the same non-zero eigenvalues (same rank).

## 1.7 Why does this matter?

Our data,  $\mathbf{X}$ , is organized into rows of feature values with one observation per row and one feature per column. We write this as

$$\mathbf{X} = (\mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_D)$$

If  $D = 3$  (four features)...

... we recall that the **covariance matrix** takes this form:

$$\mathbf{S} = \begin{pmatrix} \text{Var}(X_0) & \text{Cov}(X_0, X_1) & \text{Cov}(X_0, X_2) & \text{Cov}(X_0, X_3) \\ \text{Cov}(X_1, X_0) & \text{Var}(X_1) & \text{Cov}(X_1, X_2) & \text{Cov}(X_1, X_3) \\ \text{Cov}(X_2, X_0) & \text{Cov}(X_2, X_1) & \text{Var}(X_2) & \text{Cov}(X_2, X_3) \\ \text{Cov}(X_3, X_0) & \text{Cov}(X_3, X_1) & \text{Cov}(X_3, X_2) & \text{Var}(X_3) \end{pmatrix}$$

because  $\text{Cov}(X, X) = \text{Var}(X)$ . Since  $\text{Cov}(X, Y) = \text{Cov}(Y, X)$ , this matrix is **symmetric** and so has real eigenvalues.

We have seen that if the data are already centred then,

$$(N-1)\mathbf{S} = \begin{pmatrix} \mathbf{X}_0 \cdot \mathbf{X}_0 & \mathbf{X}_0 \cdot \mathbf{X}_1 & \mathbf{X}_0 \cdot \mathbf{X}_2 & \mathbf{X}_0 \cdot \mathbf{X}_3 \\ \mathbf{X}_1 \cdot \mathbf{X}_0 & \mathbf{X}_1 \cdot \mathbf{X}_1 & \mathbf{X}_1 \cdot \mathbf{X}_2 & \mathbf{X}_1 \cdot \mathbf{X}_3 \\ \mathbf{X}_2 \cdot \mathbf{X}_0 & \mathbf{X}_2 \cdot \mathbf{X}_1 & \mathbf{X}_2 \cdot \mathbf{X}_2 & \mathbf{X}_2 \cdot \mathbf{X}_3 \\ \mathbf{X}_3 \cdot \mathbf{X}_0 & \mathbf{X}_3 \cdot \mathbf{X}_1 & \mathbf{X}_3 \cdot \mathbf{X}_2 & \mathbf{X}_3 \cdot \mathbf{X}_3 \end{pmatrix} = \begin{pmatrix} \mathbf{X}_0^T \\ \mathbf{X}_1^T \\ \mathbf{X}_2^T \\ \mathbf{X}_3^T \end{pmatrix} (\mathbf{X}_0 \ \mathbf{X}_1 \ \mathbf{X}_2 \ \mathbf{X}_3)$$

and, hence (in general), the (sample) covariance matrix for  $N$  observations is

$$\mathbf{S} = \frac{1}{(N-1)} \mathbf{X}^T \mathbf{X}.$$

## 1.8 Terminology

We just introduced the **sample** covariance matrix:

$$\mathbf{S} = \frac{1}{(N-1)} \mathbf{X}^T \mathbf{X}.$$

The  $N-1$  in the denominator makes this an **unbiased** estimate of the population statistics. When  $N$  is large we can just work with

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$$

and call it the **empirical** covariance matrix.

This terminology is discussed in [MML, Section 6.4.2].

## 1.9 Conventions

We have now adopted a convention that our data matrix  $\mathbf{X}$  has features varying along the rows, and observations varying down the columns so that:

$$\mathbf{X} = \left( \mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_D \right)$$

gives a data set with  $D + 1$  features. The length (they must all be the same) of the column vectors  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_D$  tell us how many observations there are. We've been denoting this by  $N$ .

**HOWEVER:** in some sources this convention is transposed. The different features occupy their own rows of the matrix, with the observations recorded along the rows.

This is the case in [MML]. It means that

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad \text{for us, becomes} \quad \mathbf{S} = \frac{1}{N} \mathbf{X} \mathbf{X}^T \quad \text{for them}$$

because our  $\mathbf{X}$  is their  $\mathbf{X}^T$ .

**BE CAREFUL: this is not uncommon**

## 1.10 Features and Observations

Let's say there are  $D + 1$  features (columns) in our data set  $\mathbf{X}$  and  $N$  observations (rows).

An observation takes the form  $\mathbf{x}_j = (x_1, x_2, \dots, x_d)^T$ , a **column vector**, for  $j = 1, 2, \dots, N$ .

Hence,

$$\mathbf{X} = \left( \mathbf{X}_0, \mathbf{X}_1, \dots, \mathbf{X}_D \right) = \left( \mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_N \right)^T$$

Note: we are using **upper case**  $\mathbf{X}_k$  for a column vector of observations of a feature (in the column indexed by  $k$ ), and **lower case**,  $\mathbf{x}_j$ , for a feature vector arising from a single observation (in the row indexed by  $j$ ).

We only use column vectors in these notes.

## 1.11 PCA - Principal Component Analysis

The main idea and motivation behind this is that high dimensional data often lives very close to a lower dimensional subspace.

A typical and often used example of this is that in Figure 10.1 of [MML, Chap. 10].

We can see that here: <https://mml-book.github.io>.

PCA will analyze a data set and determine the direction in which most variation occurs. If we are to approximate using a lower dimensional space then this is a good direction (subspace component) to start with.

Technically, PCA determines directions which maximize variance.

Let's go through this slowly - it can be quite confusing.

### 1.12 PCA - outline algorithm.

- Take our  $D$ -column by  $N$ -row data set  $\mathbf{X}$  and ensure that the column means are zero.
- This is referred to as **centering the data**.
- It means that  $\mathbb{E}(\mathbf{X}_d) = 0$  for columns  $d = 0, 1, 2, \dots, D - 1$ .
- We want a  $D$ -row by  $M$ -column matrix  $\mathbf{B} \in \mathbb{R}^{D,M}$ , called the **code** in [MML], such that we can define  $\mathbf{Z}$  as follows:

$$\mathbf{Z} = \mathbf{X}\mathbf{B}\mathbf{B}^T$$

- If  $M = D$  we insist that  $\mathbf{Z} = \mathbf{X}$ . Otherwise we have  $M < D$  and  $\mathbf{Z}$  is an approximation (a projection) of  $\mathbf{X}$  in a lower dimensional subspace.
- We determine  $\mathbf{B}$  by minimizing the reconstruction error:

$$\mathcal{J} = \frac{1}{N} \sum_{n=1}^N \|\mathbf{x}_n - \mathbf{z}_n\|_2^2$$

where  $\mathbf{z}_n$  is the  $n$ -th row of  $\mathbf{Z}$ .

## 2 Some Technicalities

**THINK ABOUT:** If  $\mathbf{X} = \mathbf{X}\mathbf{B}\mathbf{B}^T$  when  $M = D$  then  $\mathbf{B}\mathbf{B}^T = \mathbf{I}_D$ . Is  $\mathbf{B}$  square?

**THINK ABOUT:** If  $\mathbf{Z} = \mathbf{X}\mathbf{B}\mathbf{B}^T$  when  $M < D$  then what shape is  $\mathbf{Z}$ ? Is it the same shape as  $\mathbf{X}$ ?

**THINK ABOUT:** If we set  $\mathbf{Y} = \mathbf{X}\mathbf{B}$  then  $\mathbf{Z} = \mathbf{Y}\mathbf{B}^T$ . What shape is  $\mathbf{Y}$ ?

**THINK ABOUT:**  $\mathbf{Y}$  will be  $N$ -rows by  $M$ -columns.  $\mathbf{Y}$  is **smaller** than  $\mathbf{X}$  if  $M < D$  and represents dimensionality reduction.

**THINK ABOUT:** the reduction  $\mathbf{Y} = \mathbf{X}\mathbf{B}$ , and the subsequent enlargement  $\mathbf{Z} = \mathbf{Y}\mathbf{B}^T$  is the basis of an **autoencoder**.  $\mathbf{Z}$  is a reconstruction of  $\mathbf{X}$  resulting from a data compression step.

This minimization referred to above is a long technical excursion in multivariable calculus. The important result of it is that we need to find the eigensystem of the empirical data covariance matrix,

$$\mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X}.$$

This means we want to solve

$$\mathbf{S}\mathbf{v} = \lambda\mathbf{v}$$

for the eigenpairs  $(\lambda_1, \mathbf{v}_1)$ ,  $(\lambda_2, \mathbf{v}_2)$ , ...

Then  $\mathbf{B} = (\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_M)$  and the eigenvalues tell us how much variance of the original data set is captured by the  $M$ -dimensional projection.

Let's see this in action...

## 2.1 Worked example

Consider this set of data (already centered),

$$\mathbf{X} = \begin{pmatrix} 1 & 2 \\ 2 & 1 \\ -2 & -1 \\ -1 & -2 \end{pmatrix} \implies \mathbf{S} = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{4} \begin{pmatrix} 10 & 8 \\ 8 & 10 \end{pmatrix}$$

The eigensystem, with  $\mathbf{S}\mathbf{V} = \mathbf{V}\mathbf{D}$ , is

$$\mathbf{V} = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 & -1 \\ 1 & 1 \end{pmatrix} = (\mathbf{v}_0 \quad \mathbf{v}_1) \quad \text{and} \quad \mathbf{D} = \frac{1}{2} \begin{pmatrix} 9 & 0 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} \lambda_0 & 0 \\ 0 & \lambda_1 \end{pmatrix}.$$

Let's see this in python...

```
[1]: import numpy as np
import matplotlib.pyplot as plt
# set up the feature matrix and check the column means are zero
D=2; N=4
X = np.array([[1, 2], [2, 1], [-2, -1], [-1, -2]])
print(f'Column means: col 1, {X[:,0].mean()} and col 2, {X[:,1].mean()}')
# and the empirical covariance matrix
S = 1/N*X.T @ X
print('S = \n',S)
# solve the eigenvalue problem
lmda, V = np.linalg.eig(S)
print('evals = ', lmda)
print('evecs = \n', V)
```

Column means: col 1, 0.0 and col 2, 0.0

```
S =
[[2.5 2. ]
 [2.  2.5]]
evals = [4.5 0.5]
evecs =
[[ 0.70710678 -0.70710678]
 [ 0.70710678  0.70710678]]
```

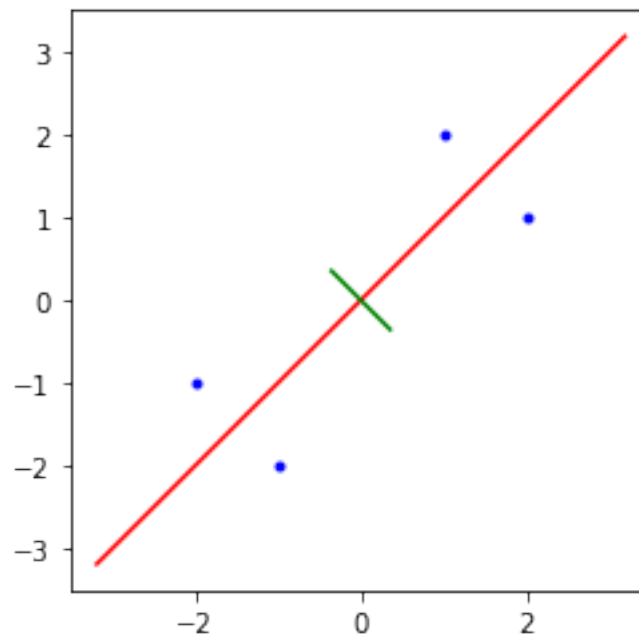
```
[2]: # a picture will tell us much more...
plt.figure(figsize=(4,4)); plt.gca().set_aspect('equal')
# plot the data in blue
```

```

plt.plot(X[:,0], X[:,1], '.', color='b')
# plot the eigenvectors...
# the first in red from -v0 to +v0 with length 2*lambda_0
x0 = lmda[0]*V[0,[0]]; y0 = lmda[0]*V[1,[0]]
plt.plot([-x0,x0],[-y0,y0], '-', color='r')
# the second in green from -v1 to +v1 with length 2*lambda_1
x1 = lmda[1]*V[0,[1]]; y1 = lmda[1]*V[1,[1]]
plt.plot([-x1,x1],[-y1,y1], '-', color='g')

```

[2]: [[matplotlib.lines.Line2D at 0x7f9fd80ce588](#)]



We see that the direction of maximum variance is given by the dominant eigenpair. The next eigenpair is **orthogonal**.

Now, we know that  $\mathbf{V}$  is an orthogonal matrix, so that  $\mathbf{V}\mathbf{V}^T = \mathbf{I}$ .

It is therefore clear that  $\mathbf{X} = \mathbf{X}\mathbf{V}\mathbf{V}^T$ . Now, with  $\mathbf{V} = (\mathbf{v}_0, \mathbf{v}_1)$ , we observe that

$$\mathbf{Z}_0 = \mathbf{X}\mathbf{v}_0\mathbf{v}_0^T = \frac{3}{2} \begin{pmatrix} 1 & 1 \\ 1 & 1 \\ -1 & -1 \\ -1 & -1 \end{pmatrix} \quad \text{and} \quad \mathbf{Z}_1 = \mathbf{X}\mathbf{v}_1\mathbf{v}_1^T = \frac{1}{2} \begin{pmatrix} -1 & 1 \\ 1 & -1 \\ -1 & 1 \\ 1 & -1 \end{pmatrix}$$

The rows give us the projections of the original rows (features) onto the lower dimensional subspaces.

Let's see it in code, and then in pictures (building on the picture above)...



```
[3]: print('X - X V V.T = \n', X - X @ V @ V.T )
```

```
v0 = V[:,[0]]; Z0 = X @ v0 @ v0.T
```

```
print('X @ v0 @ v0.T = \n', Z0)
```

```
v1 = V[:,[1]]; Z1 = X @ v1 @ v1.T
```

```
print('X @ v1 @ v1.T = \n', Z1)
```

```
X - X V V.T =
```

```
[[ 1.11022302e-16  2.22044605e-16]
```

```
 [ 4.44089210e-16  2.22044605e-16]
```

```
[-4.44089210e-16 -2.22044605e-16]
```

```
[-1.11022302e-16 -2.22044605e-16]]
```

```
X @ v0 @ v0.T =
```

```
[[ 1.5  1.5]
```

```
 [ 1.5  1.5]
```

```
[-1.5 -1.5]
```

```
[-1.5 -1.5]]
```

```
X @ v1 @ v1.T =
```

```
[[ -0.5  0.5]
```

```
 [ 0.5 -0.5]
```

```
[-0.5  0.5]
```

```
 [ 0.5 -0.5]]
```

```
[4]: plt.figure(figsize=(3,3)); plt.gca().set_aspect('equal')
```

```
plt.plot(X[:,0], X[:,1], '.', color='b')
```

```
plt.plot([-x0,x0],[-y0,y0], '-', color='r')
```

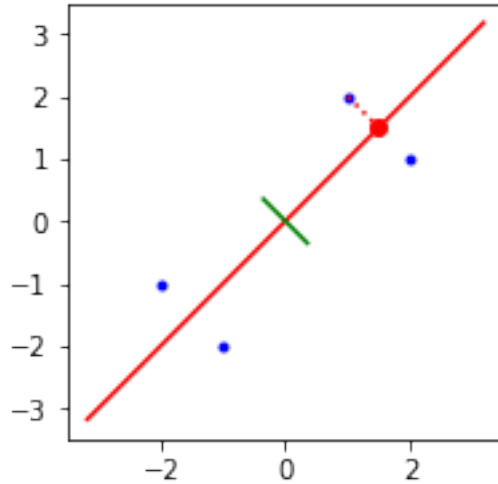
```
plt.plot([-x1,x1],[-y1,y1], '-', color='g')
```

```
# just the first row of Z0 for the moment
```

```
plt.plot([X[0,0], Z0[0,0]], [X[0,1], Z0[0,1]], ':', marker='o', color='r',
```

```
↪markevery=[1])
```

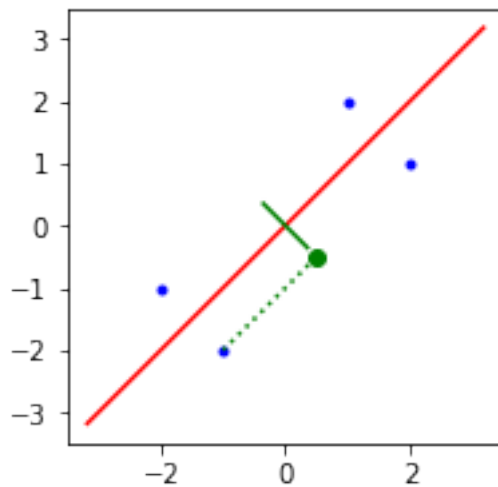
```
[4]: [<matplotlib.lines.Line2D at 0x7f9fc008e3c8>]
```



In the first row of  $\mathbf{Z}_0 = \mathbf{X}\mathbf{v}_0\mathbf{v}_0^T$ , the original point at (1, 2) is projected to the new point (1.5, 1.5) on the dominant lower dimensional subspace. What about  $\mathbf{Z}_1 = \mathbf{X}\mathbf{v}_1\mathbf{v}_1^T$ ?

```
[5]: plt.figure(figsize=(3,3)); plt.gca().set_aspect('equal')
plt.plot(X[:,0], X[:,1], '.', color='b')
plt.plot([-x0,x0],[-y0,y0], '-', color='r')
plt.plot([-x1,x1],[-y1,y1], '-', color='g')
# just the last row of Z1 for the moment
plt.plot([X[3,0], Z1[3,0]], [X[3,1], Z1[3,1]], ':', marker='o', color='g',
        ↳markevery=[1])
```

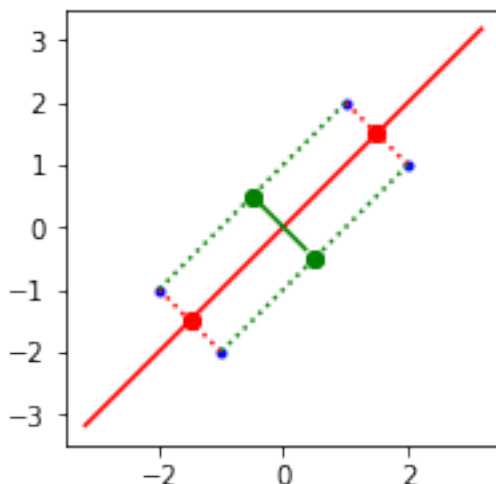
[5]: [[matplotlib.lines.Line2D](#) at 0x7fa0002e5240>]



Now in the last row of  $\mathbf{Z}_1 = \mathbf{X}\mathbf{v}_1\mathbf{v}_1^T$ , the original point at  $(-1, -2)$  is projected to the new point  $(0.5, -0.5)$  on the next-dominant lower dimensional subspace.

Let's see all the projections in one picture.

```
[6]: plt.figure(figsize=(3,3)); plt.gca().set_aspect('equal')
plt.plot(X[:,0], X[:,1], '.', color='b')
plt.plot([-x0,x0],[-y0,y0], '-', color='r')
plt.plot([-x1,x1],[-y1,y1], '-', color='g')
for k in range(4):
    plt.plot([X[k,0], Z0[k,0]], [X[k,1], Z0[k,1]], ':', marker='o', color='r',
             ↳markevery=[1])
    plt.plot([X[k,0], Z1[k,0]], [X[k,1], Z1[k,1]], ':', marker='o', color='g',
             ↳markevery=[1])
```



## 2.2 The Local Coordinate System

If we treat the eigenvectors as subspaces then the length along each eigenvector is the local coordinate in that subspace. How can we get these coordinates? Well, look at this:

$$\mathbf{X} = \begin{pmatrix} 1 & 2 \\ 2 & 1 \\ -2 & -1 \\ -1 & -2 \end{pmatrix} \implies \mathbf{X}\mathbf{v}_0 = \frac{1}{\sqrt{2}} \begin{pmatrix} 3 \\ 3 \\ -3 \\ -3 \end{pmatrix} \quad \text{and} \quad \mathbf{X}\mathbf{v}_1 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ -1 \\ 1 \\ -1 \end{pmatrix}.$$

These tell us that the first two points (rows) in  $\mathbf{X}$  project to coincident points a distance  $\frac{3}{\sqrt{2}}$  along the dominant eigenvector, and that the second two points project to coincident points a distance  $\frac{-3}{\sqrt{2}}$  along the dominant eigenvector.

On the other hand, the first and third points project to distances  $\frac{1}{\sqrt{2}}$  along the second eigenvector, while the second and fourth project to distances  $\frac{-1}{\sqrt{2}}$ .

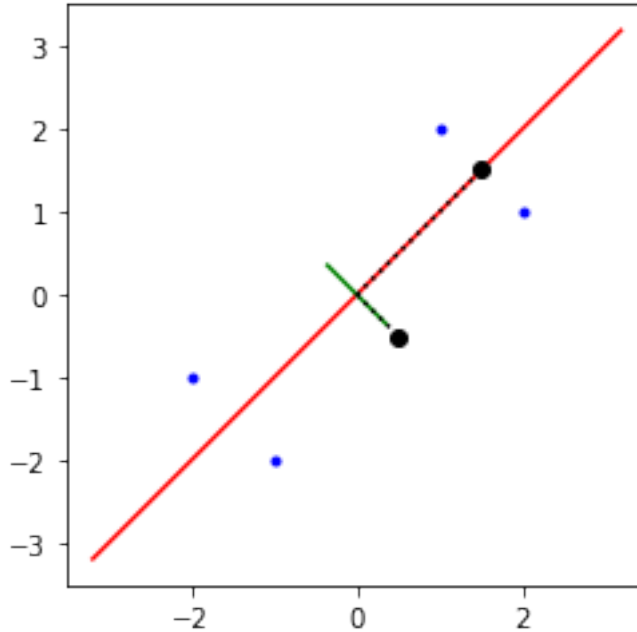
Here is some code to illustrate this...

```
[7]: Xv0 = X @ v0
      Xv1 = X @ v1
      # multiply by root 2 to tidy up the output
      print('sqrt(2) X v0 = \n', np.sqrt(2) * Xv0)
      print('sqrt(2) X v1 = \n', np.sqrt(2) * Xv1)
```

```
sqrt(2) X v0 =
[[ 3.]
 [ 3.]
 [-3.]
 [-3.]]
sqrt(2) X v1 =
[[ 1.]
 [-1.]
 [ 1.]
 [-1.]]
```

```
[8]: plt.figure(figsize=(4,4)); plt.gca().set_aspect('equal')
      plt.plot(X[:,0], X[:,1], '.', color='b')
      plt.plot([-x0,x0],[-y0,y0], '-', color='r')
      plt.plot([-x1,x1],[-y1,y1], '-', color='g')
      cos45 = sin45 = 1/np.sqrt(2)
      # for data point 0 along v0
      plt.plot([0, Xv0[0,0]*cos45], [0, Xv0[0,0]*sin45], ':', marker='o', color='k',
               ↳markevery=[1])
      # for data point 3 along v1
      plt.plot([0, -Xv1[3,0]*cos45], [0, Xv1[3,0]*sin45], ':', marker='o', color='k',
               ↳markevery=[1])
```

```
[8]: [<matplotlib.lines.Line2D at 0x7fa00030ee10>]
```



### 2.3 Explained Variance

There is yet more to see... The column-wise variances in the original data set,

$$\mathbf{X} = \begin{pmatrix} 1 & 2 \\ 2 & 1 \\ -2 & -1 \\ -1 & -2 \end{pmatrix}$$

are

$$\begin{cases} \text{Var}(\mathbf{X}_0) = \mathbb{E}(\mathbf{X}_0 \cdot \mathbf{X}_0) = \frac{1}{4}(1 + 2^2 + (-2)^2 + 1) = \frac{5}{2}, \\ \text{Var}(\mathbf{X}_1) = \mathbb{E}(\mathbf{X}_1 \cdot \mathbf{X}_1) = \frac{1}{4}(2^2 + 1 + 1 + (-2)^2) = \frac{5}{2} \end{cases}$$

and the total variance in the data set is  $\frac{5}{2} + \frac{5}{2} = 5$  (Note - this isn't what you get by stacking the data and taking a single variance because the mean(s) may get altered).

The sum of the individual feature variances is  $\frac{5}{2} + \frac{5}{2} = 5$ .

In PCA, each eigenvalue gives the variance in the direction of its eigenvector.

Our eigenvalues were  $\lambda_0 = \frac{9}{2}$  and  $\lambda_1 = \frac{1}{2}$ . The total variance is therefore  $\frac{9}{2} + \frac{1}{2} = 5$ .

**THINK ABOUT:** the trace of a matrix is the sum of its eigenvalues. Relevance?

We talk about each eigenvalue **explaining** variance in the original data set.

Here the first eigenvalue explains  $\frac{9}{2} \div 5 = 90\%$  of the original variance. The remaining 5% is in the orthogonal direction of the second eigenvector.

## 2.4 The connection to SVD

We just performed PCA using an eigenvalue analysis of the empirical covariance matrix

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} \quad \text{leading to} \quad \frac{1}{N} \mathbf{X}^T \mathbf{X} \mathbf{v} = \lambda \mathbf{v}.$$

Earlier, given the SVD  $\mathbf{B} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ , we saw that putting  $\mathbf{A} = \mathbf{B}^T \mathbf{B}$  (symmetric) and  $\mathbf{D} = \mathbf{\Sigma}^T \mathbf{\Sigma}$  led to

$$\mathbf{B}^T \mathbf{B} = \mathbf{V} \mathbf{\Sigma}^T \mathbf{\Sigma} \mathbf{V}^T \quad \text{becoming} \quad \mathbf{A} = \mathbf{V} \mathbf{D} \mathbf{V}^T.$$

Therefore, for PCA we could also obtain the SVD of  $\mathbf{X}$  and use the right singular vectors. The eigenvalues will be the squares of the singular values divided by  $N$ .

Let's check this:

```
[9]: # re-solve the eigenvalue problem
lmda, V = np.linalg.eig(S)
print(f'evals = {lmda} and V = ')
print(V)
# take the SVD of X
U, Sig, VT = np.linalg.svd(X)
print(f'singular values Sigma = {Sig}')
print(f'Sigma^2/N = {Sig*Sig/N} and V = ')
print(VT.T)
```

```
evals = [4.5 0.5] and V =
[[ 0.70710678 -0.70710678]
 [ 0.70710678  0.70710678]]
singular values Sigma = [4.24264069 1.41421356]
Sigma^2/N = [4.5 0.5] and V =
[[-0.70710678 -0.70710678]
 [-0.70710678  0.70710678]]
```

Let's see how to do PCA with sklearn...

```
[10]: from sklearn.decomposition import PCA
pca = PCA(n_components=2)
pca.fit(X)
print('Here is the explained variance as ratios...')
print(f'XV ratio = {pca.explained_variance_ratio_}')
print(f'Here are the singular values = {pca.singular_values_}')
print('the eigenvalues are squares of the singular values divided by N')
print(f'eigenvalues = {pca.singular_values_**2/N}')
```

Here is the explained variance as ratios...

XV ratio = [0.9 0.1]

Here are the singular values = [4.24264069 1.41421356]

the eigenvalues are squares of the singular values divided by N

eigenvalues = [4.5 0.5]

The principal components can be compared with the eigenvectors from above.

Beware: they are given to us in **rows**.

```
[11]: # the principal components are stored as row vectors, so transpose
      B = pca.components_.T
      print('Principal Components (transposed) B = \n', B)
      print('Compare with our earlier V = \n', V)
```

Principal Components (transposed) B =

```
[[ 0.70710678  0.70710678]
```

```
 [ 0.70710678 -0.70710678]]
```

Compare with our earlier V =

```
[[ 0.70710678 -0.70710678]
```

```
 [ 0.70710678  0.70710678]]
```

Do you remember that above we used  $Xv_0$  and  $Xv_1$  to get the positions (lengths) along the principal axes of the projected data points?

sklearn can do this for us:

```
[12]: # multiply by sqrt(2) to tidy the output.
      Xf = pca.fit_transform(X)
      print('np.sqrt(2)*Xf = \n', np.sqrt(2)*Xf)
```

np.sqrt(2)\*Xf =

```
[[ 3. -1.]
```

```
 [ 3.  1.]
```

```
[-3. -1.]
```

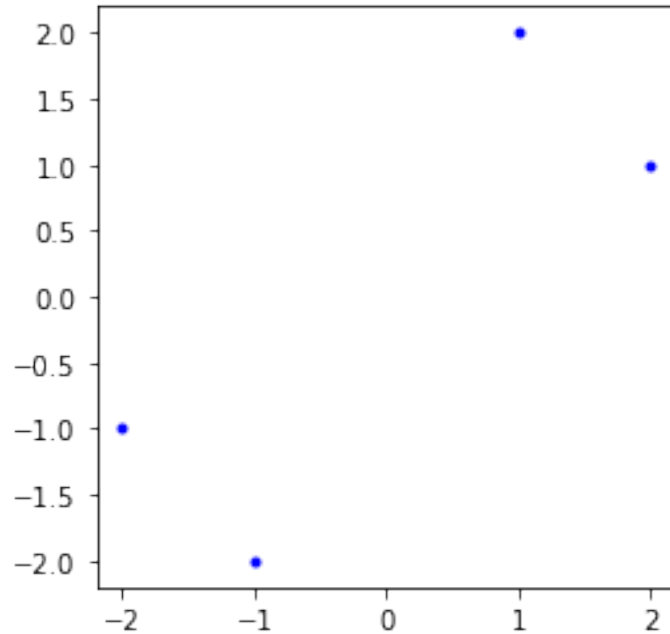
```
[-3.  1.]
```

Technically, the principal components give us a new **basis** for the data space. These **transformed** coordinates gives us the coordinates of the data in the new basis. Let's see this in pictures...

Here are the original data points in the original feature basis. This is where each axis is labelled with the feature name.

```
[13]: plt.figure(figsize=(4,4)); plt.gca().set_aspect('equal')
      plt.plot(X[:,0], X[:,1], '.', color='b')
```

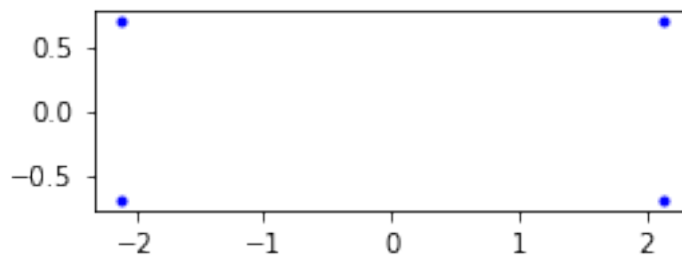
```
[13]: [<matplotlib.lines.Line2D at 0x7f9ff07e8eb8>]
```



Here are the data points in the PCA basis.

```
[14]: plt.figure(figsize=(4,4)); plt.gca().set_aspect('equal')
plt.plot(Xf[:,0], Xf[:,1], '.', color='b')
```

```
[14]: [<matplotlib.lines.Line2D at 0x7f9fd9716e48>]
```



- It is immediately apparent that the data has become more *one dimensional*.
- An issue though is that the axes are no longer easily interpreted.
- This is relevant to the **explainability agenda** in AI and Data Science.

## 2.5 Reflection

That was a long journey - and we didn't even derive the results, we just quoted and illustrated them. This, again, is because we are doing *just enough* to make *progress at pace*.



We're now going to embark on a much more realistic (well, in 2D at least) example of how this works. We'll go faster because all the work has been done.

## 2.6 Standard Example

The following example is used a lot in account of PCA.

The idea is to generate and plot a lozenge of Gaussian distributed data. It will have unequal variances (otherwise it would be a circle).

```
[15]: # generate this number of sample points
Ns=50
# The Gaussian lozenge will be centered with non-unit covariance
mean = [0, 0]
cov = [[30, 15], [15, 15]]
# generate Ns random points (x,y)
x, y = np.random.multivariate_normal(mean, cov, Ns).T
# reshape them to columns and stack them next to each other
X = np.hstack((x.reshape(-1,1),y.reshape(-1,1)))
# We can see it if the data matrix is small - otherwise little point
if Ns < 8: print(X)
```

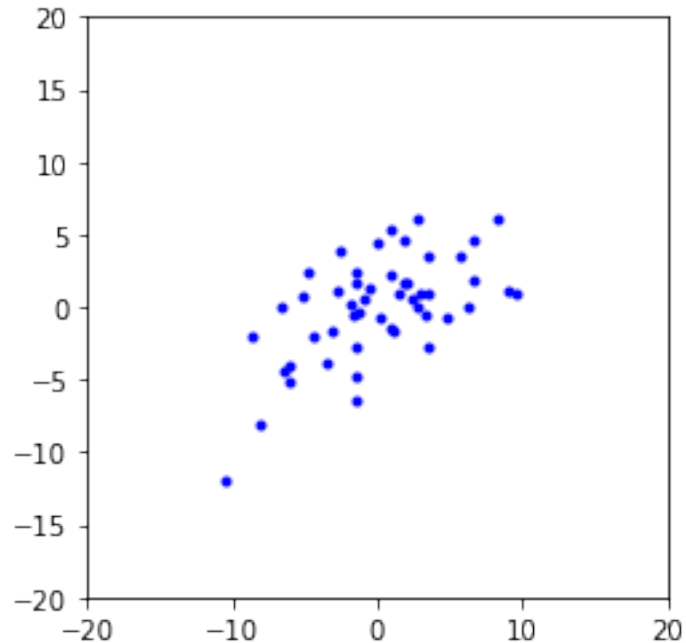
```
[16]: # let's look at the empirical column means - they wont be exactly zero
print(f'Column means 1, {X[:,0].mean()} and 2, {X[:,1].mean()}')
# so let's center this sample data
X[:,0] -= X[:,0].mean()
X[:,1] -= X[:,1].mean()
print(f'Centered column means 1, {X[:,0].mean()} and 2, {X[:,1].mean()}')
```

Column means 1, -1.128231495797423 and 2, -0.9316252071162511

Centered column means 1, -1.5987211554602254e-16 and 2, 1.509903313490213e-16

```
[17]: # let's plot our data set
plt.figure(figsize=(4,4))
plt.xlim(-20,20)
plt.ylim(-20,20)
plt.gca().set_aspect('equal')
plt.plot(X[:,0], X[:,1], '.', color='b')
```

```
[17]: [<matplotlib.lines.Line2D at 0x7f9fe08ee4a8>]
```



```
[18]: # perform the PCA
pca = PCA(n_components=2)
pca.fit(X)
print(f'XV ratio = {pca.explained_variance_ratio_}')
print(f'sing vals = {pca.singular_values_}')
# the component are stored as row vectors, so transpose
V = pca.components_.T
print('V = \n', V)
```

```
XV ratio = [0.81242882 0.18757118]
sing vals = [36.93179951 17.74560735]
V =
[[-0.84152666 -0.54021558]
 [-0.54021558  0.84152666]]
```

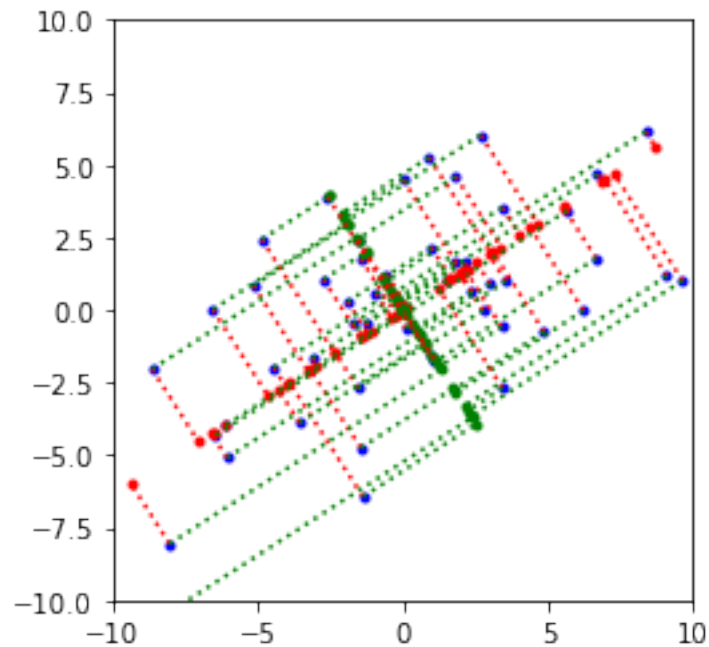
```
[19]: # project X to the singular components
Z1 = X @ V[:,[0]] @ V[:,[0]].T
Z2 = X @ V[:,[1]] @ V[:,[1]].T
```

```
[20]: # re-plot, and include all the projected data
plt.figure(figsize=(4,4))
plt.gca().set_aspect('equal')
plt.plot(X[:,0], X[:,1], '.', color='b')
# plot the projections - these illustrate the directions
plt.plot(Z1[:,0], Z1[:,1], '.', color='r')
```

```

plt.plot(Z2[:,0], Z2[:,1], '.', color='g')
# Now loop over each point and dot-line the projection onto v0 and v1
for k in range(Ns):
    plt.plot([X[k,0], Z1[k,0]], [X[k,1],Z1[k,1]], ':', color='r')
    plt.plot([X[k,0], Z2[k,0]], [X[k,1],Z2[k,1]], ':', color='g')
# zoom in or out with this...
window=10; plt.xlim(-window,window); plt.ylim(-window,window); plt.show()

```



```

[21]: # Let's see what the data looks like in the new coordinate system
Xf = pca.fit_transform(X)

```

```

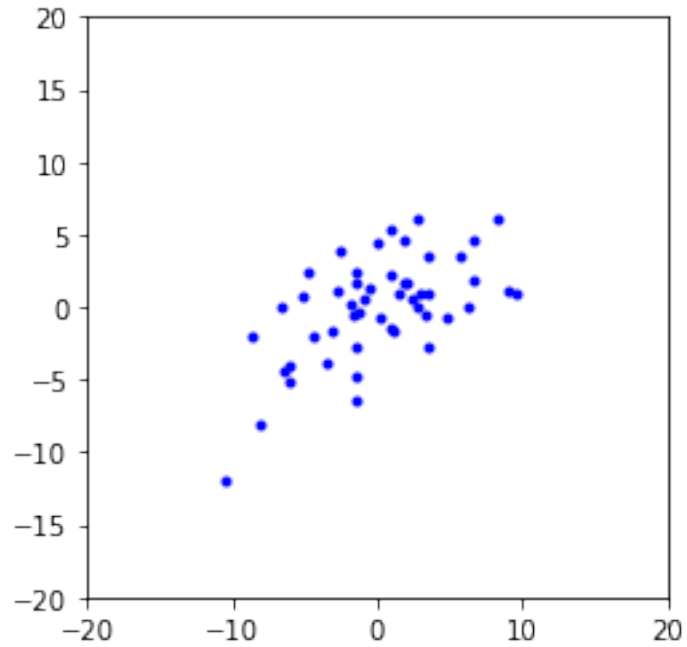
[22]: # here is the original...
plt.figure(figsize=(4,4))
plt.xlim(-20,20)
plt.ylim(-20,20)
plt.gca().set_aspect('equal')
plt.plot(X[:,0], X[:,1], '.', color='b')

```

```

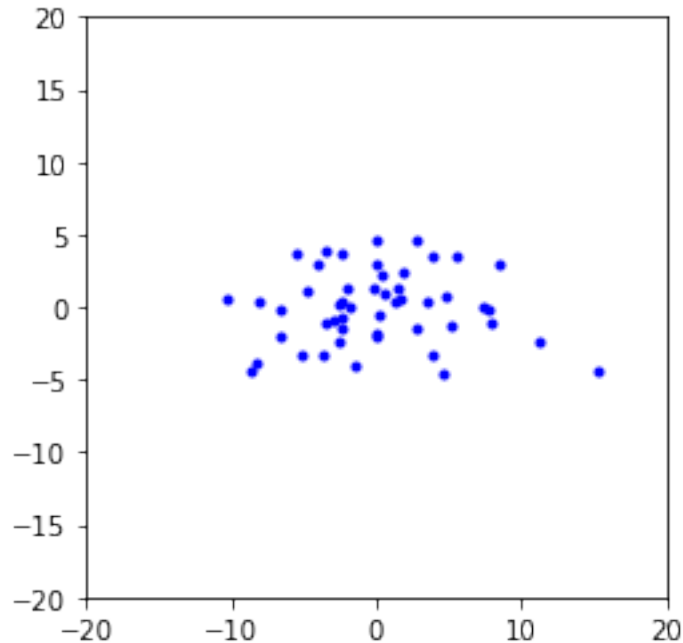
[22]: [<matplotlib.lines.Line2D at 0x7f9ff0856080>]

```



```
[23]: # here is the transformed data
plt.figure(figsize=(4,4))
plt.xlim(-20,20)
plt.ylim(-20,20)
plt.gca().set_aspect('equal')
plt.plot(Xf[:,0], Xf[:,1], '.', color='b')
```

```
[23]: [<matplotlib.lines.Line2D at 0x7f9fd133ee10>]
```



### 2.6.1 Review

We covered *just enough*, to make *progress at pace*. We looked at

- How the SVD and eigenvalue decomposition are related.
- How this becomes relevant to the data covariance matrix.
- PCA and its use in variance maximization.

Now we can start putting all of this material to work.

## 2.7 Technical Notes, Production and Archiving

Ignore the material below. What follows is not relevant to the material being taught.

### Production Workflow

- Finalise the notebook material above
- Clear and fresh run of entire notebook
- Create html slide show:
  - `jupyter nbconvert --to slides 10_pca.ipynb`
- Set `OUTPUTTING=1` below
- Comment out the display of web-sourced diagrams
- Clear and fresh run of entire notebook
- Comment back in the display of web-sourced diagrams
- Clear all cell output
- Set `OUTPUTTING=0` below
- Save

- git add, commit and push to FML
- copy PDF, HTML etc to web site
  - git add, commit and push
- rebuild binder

Some of this originated from

<https://stackoverflow.com/questions/38540326/save-html-of-a-jupyter-notebook-from-within-the-r>

These lines create a back up of the notebook. They can be ignored.

At some point this is better as a bash script outside of the notebook

```
[24]: %%bash
NBROOTNAME=10_pca
OUTPUTTING=1

if [ $OUTPUTTING -eq 1 ]; then
  jupyter nbconvert --to html $NBROOTNAME.ipynb
  cp $NBROOTNAME.html ../backups/$(date +"%m_%d_%Y-%H%M%S")_$NBROOTNAME.html
  mv -f $NBROOTNAME.html ../formats/html/

  jupyter nbconvert --to pdf $NBROOTNAME.ipynb
  cp $NBROOTNAME.pdf ../backups/$(date +"%m_%d_%Y-%H%M%S")_$NBROOTNAME.pdf
  mv -f $NBROOTNAME.pdf ../formats/pdf/

  jupyter nbconvert --to script $NBROOTNAME.ipynb
  cp $NBROOTNAME.py ../backups/$(date +"%m_%d_%Y-%H%M%S")_$NBROOTNAME.py
  mv -f $NBROOTNAME.py ../formats/py/
else
  echo 'Not Generating html, pdf and py output versions'
fi
```

```
[NbConvertApp] Converting notebook 10_pca.ipynb to html
[NbConvertApp] Writing 691045 bytes to 10_pca.html
[NbConvertApp] Converting notebook 10_pca.ipynb to pdf
[NbConvertApp] Support files will be in 10_pca_files/
[NbConvertApp] Making directory ./10_pca_files
[NbConvertApp] Making directory ./10_pca_files
[NbConvertApp] Making directory ./10_pca_files
[NbConvertApp] Making directory ./10_pca_files
[NbConvertApp] Writing 86353 bytes to notebook.tex
[NbConvertApp] Building PDF
[NbConvertApp] Running xelatex 3 times: ['xelatex', 'notebook.tex', '-quiet']
[NbConvertApp] Running bibtex 1 time: ['bibtex', 'notebook']
[NbConvertApp] WARNING | bibtex had problems, most likely because there were no
citations
[NbConvertApp] PDF successfully created
[NbConvertApp] Writing 140519 bytes to 10_pca.pdf
[NbConvertApp] Converting notebook 10_pca.ipynb to script
```

[NbConvertApp] Writing 32259 bytes to 10\_pca.py