

FUNDAMENTALS OF MACHINE LEARNING IN DATA SCIENCE

CSIS 3290

DIMENSION REDUCTION

PRINCIPAL COMPONENT ANALYSIS (PCA)

FATEMEH AHMADI

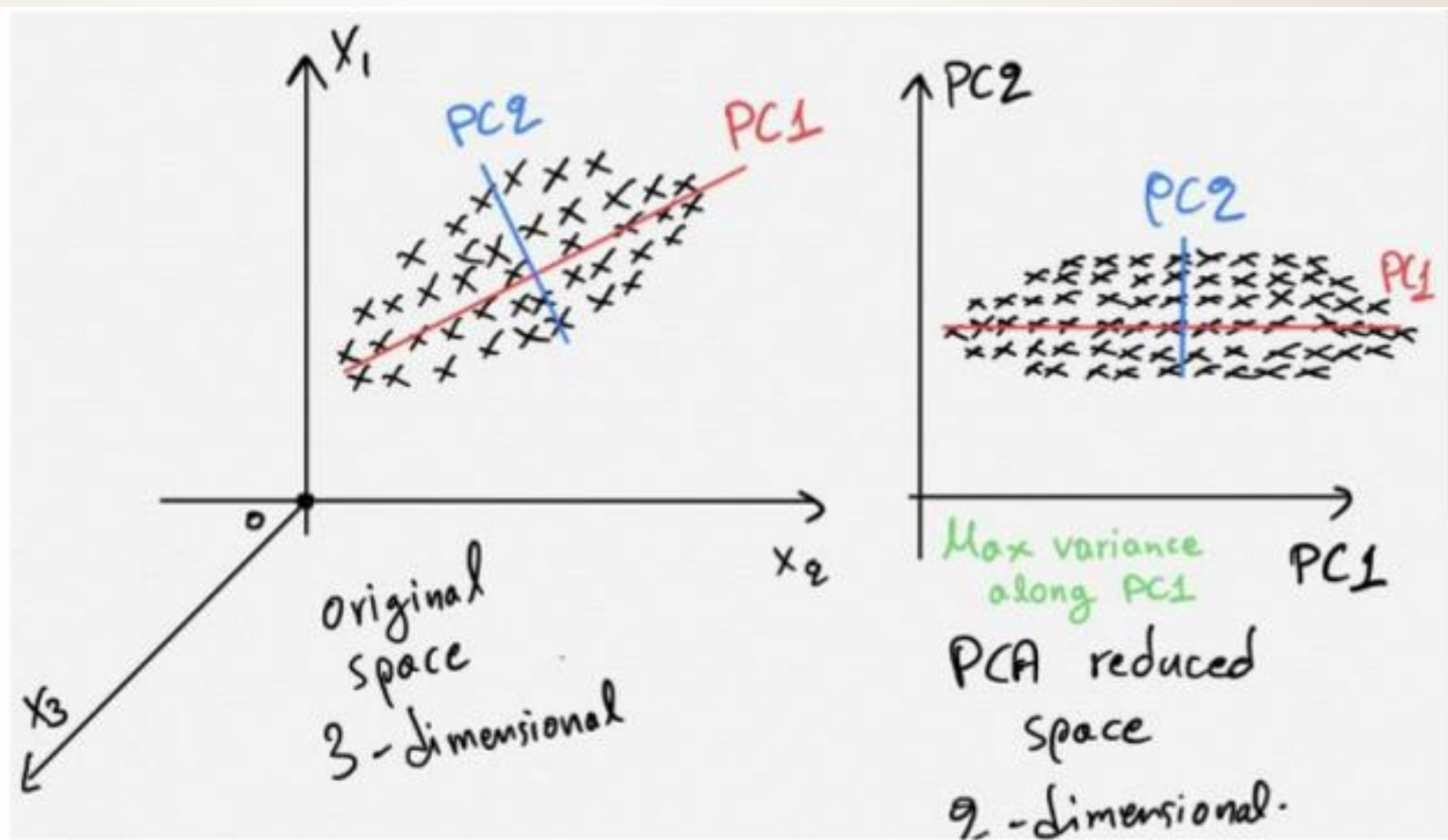
Principal Component Analysis (PCA)

- Principal Component Analysis (PCA) is a **linear dimensionality reduction** technique that can be utilized for extracting information from a high-dimensional space by projecting it into a lower-dimensional sub-space. It tries to preserve the essential parts that have more variation of the data and remove the non-essential parts with fewer variation.
- Dimensions are nothing but features that represent the data. PCA is a **statistical** procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables (entities each of which takes on various numerical values) **into a set of values of linearly uncorrelated variables called principal components**.

PCA

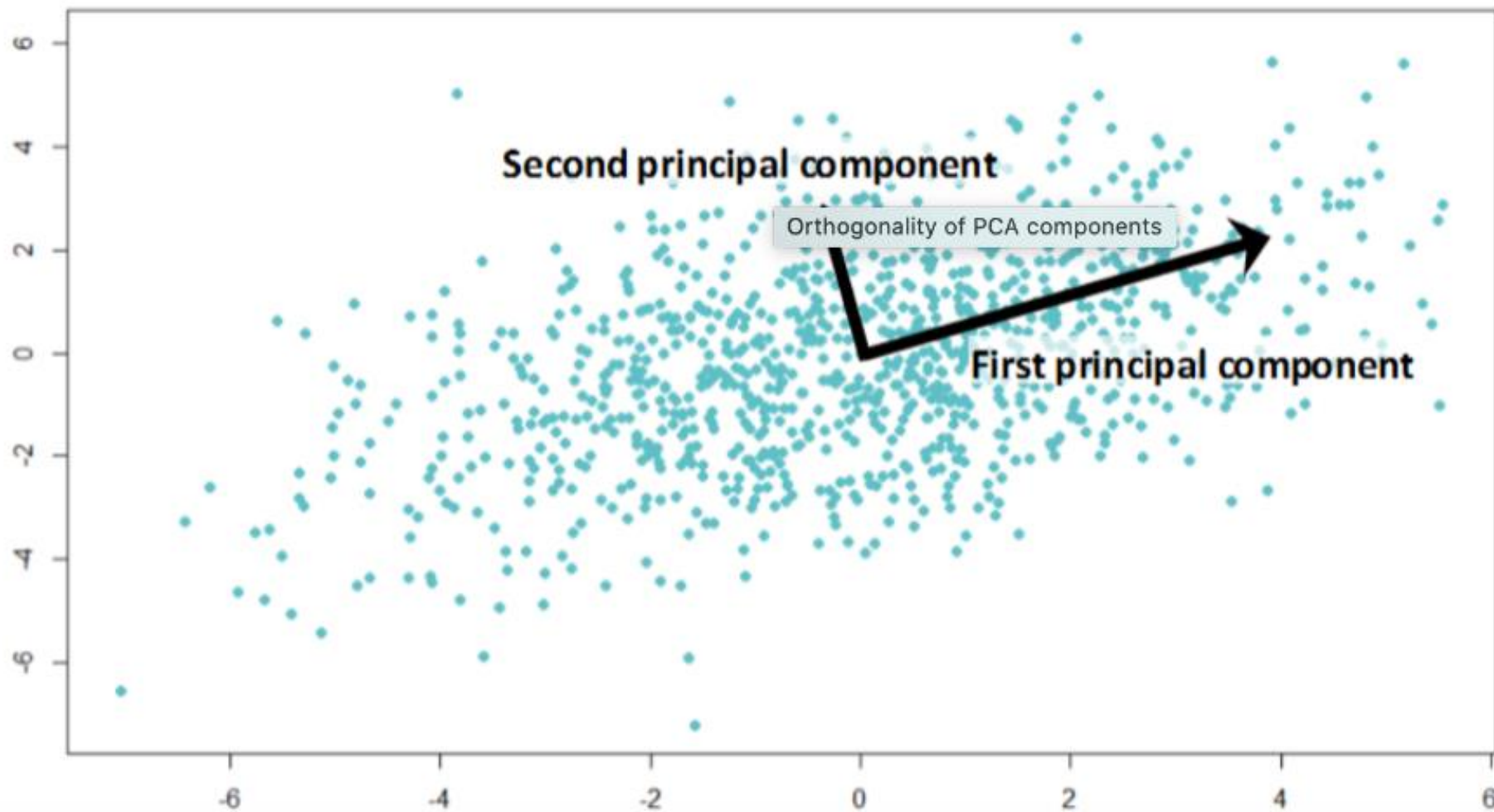
When/Why to use PCA

- PCA technique is particularly useful in processing data where **multicollinearity** exists between the **features/variables**.
- PCA can be used when **the dimensions of the input features are high** (e.g. a lot of variables).
- PCA can be also used for **denoising** and **data compression**.



PCA

https://www.bing.com/images/search?view=detailV2&ccid=77HVF15L&id=D2DF62D01D8BC2EF468A4F84E7425B3E6103AEB6&thid=OIP.77HVF15LCgFrE6qF1WmC0gHaD3&mediaurl=https%3a%2f%2fmiro.medium.com%2fmax%2f3264%2f1*2S1CN1w6FtbAlkHOb1fcsg.png&cdnurl=https%3a%2f%2fh.bing.com%2fth%2fid%2fR.efb1d5148e4b0a016b13aa85d56982d2%3fr%3dq4DYT5bQueETw%26pid%3dImgRaw%26r%3d0&exph=852&expw=1632&q=principal+component+analysis&simid=608013231911107578&FORM=IRPRST&ck=008DC95F18D5C7D31F3E9A19F8EE5B72&selectedIndex=43&ajaxhist=0&ajaxserp=0



PCA in Python

```
from sklearn.datasets import load_iris
from sklearn.decomposition import PCA
```

```
In [33]: iris1=load_iris()
```

```
In [34]: iris2=iris1.data[:,[0,2]]
```

```
In [36]: pca1=PCA()
```

```
In [37]: pca1.fit(iris2)
```

```
Out[37]:
```

▼ PCA

PCA()

PCA in Python

First and second components (by default two components will be built)

```
In [84]: pc=pca1.components_[0]  
         pcc=pca1.components_[1]
```

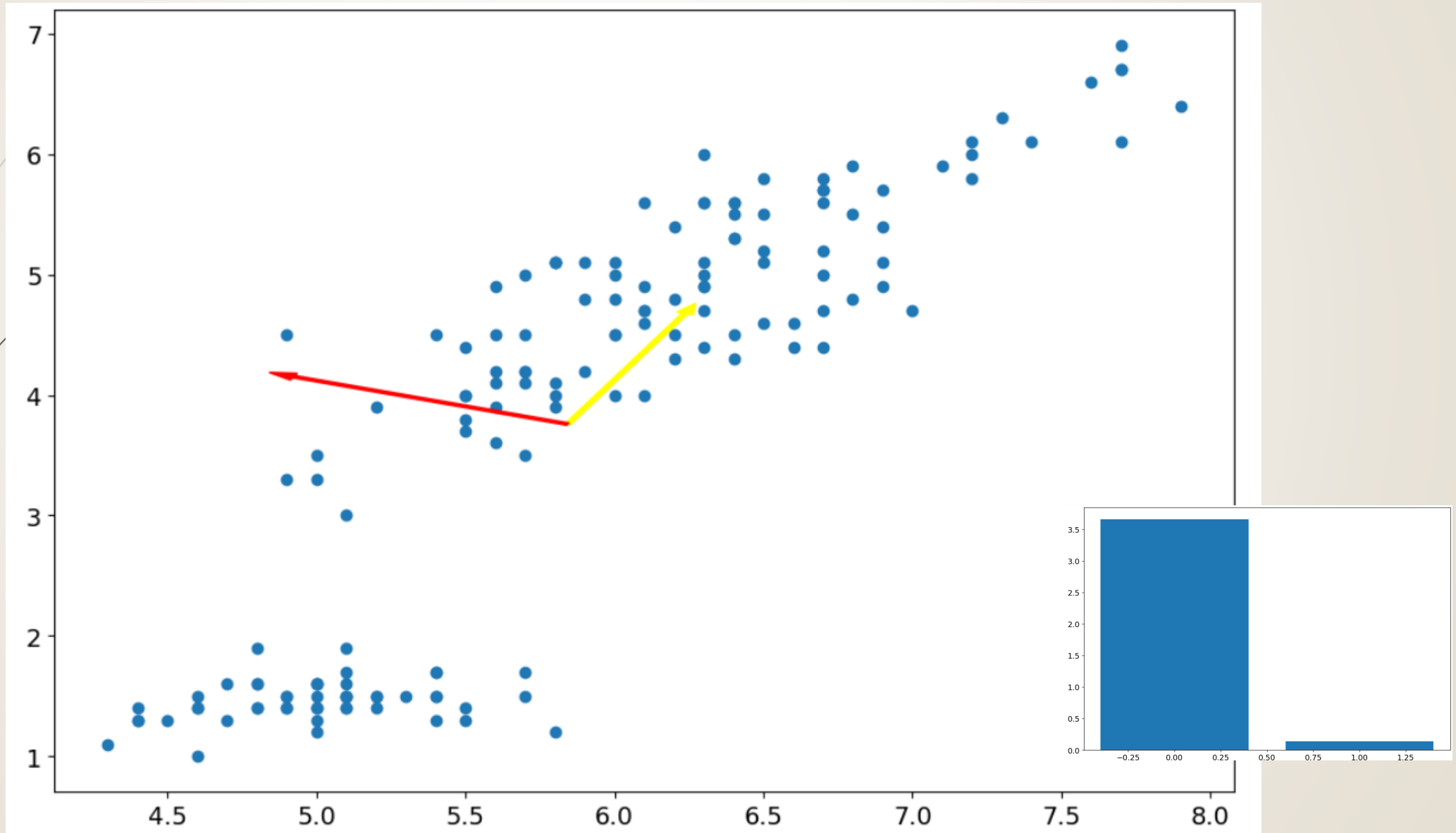
```
In [85]: mean1=pca1.mean_
```

```
In [87]: plt.scatter(iris2[:,0],iris2[:,1])  
         plt.arrow(mean1[0],mean1[1],pc[0],pc[1],color='yellow',width=0.02)  
         plt.arrow(mean1[0],mean1[1],pcc[0],pcc[1],color='red',width=0.02)  
         plt.show()
```

x and y of the start point of the arrow

x and y of the end point of the arrow

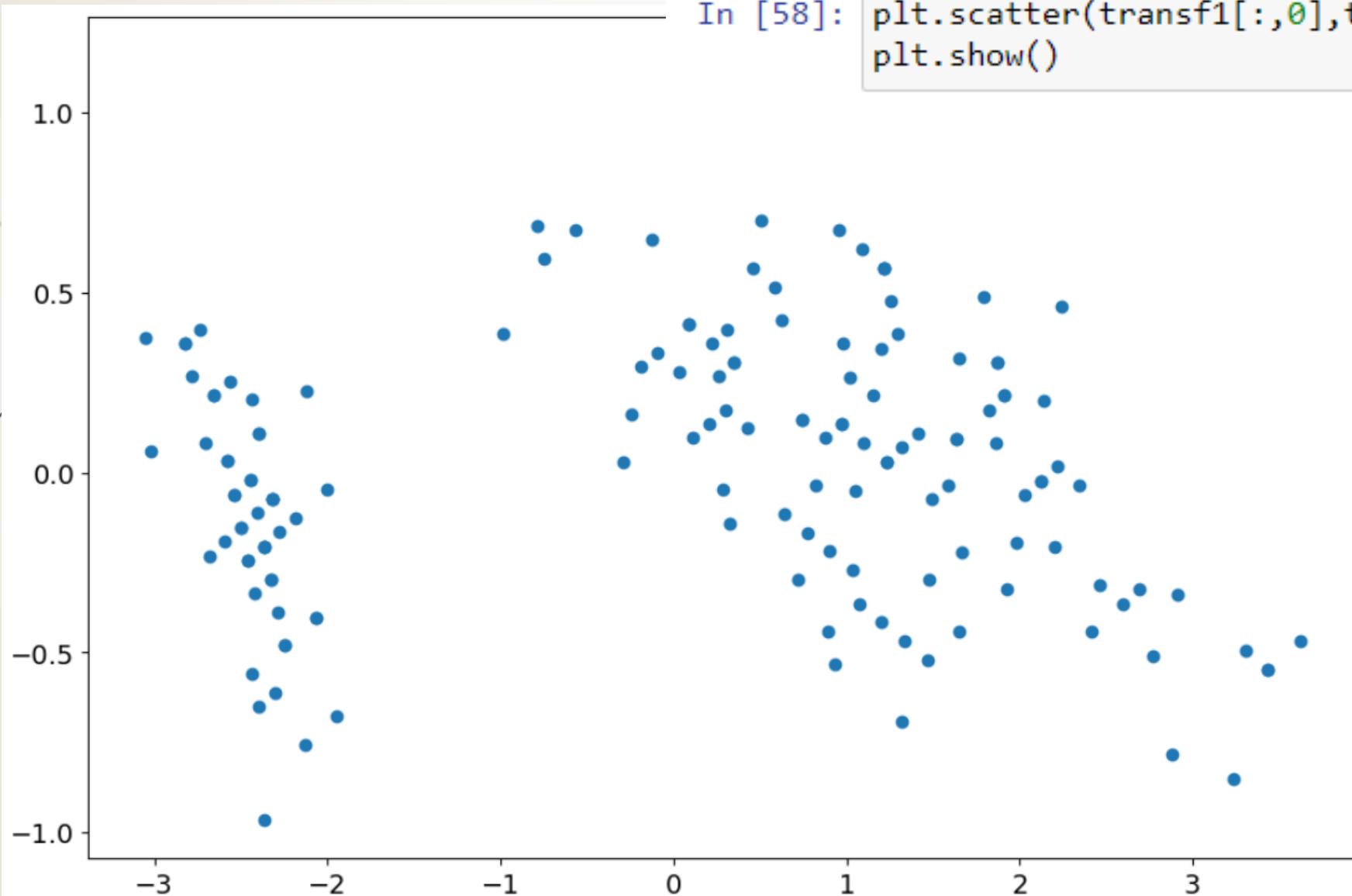
PCA in Python



PCA in Python

```
In [56]: transf1=pca1.transform(iris2)
```

```
In [58]: plt.scatter(transf1[:,0],transf1[:,1])  
plt.show()
```



PCA in Python

```
In [67]: pca1.n_components_
```

```
Out[67]: 2
```

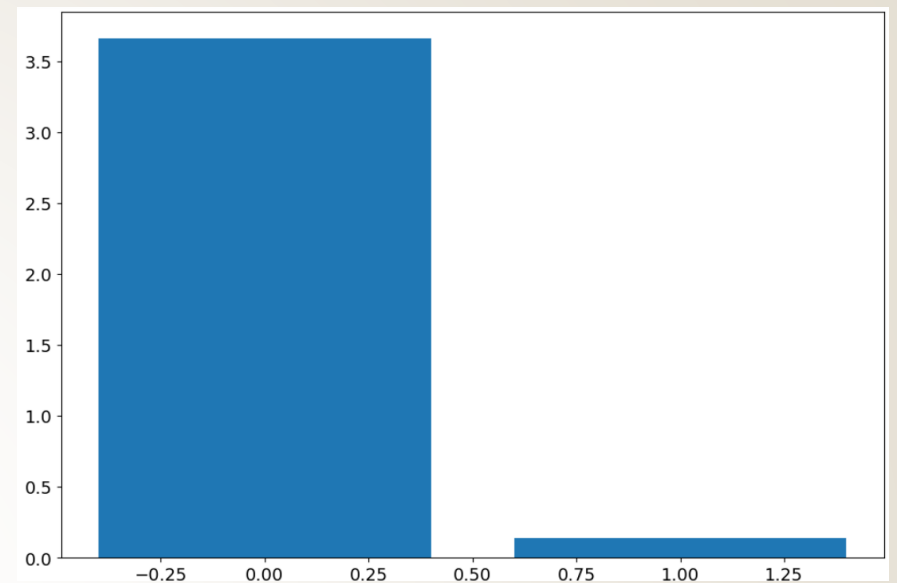
```
In [68]: nf=range(pca1.n_components_)
```

```
In [69]: var_nf=pca1.explained_variance_
```

```
In [70]: var_nf
```

```
Out[70]: array([3.66189877, 0.1400726 ])
```

```
In [71]: plt.bar(nf,var_nf)  
plt.show()
```



```
In [59]: pca1.explained_variance_ratio_
```

```
Out[59]: array([0.9631579, 0.0368421])
```

PCA in Python

```
In [72]: pca2=PCA(n_components=2)
```

```
In [74]: pca2.fit(iris1.data)
```

```
Out[74]:
```

PCA

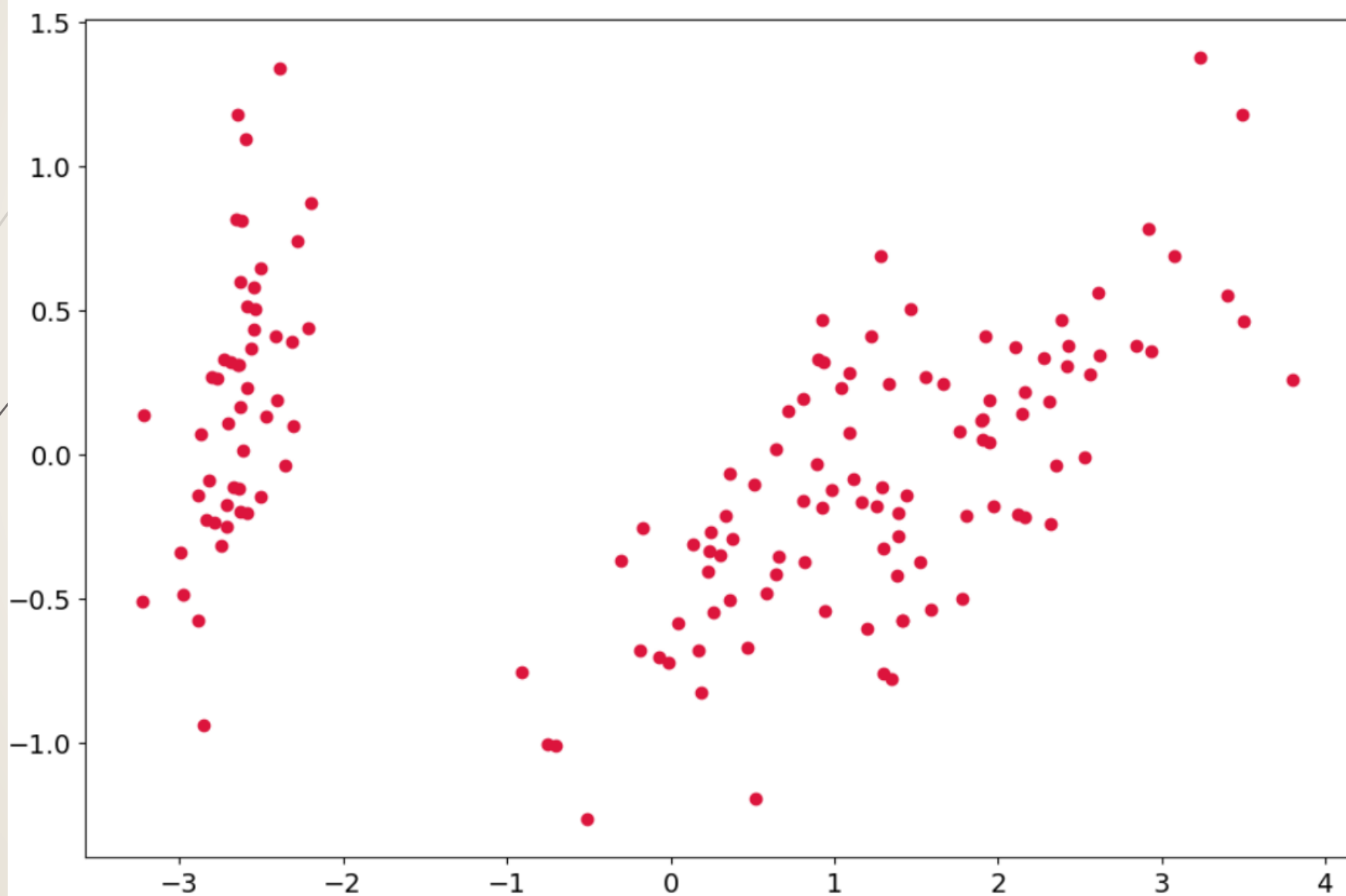
PCA(n_components=2)

```
In [76]: transf3=pca2.transform(iris1.data)  
transf3
```


PCA in Python

13

```
In [79]: plt.scatter(transf3[:,0],transf3[:,1], c='crimson')  
plt.show()
```



Reference

Data Mining, Concepts and Techniques,
Jiawei Han, Micheline Kamber, Jian Pei.
MK. Chapter 9.

