# ESS2222

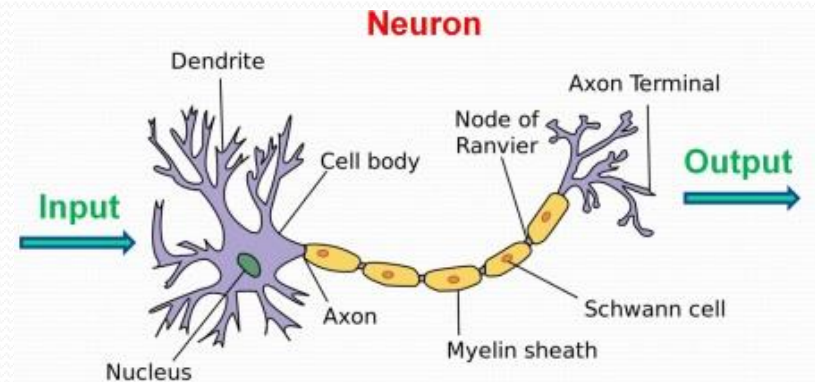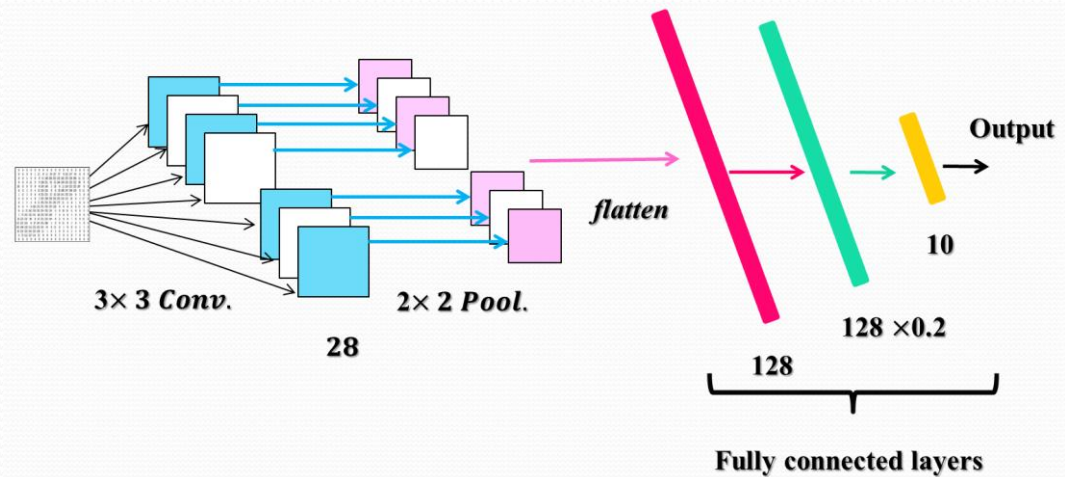# Lecture 8 – Mining of Massive Data - Dimensionality Reduction

*Hosein Shahnas*

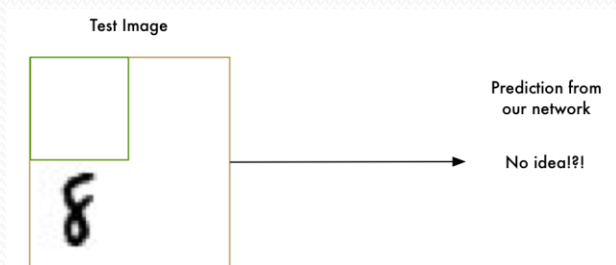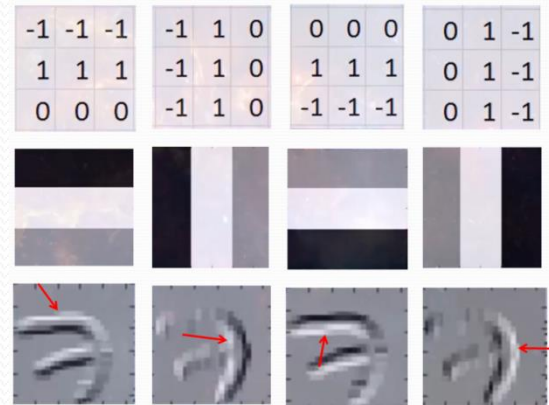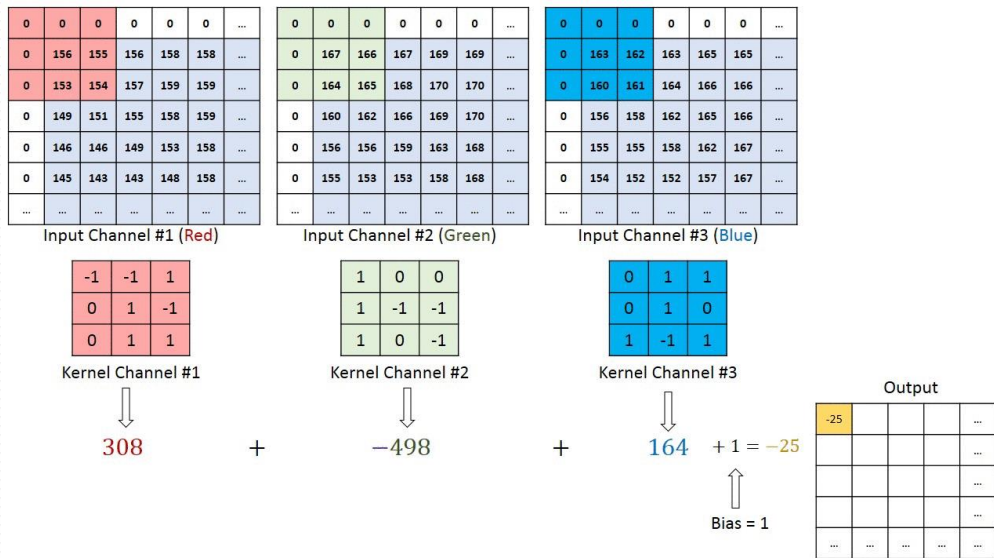*University of Toronto, Department of Earth Sciences,*

# Outline

❑ **Principal Component Analysis**
❑ **Dimensionality Reduction**

Input Channel #1 (Red)  Input Channel #2 (Green)  Input Channel #3 (Blue)

Kernel Channel #1  Kernel Channel #2  Kernel Channel #3

$$308 + -498 + 164 + 1 = -25$$

Bias = 1

Output



Test Image

Prediction from
our network

No idea!?!

$3 \times 3$ *Conv.*  $2 \times 2$ *Pool.*

*flatten*

28

Output

10

$128 \times 0.2$
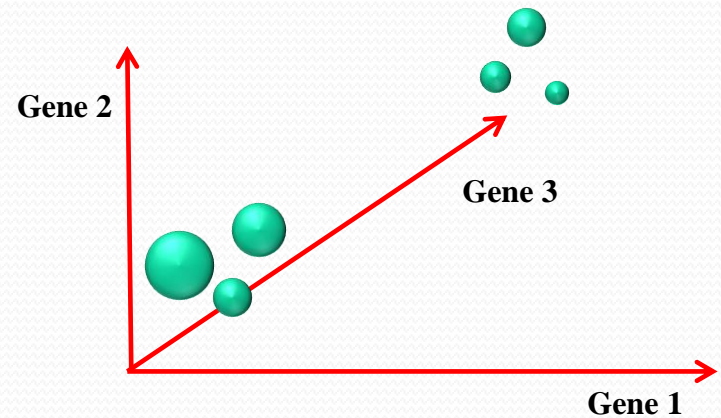
128

**Fully connected layers**

3

# Principal Component Analysis

**Principal component analysis (PCA)**
A statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly **correlated variables** into a set of values of **linearly uncorrelated variables** called principal components.

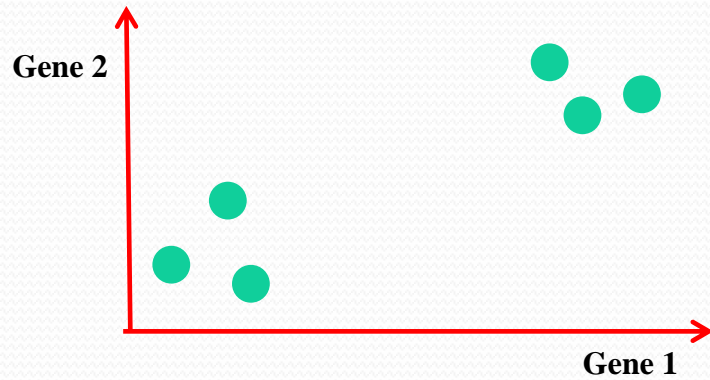| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |
| Gene 3 | 12 | 9 | 10 | 2.5 | 1.3 | 2 |

**Feature 1**
**Feature 2**
**Feature 3**

**How can we take three or more features (three or more dimensional feature data) and make a lower PC representation (lower dimension)?**

# Principal Component Analysis

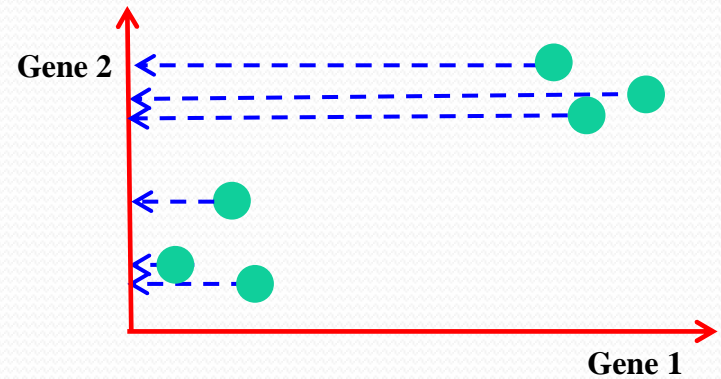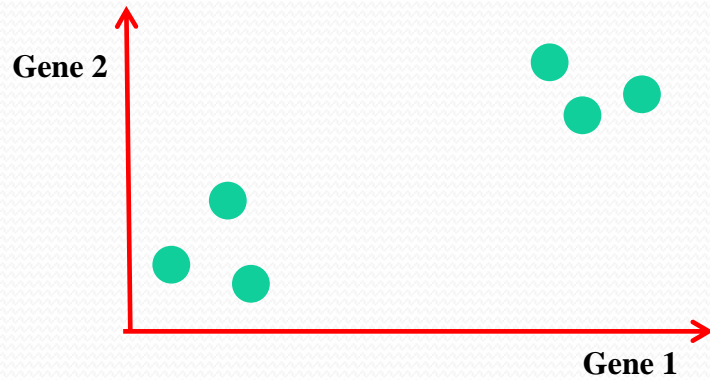**Let's assume we have only two genes (two features):**

| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

**Feature 1**
**Feature 2**



5

# Principal Component Analysis

**Let's assume we have only two genes (two features):**

| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| **Feature 1** Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| **Feature 2** Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

**Let's assume we have only two genes (two features):**

| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

**Feature 1**
**Feature 2**

# Principal Component Analysis

**Let's assume we have only two genes (two features):**

| | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| **Feature 1** — Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| **Feature 2** — Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

# Principal Component Analysis

**Let's assume we have only two genes (two features):**

|  | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

Feature 1
Feature 2

**a) Find the center of mean**

# Principal Component Analysis
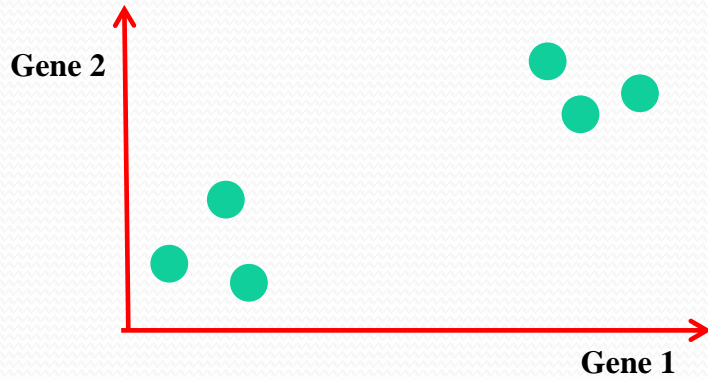
**Let's assume we have only two genes (two features):**

|          | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|----------|---------|---------|---------|---------|---------|---------|
| Gene 1   | 10      | 11      | 8       | 3       | 2       | 1       |
| Gene 2   | 6       | 4       | 5       | 3       | 2.8     | 1       |

**Feature 1**
**Feature 2**

a) Find the center of mean
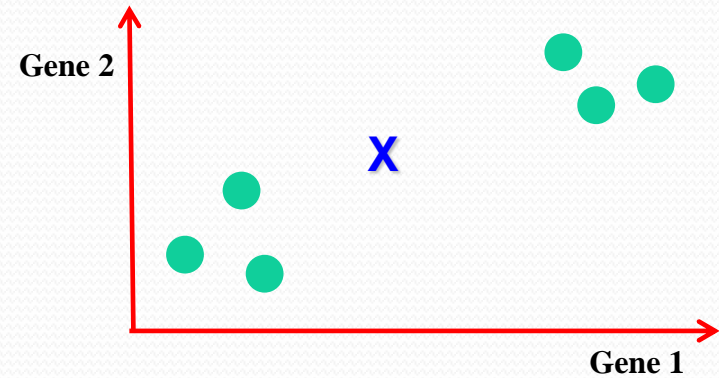b) Move the center to the origin

# Principal Component Analysis

**Let's assume we have only two genes (two features):**

|  | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

**Feature 1** is Gene 1
**Feature 2** is Gene 2

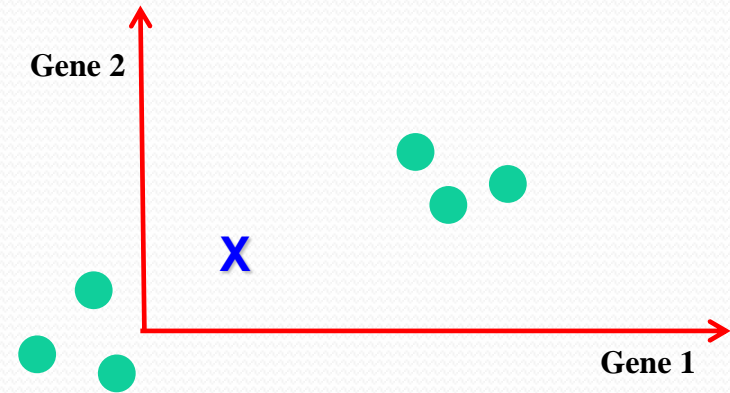a) **Find the center of mean**
b) **Move the center to the origin**

# Principal Component Analysis

**Let's assume we have only two genes (two features):**

|  | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---|---|---|---|---|---|---|
| **Feature 1** Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| **Feature 2** Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

a) Find the center of mean
b) Move the center to the origin
c) Fit the data points to a line (minimize b or maximize c)

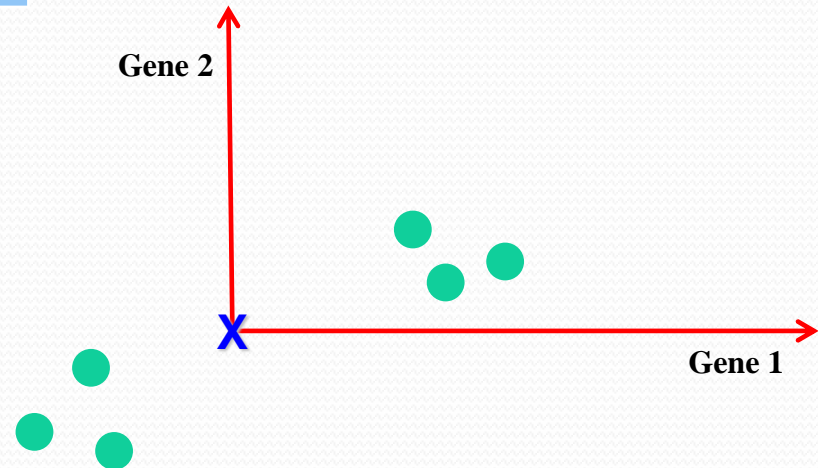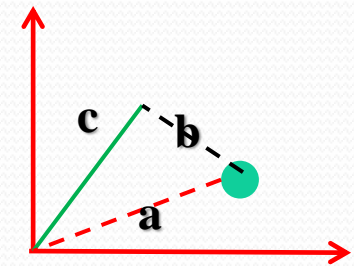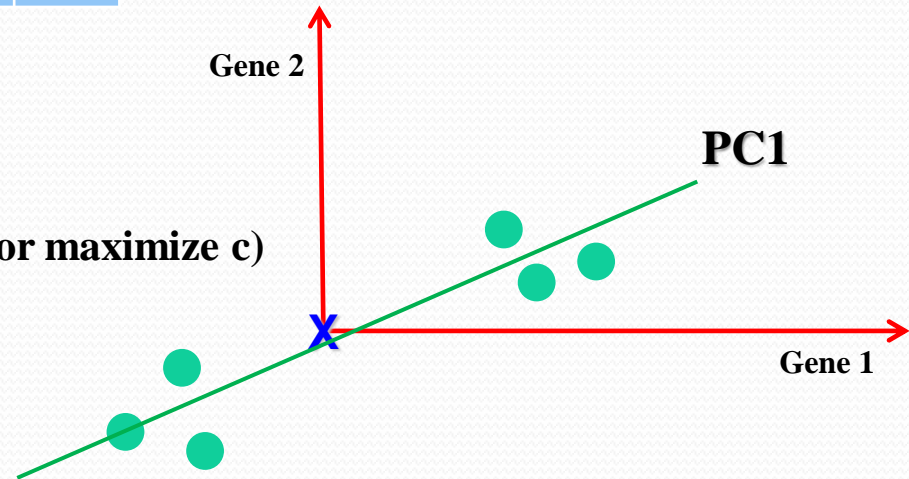**Max** $\sum_i c_i^2$



12

# Principal Component Analysis

**Let's assume we have only two genes (two features):**

|         |        | Mouse 1 | Mouse 2 | Mouse 3 | Mouse 4 | Mouse 5 | Mouse 6 |
|---------|--------|---------|---------|---------|---------|---------|---------|
| **Feature 1** | Gene 1 | 10 | 11 | 8 | 3 | 2 | 1 |
| **Feature 2** | Gene 2 | 6 | 4 | 5 | 3 | 2.8 | 1 |

a) Find the center of mean
b) Move the center to the origin
c) Fit the data points to a line (minimize b or maximize c)
**Max** $\sum_i c_i^2$
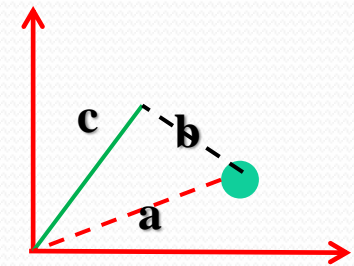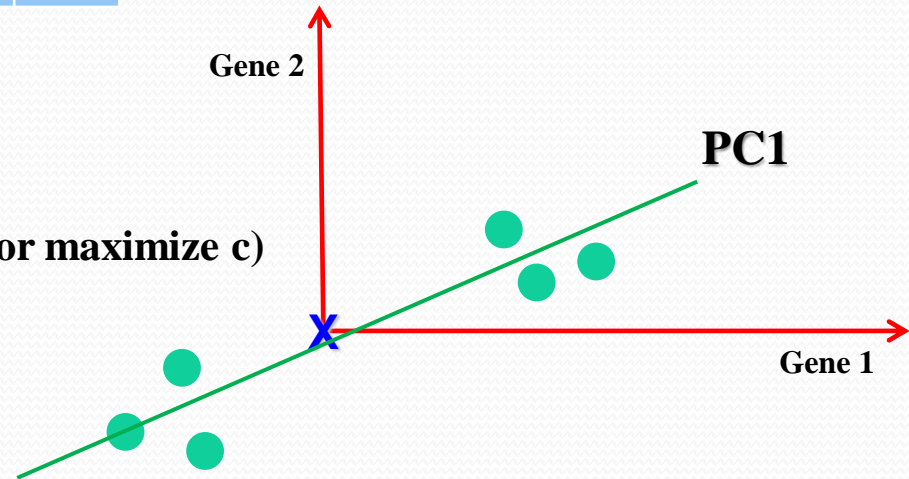
PC1: principal component 1, which is a linear combination of feature 1 and 2.

A unit vector along PC1 is an Eigenvector for PC1.
SSD = **Max** $\sum_i c_i^2$    Eigenvalue for PC1
$\sqrt{SSD}$    Singular value for PC1

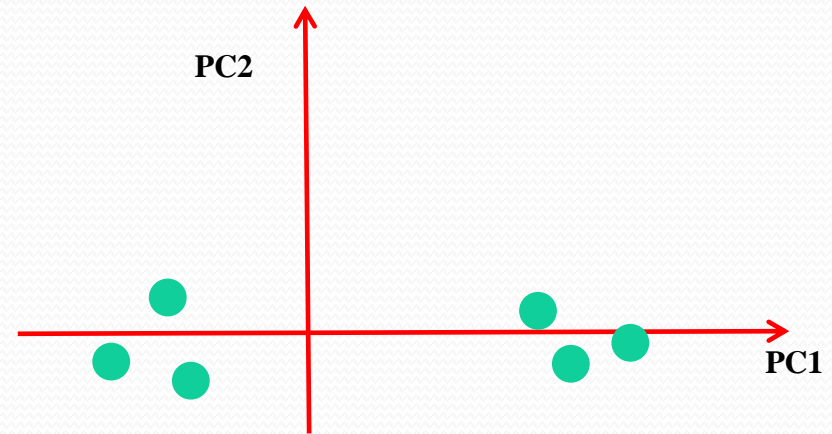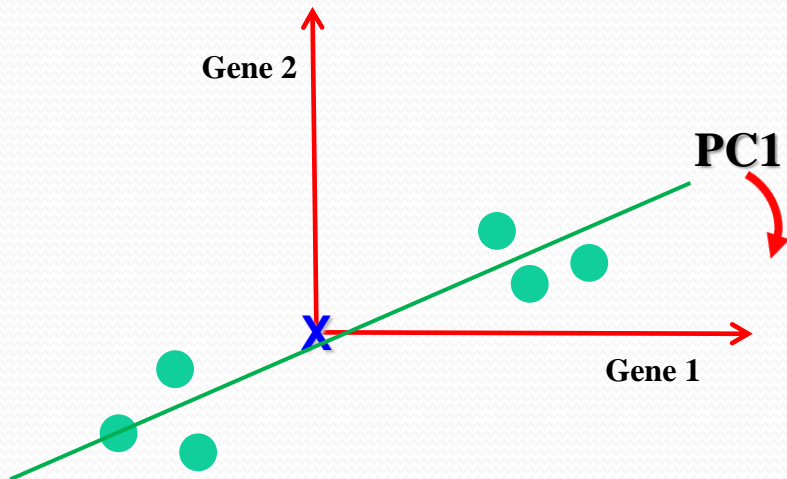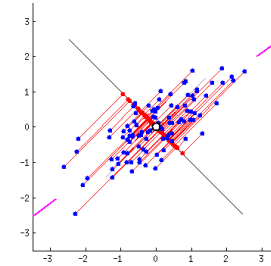SSD: Sum of squares od distances

# Principal Component Analysis

$$\frac{SSD_{PC1}}{n-1}$$  **variation for PC1**

$$\frac{SSD_{PC2}}{n-1}$$  **variation for PC2**



**For example if:**

$$\frac{SSD_{PC1}}{n-1} = 16, \qquad \frac{SSD_{PC2}}{n-1} = 2 \qquad \textbf{Var(PC1) + Var(PC2) = 18,}$$

**PC1 accounts for 16/18 = 89% of variations (information)  around PC's .**

# Dimensionality Reduction

**Singular Value Decomposition (SVD)**

$$A_{[mn]} = U_{[mr]} \ \Sigma_{[rr]} \ V_{[nr]}{}^{T}$$

m: number of rows
n: number of columns

A: Input data
U: The <span style="color:red">left-singular</span> vectors of A are a set of orthonormal eigenvectors of $AA^{T}$.
V: The <span style="color:green">right-singular</span> vectors of A are a set of orthonormal eigenvectors of $A^{T}A$.
$\Sigma$ : Singular values

The non-zero singular values of A (found on the diagonal entries of $\Sigma$) are the square roots of the non-zero eigenvalues of both $AA^{T}$ and $A^{T}A$.

# Dimensionality Reduction

It is always possible to decompose a real matrix A into $A = U \Sigma V^T$ ,where

U, $\Sigma$, V: unique (decomposition is unique)
U,V: orthonormal ($U^T U = I$; $V^T V = I$)
$\Sigma$ : diagonal
Entries are (singular values) are positive ($\sigma_1 \geq \sigma_2 \geq \ldots \sigma_r \geq 0$).

**Example:**
Netflix users ranking the movies

# Dimensionality Reduction

**What do we learn from SVD decomposition?**



The **first four users** **strongly** correspond to SciFi – concept,
The **last three users** **heavily** correspond to Romance – concept,

U: **"user-to-concept"** similarity matrix
V: **"movie-to-concept"** similarity matrix

# Principal Component Analysis

| 784 | 331 | 154 | 87 | 59 | **Components** |
|-----|-----|-----|-----|-----|-----|



| original | 99% | 95% | 90% | 85% | **Explained Variance** |
|----------|-----|-----|-----|-----|-----|

**Example 1:** Reduction of the number of features in Iris-problem from 4 to 2:

## Standardization



## PCA



## Concatenating



principalDf      df[['target']]      finalDf



2 Component PCA Graph

# PCA for Visualization

```
1 #https://towardsdatascience.com/pca-using-python-scikit-learn-e653f8989e60
2 #from sklearn.datasets import fetch_mldata
3 #mnist = fetch_mldata('MNIST original')
4 import sys
5 #========================================= PCA for Data Visualization
6 #========================================= import data
7 print('========================================= 1')
8 import pandas as pd
9 url = "https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data"
10 # Load dataset into Pandas DataFrame
11 df = pd.read_csv(url, names=['sepal length','sepal width','petal length','petal width','target'])
12 #df = pd.read_csv(url, names=['sepal length','sepal width','target'])
13
14 print ('df = ', df)
15 print('========================================= 1')
16 #========================================= import data
17 #========================================= set x and y and Standardize
18 print('========================================= 2')
19 from sklearn.preprocessing import StandardScaler
20 features = ['sepal length', 'sepal width', 'petal length', 'petal width']
21 # Separating out the features
22 x = df.loc[:, features].values
23 # Separating out the target
24 y = df.loc[:,['target']].values
25 # Standardizing the features
26 x = StandardScaler().fit_transform(x)
27 #print ('y = ', y)
28 #print ('x = ', x)
29 print ('y.shape = ', y.shape)
30 print ('x.shape = ', x.shape)
31 print('========================================= 2')
32 #========================================= set X and y and Standardize
33
34 #================================================================ Principal Component Analysis 1
35
36 print('========================================= 3')
37 from sklearn.decomposition import PCA
38 pca = PCA(n_components=2)
39 principalComponents = pca.fit_transform(x)
40 print('pca.explained_variance_ratio_ = ', pca.explained_variance_ratio_ )
```

```python
41 principalDf = pd.DataFrame(data = principalComponents
42                , columns = ['principal component 1', 'principal component 2'])
43 #print('principalComponents = ',principalComponents)
44 print ('principalComponents.shape = ', principalComponents.shape)
45 print ('principalDf.shape = ', principalDf.shape)
46
47 #print ('principalComponents = ', principalComponents)
48 print ('principalDf.head(3) = ', principalDf.head(3))
49
50 print('====================================== 3')
51
52 #======================================================= Principal Component Analysis 1
53
54 #======================================================= Principal Component Analysis 2
55 '''
56 print('====================================== 4')
57 # scikit-learn choose the minimum number of principal components such that 95% of the variance is retained.
58 from sklearn.decomposition import PCA
59 pca = PCA(0.98)
60 principalComponents = pca.fit_transform(x)
61 n_dim = principalComponents[1].shape
62 print('pca.explained_variance_ratio_ = ', pca.explained_variance_ratio_ )
63 Feature_size = int(principalComponents.size/y.size)
64 print('Feature_size = ', Feature_size)
65
66 if (Feature_size==1):
67     columns0 = ['principal component 1']
68 if (Feature_size==2):
69     columns0 = ['principal component 1', 'principal component 2']
70 if (Feature_size==3):
71     columns0 = ['principal component 1', 'principal component 2', 'principal component 3']
72 if (Feature_size==4):
73     columns0 = ['principal component 1', 'principal component 2', 'principal component 3', 'principal component 4']
74
75 principalDf = pd.DataFrame(data = principalComponents
76                , columns = columns0)
77 #print('principalComponents = ',principalComponents)
78 print ('principalComponents.shape = ', principalComponents.shape)
79 print ('principalDf.shape = ', principalDf.shape)
80
```

```python
80
81 #print ('principalComponents = ', principalComponents)
82 print ('principalDf.head(3) = ', principalDf.head(3))
83 print('principalComponents.shape = ', len(principalDf.columns))
84 print('======================================= 4')
85 '''
86 #======================================= Principal Component Analysis 2
87
88 #======================================= Concatenating along axis = 1 (column)
89 print('======================================= 5')
90  # combine reduced features and targets (3 columns: 2 reduced features + 1 taget)
91 finalDf = pd.concat([[principalDf, df[['target']]], axis = 1)
92 #print('finalDf = ', finalDf)
93 print('finalDf.shape = ', finalDf.shape)
94 print('======================================= 5')
95 #======================================= Concatenating along axis = 1 (column)
96 #======================================= Visualize 2D Projection
97 import matplotlib.pyplot as plt
98 fig = plt.figure(figsize = (8,8))
99 ax = fig.add_subplot(1,1,1)
100 ax.set_xlabel('Principal Component 1', fontsize = 15)
101 ax.set_ylabel('Principal Component 2', fontsize = 15)
102 ax.set_title('2 component PCA', fontsize = 20)
103 targets = ['Iris-setosa', 'Iris-versicolor', 'Iris-virginica']
104 colors = ['r', 'g', 'b']
105 for target, color in zip(targets,colors):
106     indicesToKeep = finalDf['target'] == target
107     ax.scatter(finalDf.loc[indicesToKeep, 'principal component 1']
108                , finalDf.loc[indicesToKeep, 'principal component 2']
109                , c = color
110                , s = 50)
111 ax.legend(targets)
112 ax.grid()
113 #======================================= Visualize 2D Projection
114 explained_var = pca.explained_variance_ratio_
115 print('explained_var = ', explained_var)
116
117
```