

# Machine Learning

## Unsupervised Learning

**FAST** 

---

DISCOVERING  
THE FUTURE

# Topics of previous lectures

- ✓ Ingredients of Machine Learning
- ✓ Classification Basics, Basic Linear Classifier
- ✓ K-Nearest Neighbours and Naive Bayes Classifier
- ✓ Linear and Quadratic Discriminant Analysis
- ✓ Support Vector Machines (SVM)
- ✓ Decision Trees
- ✓ Ensemble Methods (Bagging, Weighted Voting, Stacking)
- ✓ Regression Methods
- ✓ Evaluation and Scoring of Classifiers
- ✓ Ensemble Methods (Boosting)
- ✓ Clustering (Hierarchical, K-Means, K-Medoids, DBSCAN)



# Topics of today's lecture

- Dimensionality Reduction
- Principal Component Analysis (PCA)
- t-Distributed Stochastic Neighbor Embedding (t-SNE)

# Motivation for Dimensionality Reduction

- Working directly with high-dimensional data, such as images, comes with some difficulties

# Motivation for Dimensionality Reduction

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.

# Motivation for Dimensionality Reduction

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.
- Dimensions in high-dimensional data are often correlated so that the data possesses an intrinsic lower-dimensional structure.

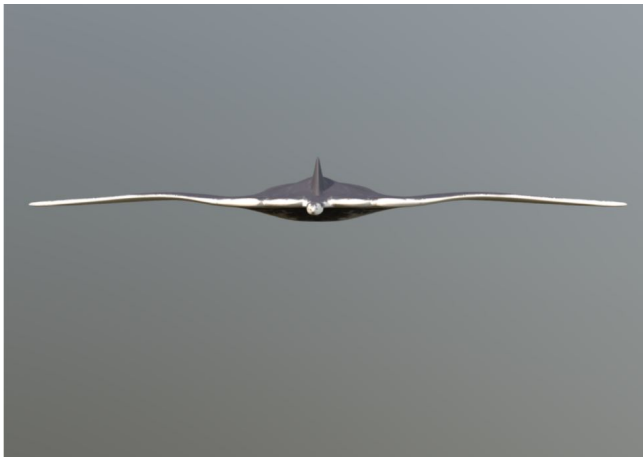
# Motivation for Dimensionality Reduction

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.
- Dimensions in high-dimensional data are often correlated so that the data possesses an intrinsic lower-dimensional structure.
- Dimensionality reduction exploits structure and correlation and allows us to work with a more compact representation of the data, ideally without losing much information.

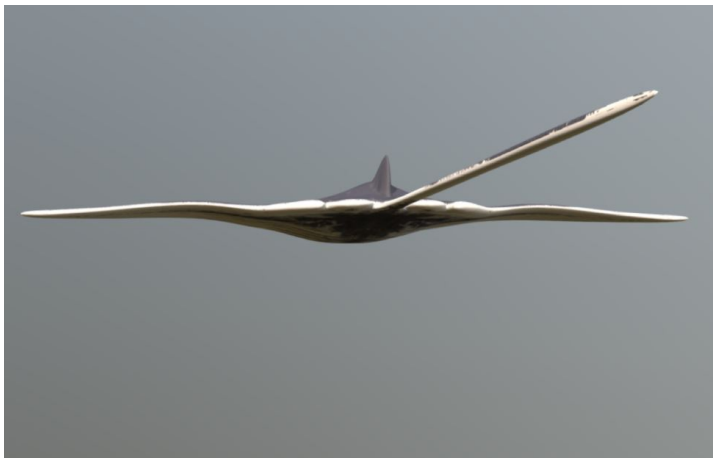
# Motivation for Dimensionality Reduction

- Working directly with high-dimensional data, such as images, comes with some difficulties
- High-dimensional data is often overcomplete, i.e., many dimensions are redundant and can be explained by a combination of other dimensions.
- Dimensions in high-dimensional data are often correlated so that the data possesses an intrinsic lower-dimensional structure.
- Dimensionality reduction exploits structure and correlation and allows us to work with a more compact representation of the data, ideally without losing much information.
- It can also be useful to detect potential patterns in high dimensional data, by visualizing the first 2-3 projections.

# Motivation for Dimensionality Reduction

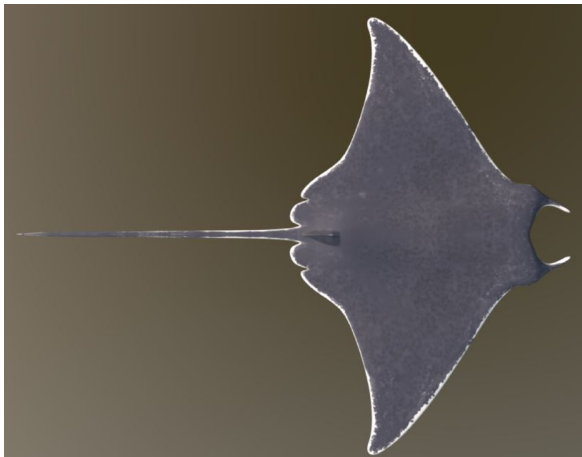


# Motivation for Dimensionality Reduction



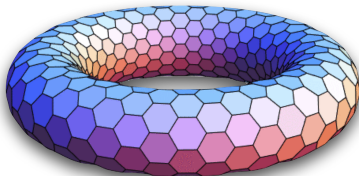


# Motivation for Dimensionality Reduction



# Principal Component Analysis

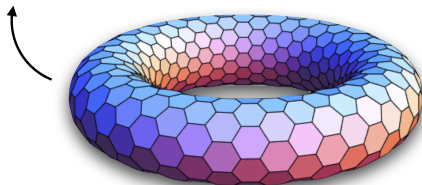
What is the problem with **high-dimensional** things?



# Principal Component Analysis

What is the problem with **high-dimensional** things?

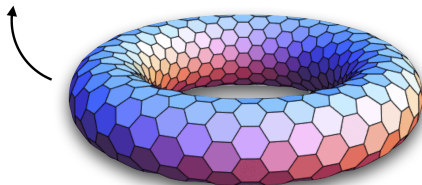
**Hard** to visualise



# Principal Component Analysis

What is the problem with **high-dimensional** things?

**Hard** to visualise

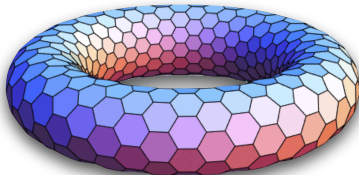


# Principal Component Analysis

What is the problem with **high-dimensional** things?

**Hard** to visualise

**Algorithms** tend to **get slow**

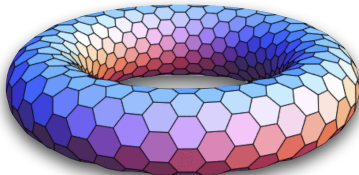


# Principal Component Analysis

What is the problem with **high-dimensional** things?

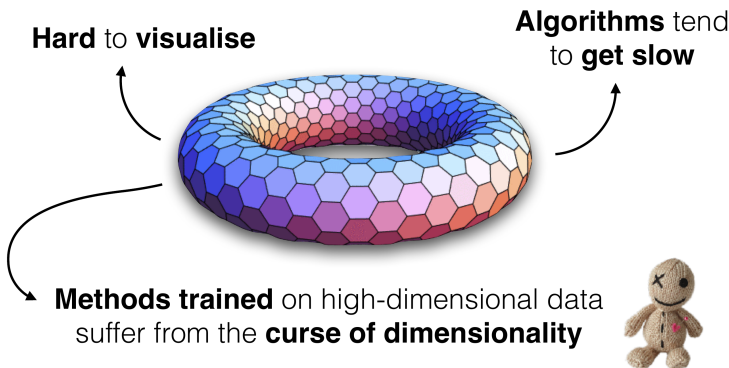
**Hard** to **visualise**

**Algorithms** tend  
to **get slow**



# Principal Component Analysis

What is the problem with **high-dimensional** things?



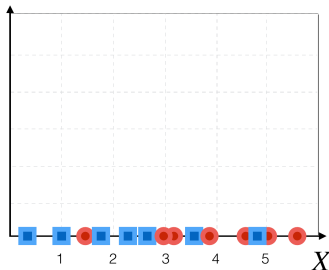
What is the **curse** of **dimensionality**?





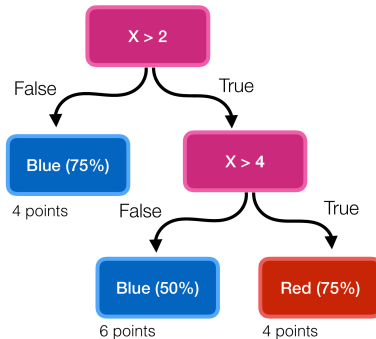
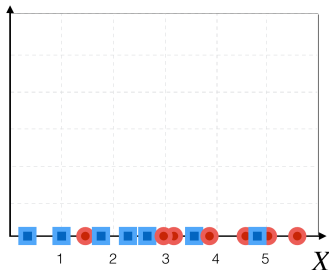
# Principal Component Analysis

What is the **curse** of **dimensionality**?



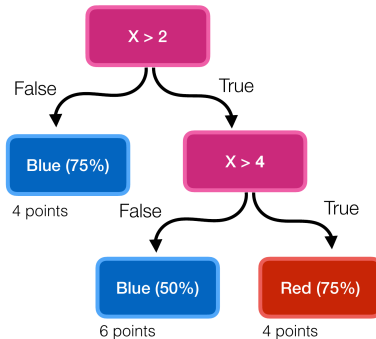
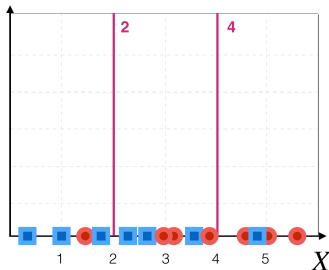
# Principal Component Analysis

What is the **curse** of **dimensionality**?



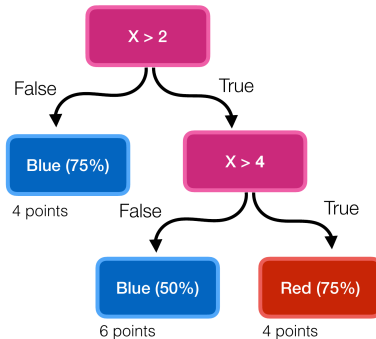
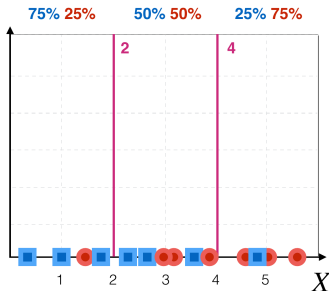
# Principal Component Analysis

What is the **curse** of **dimensionality**?



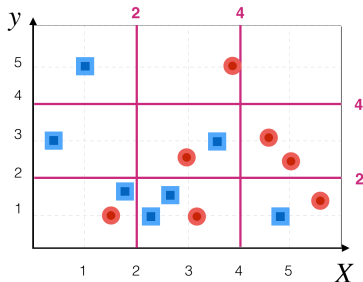
# Principal Component Analysis

What is the **curse** of **dimensionality**?



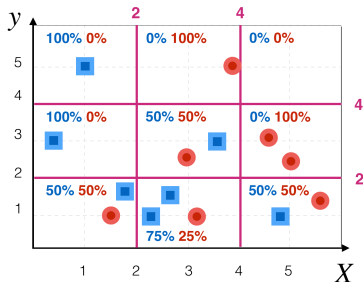
# Principal Component Analysis

What is the **curse** of **dimensionality**?



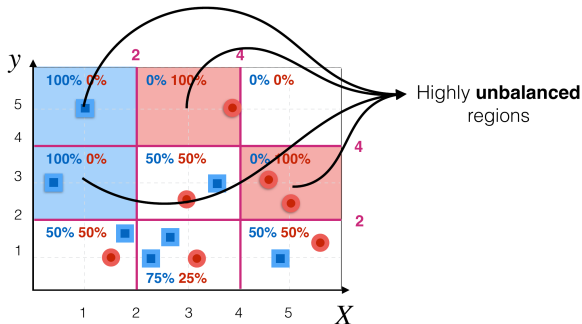
# Principal Component Analysis

What is the **curse** of **dimensionality**?



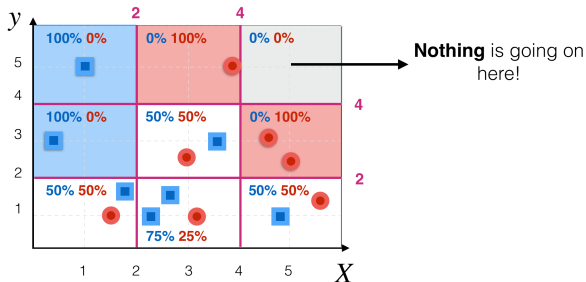
# Principal Component Analysis

What is the **curse** of **dimensionality**?



# Principal Component Analysis

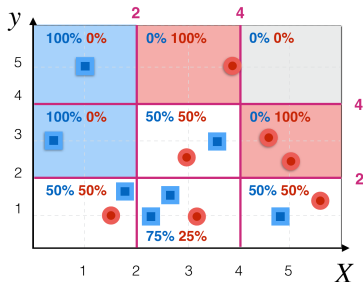
What is the **curse** of **dimensionality**?





# Principal Component Analysis

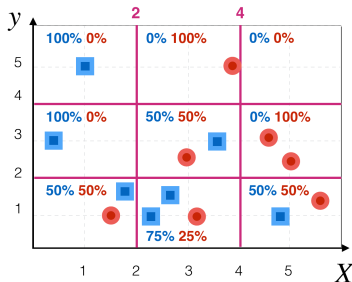
What is the **curse** of **dimensionality**?



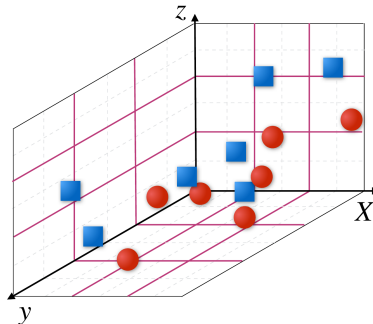
On average **55.5%** of cells will be either **empty** or **singletons**

# Principal Component Analysis

What is the **curse** of **dimensionality**?



On average **55.5%** of cells will be either **empty** or **singletons**

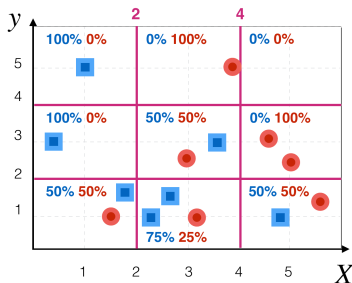


On average **92.5%** of cells will be either **empty** or **singletons**

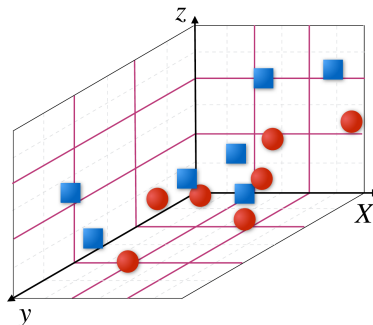
# Principal Component Analysis

## What is the **curse** of **dimensionality**?

In order to keep **high-dimensional** space reasonably covered you need **a lot more data**



On average **55.5%** of cells will be either **empty** or **singletons**



On average **92.5%** of cells will be either **empty** or **singletons**

What is the **curse** of **dimensionality**?  
(part II)

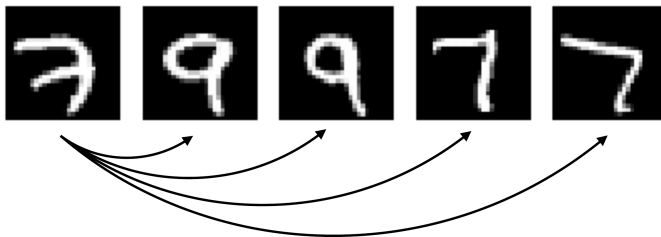


# Principal Component Analysis

What is the **curse** of **dimensionality**?  
(part II)

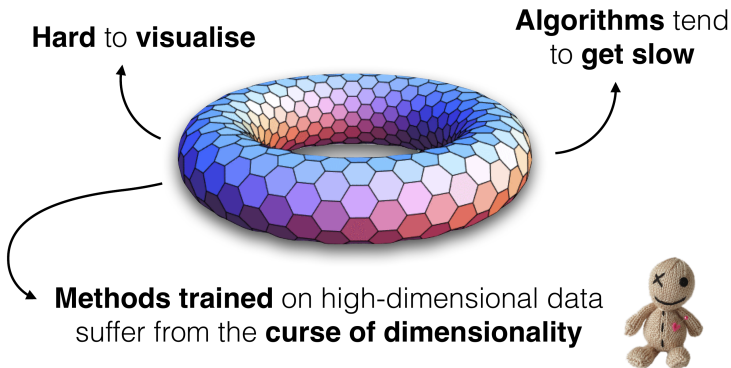


**Distances** become **similar** in high-dimensional space



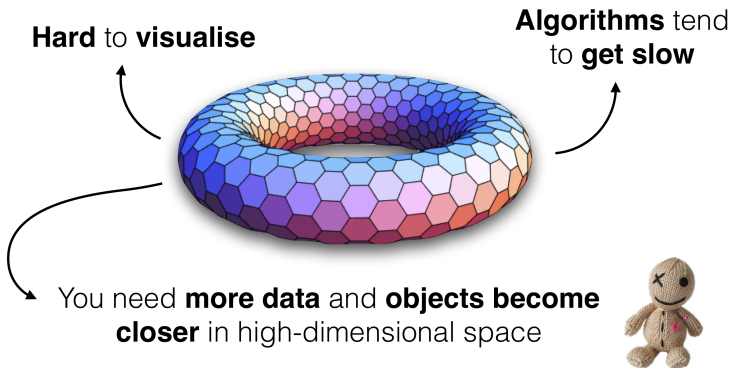
# Principal Component Analysis

What is the problem with **high-dimensional** things?

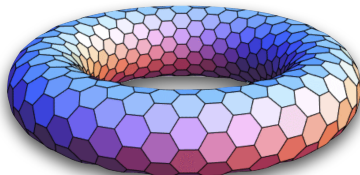


# Principal Component Analysis

What is the problem with **high-dimensional** things?



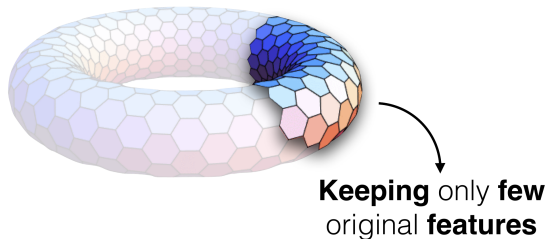
Feature **extraction** vs feature **elimination**





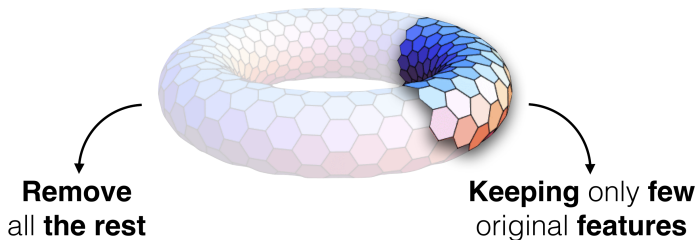
# Principal Component Analysis

Feature **extraction** vs feature **elimination**



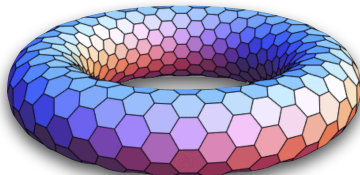
# Principal Component Analysis

Feature **extraction** vs feature **elimination**



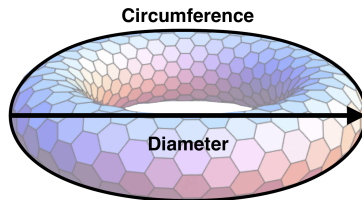
# Principal Component Analysis

Feature **extraction** vs feature **elimination**



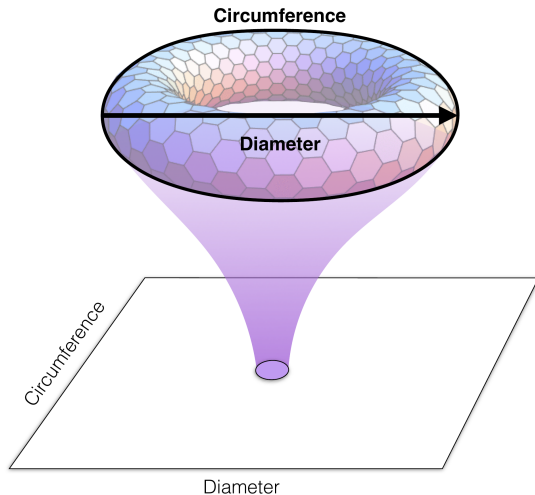
# Principal Component Analysis

Feature **extraction** vs feature **elimination**



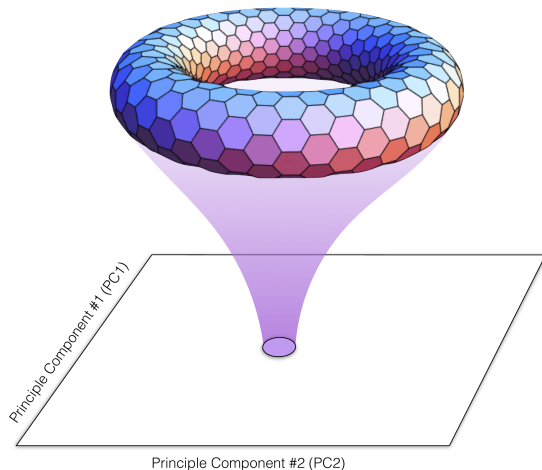
# Principal Component Analysis

Feature **extraction** vs feature **elimination**

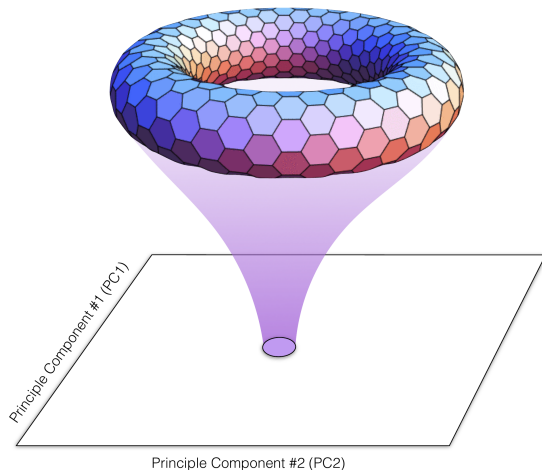


# Principal Component Analysis

## Principle Component Analysis

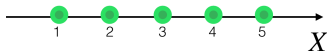


## Principle **C**omponent **A**nalysis



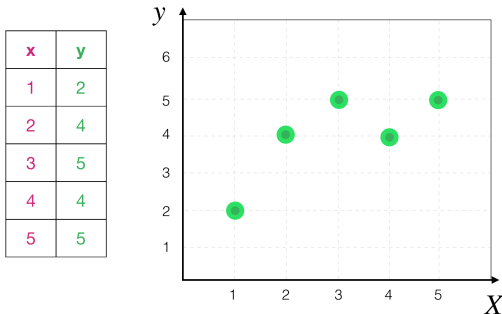
## 1-Dimensional data

$x$
1
2
3
4
5



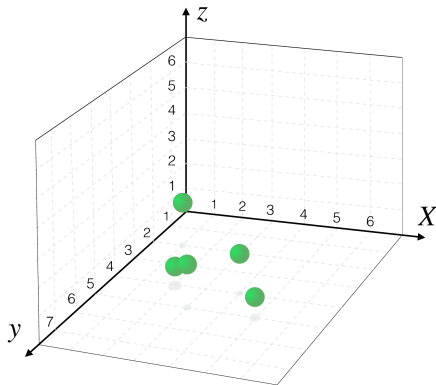


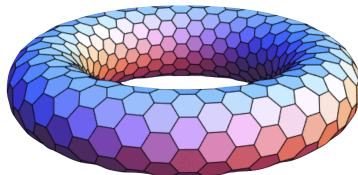
## 2-Dimensional data



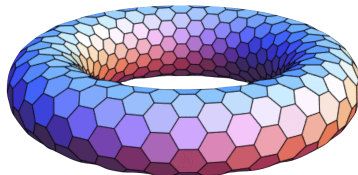
## 3-Dimensional data

x	y	z
1	2	2
2	4	0.5
3	5	1
4	4	1
5	5	0.5





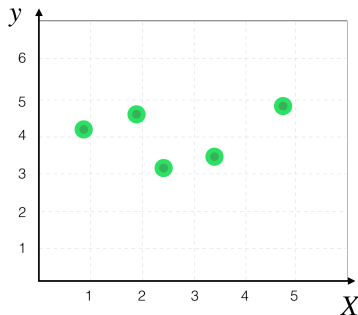
**200-D**imensional data?



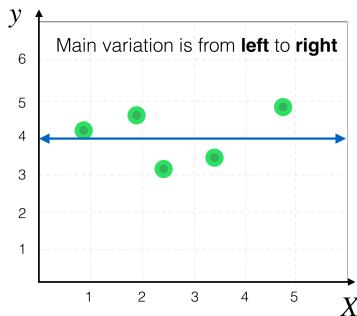
## 200-Dimensional data?

Are all of these dimensions **equally useful**?

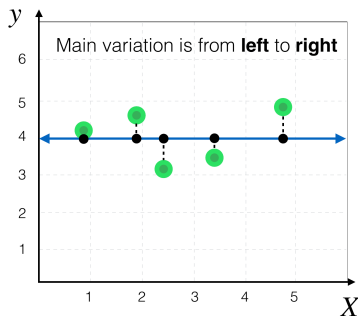
## 2-D example revisited



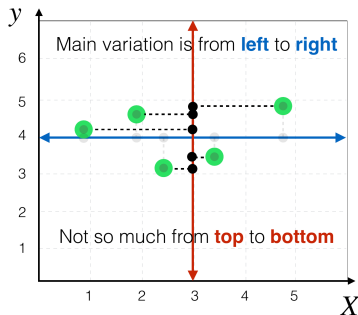
## 2-D example revisited



## 2-D example revisited

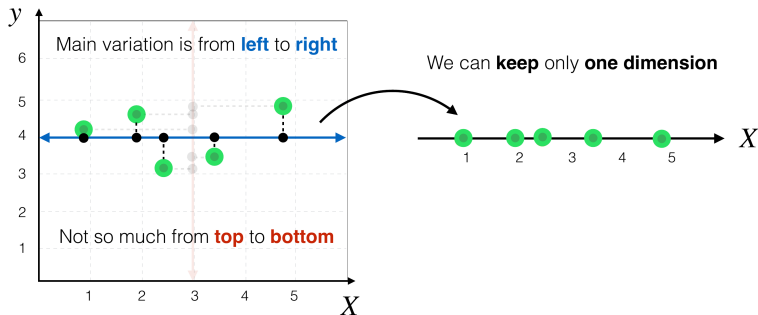


## 2-D example revisited

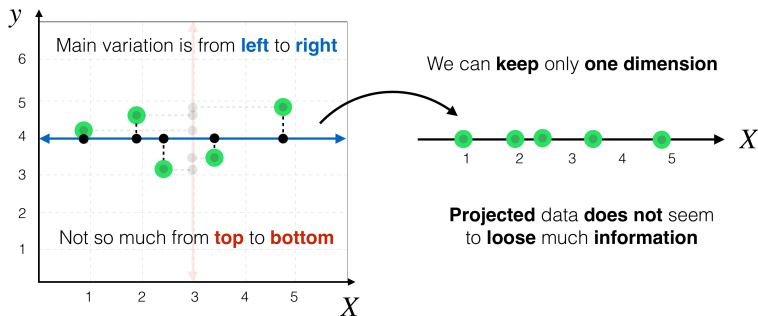




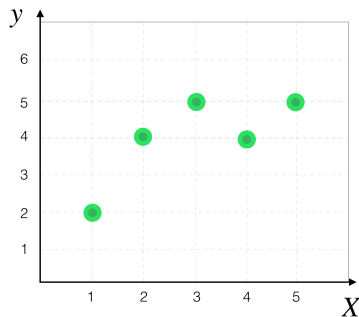
## 2-D example revisited



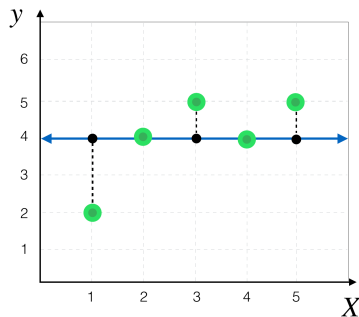
## 2-D example revisited



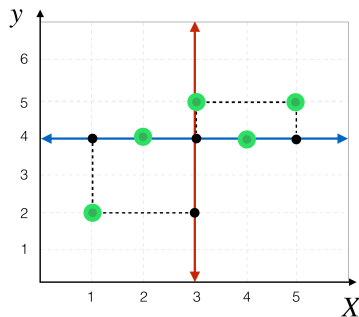
## 2-D example revisited



## 2-D example revisited

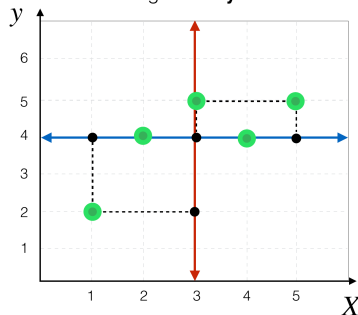


## 2-D example revisited

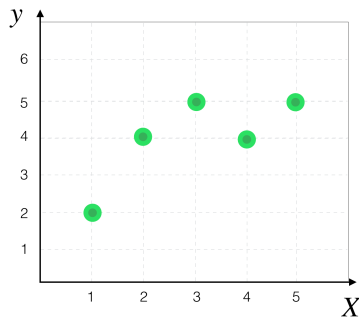


## 2-D example revisited

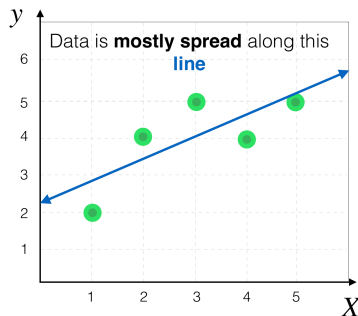
Data seem to be **spread more equally**  
along **X** and **y** axes



## 2-D example revisited

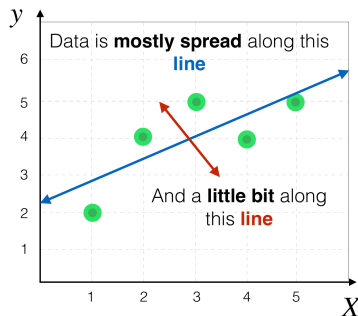


## 2-D example revisited



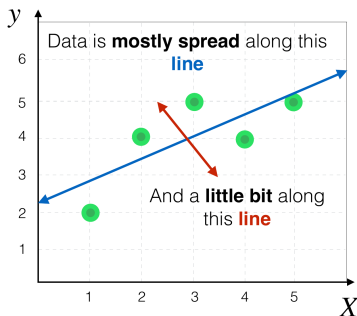


## 2-D example revisited



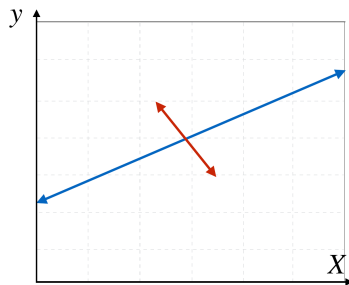
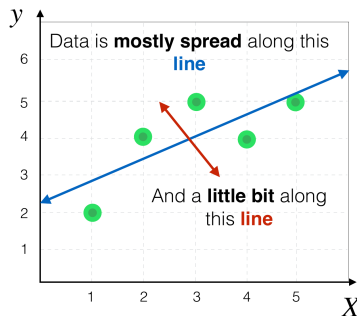
## 2-D example revisited

How about we make **new axes** from **these lines**?



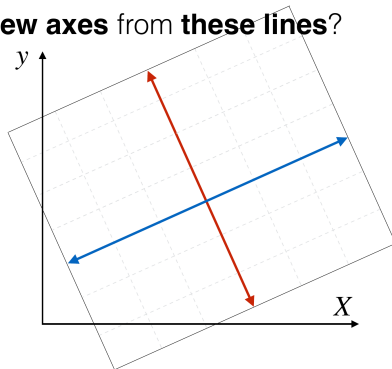
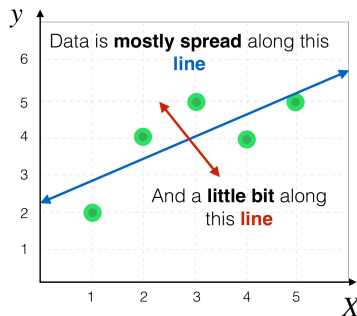
## 2-D example revisited

How about we make **new axes** from **these lines**?



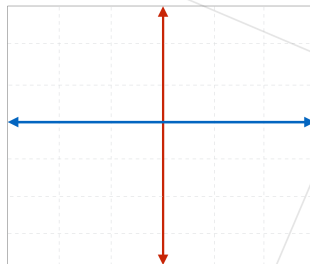
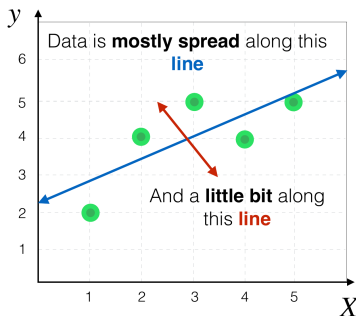
## 2-D example revisited

How about we make **new axes** from **these lines**?



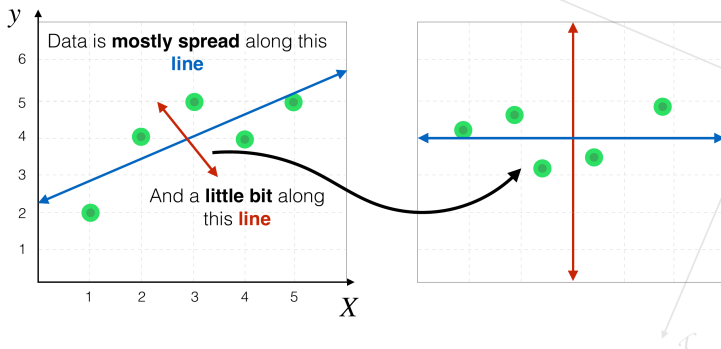
## 2-D example revisited

How about we make **new axes** from **these lines**?



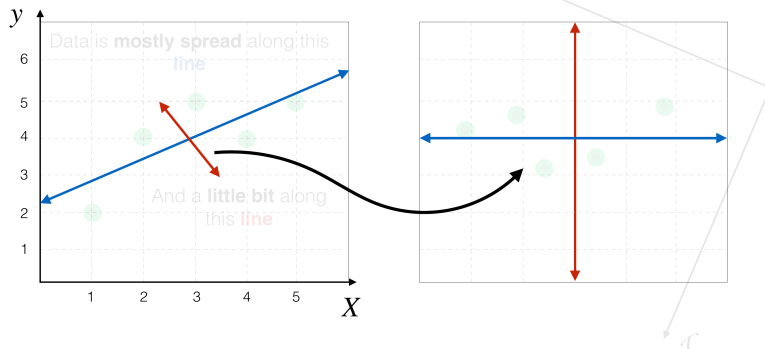
## 2-D example revisited

How about we make **new axes** from **these lines**?



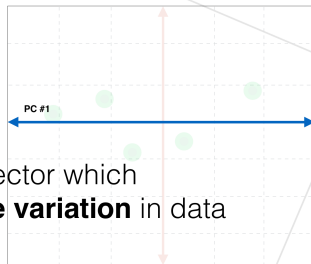
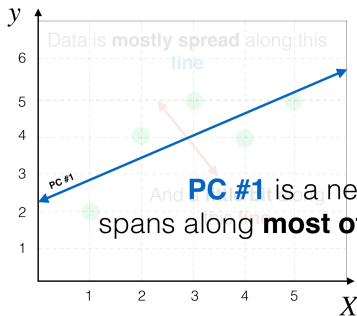
## 2-D example revisited

These **new axes** are called principle components



## 2-D example revisited

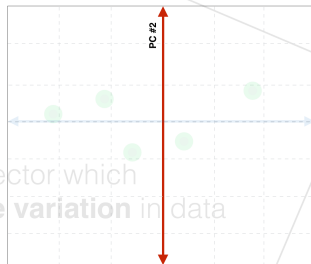
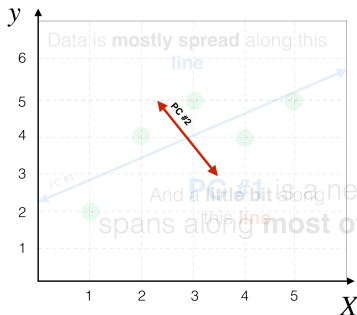
These **new axes** are called principle components





## 2-D example revisited

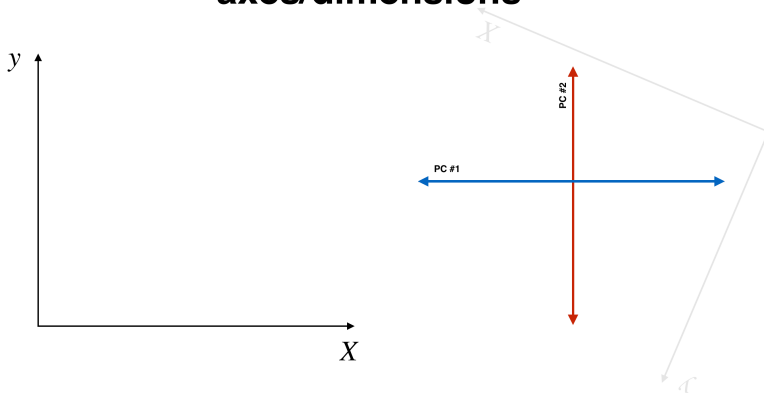
These **new axes** are called principle components



**PC #2** is another new vector which spans along the direction of the second most variation

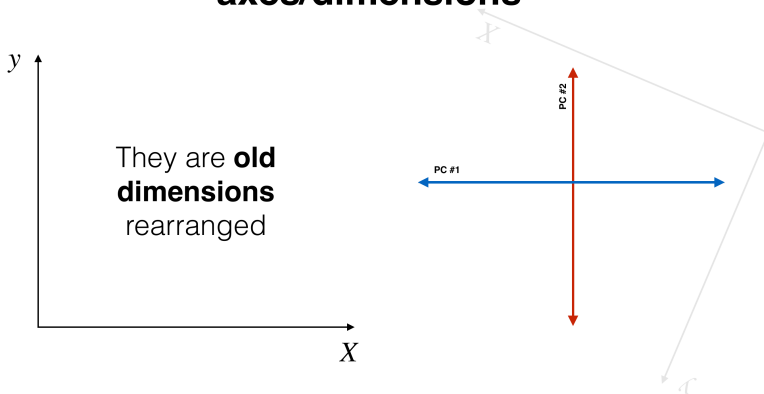
# Principal Component Analysis

Principle components are **not additional axes/dimensions**



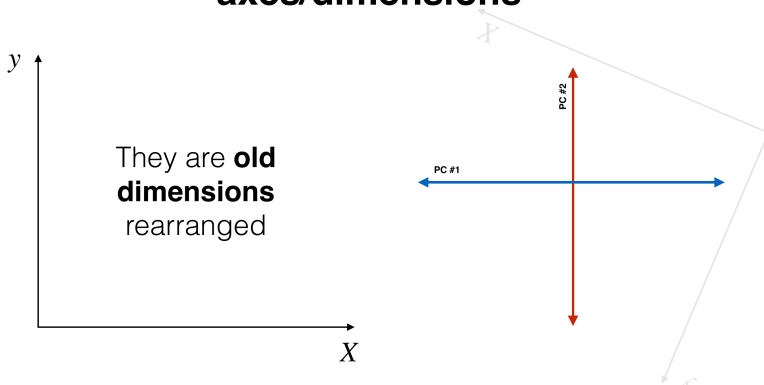
# Principal Component Analysis

Principle components are **not additional axes/dimensions**



# Principal Component Analysis

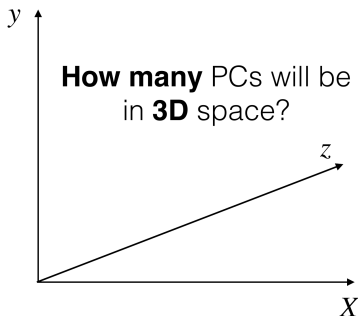
Principle components are **not additional axes/dimensions**



Such that the **first axis** now spans along **most variation**, the **second** the **second most variation** etc.

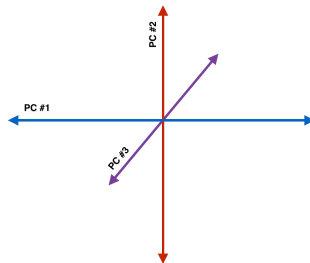
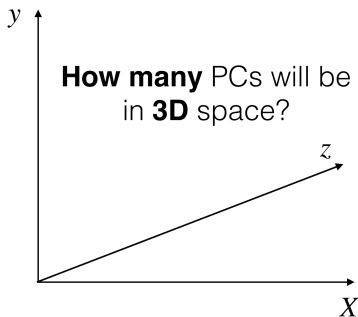
# Principal Component Analysis

Principle components are **not additional axes/dimensions**



# Principal Component Analysis

Principle components are **not additional axes/dimensions**

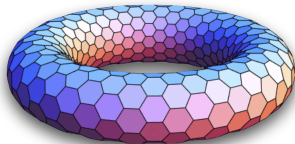


**As many** as there were  
**original dimensions**,  
hence **3** PCs

# Principal Component Analysis

Principle components are **not additional axes/dimensions**

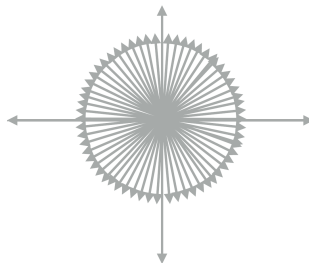
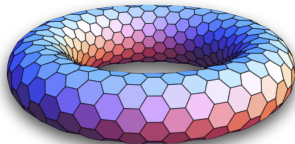
**How many** PCs will be formed in **200D** space?



# Principal Component Analysis

Principle components are **not additional axes/dimensions**

**How many** PCs will be formed in **200D** space?



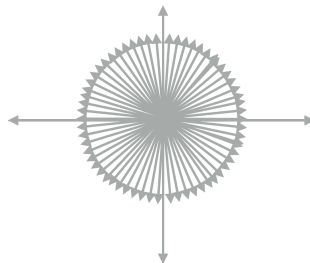
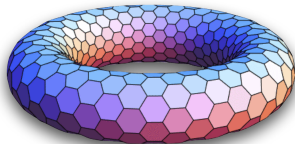
No exceptions, **200 PCs**



# Principal Component Analysis

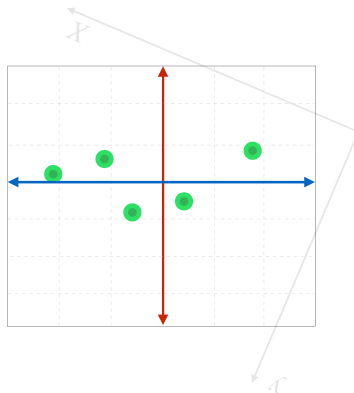
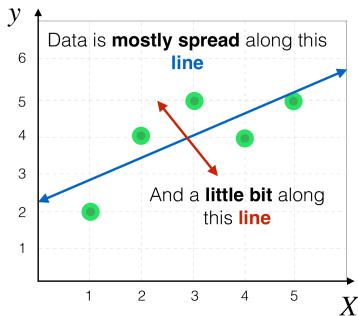
Principle components are **not additional axes/dimensions**

**How many** PCs will be formed in **200D** space?

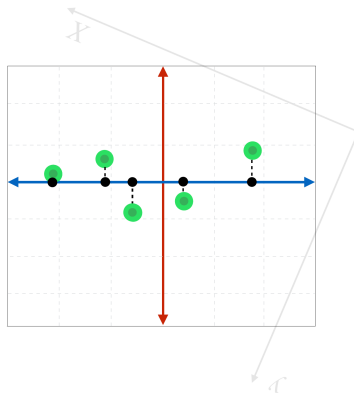
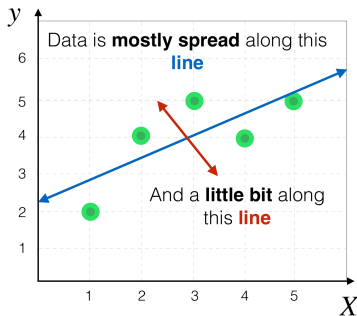


No exceptions, **200 PCs**  
But what is **the benefit** of having PCs?

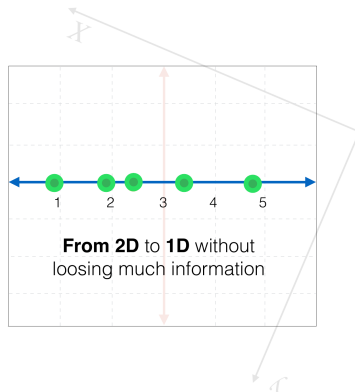
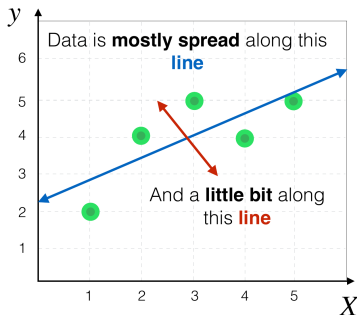
# Principal Component Analysis



# Principal Component Analysis



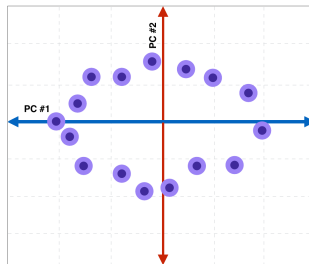
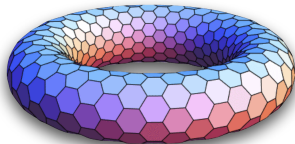
# Principal Component Analysis



# Principal Component Analysis

Principle components are **not additional axes/dimensions**

**How many** PCs will be formed in **200D** space?

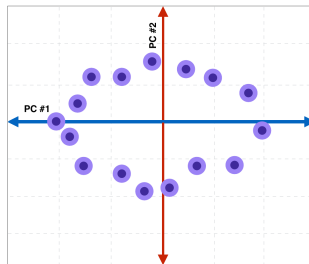
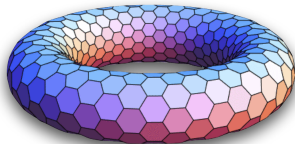


No exceptions, **200 PCs**  
But what is **the benefit** of having PCs?

# Principal Component Analysis

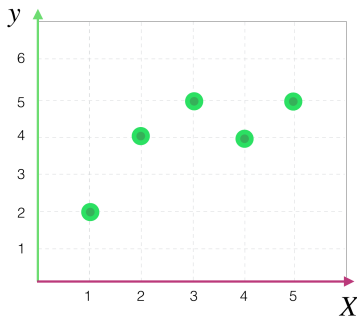
Principle components are **not additional axes/dimensions**

**How many** PCs will be formed in **200D** space?



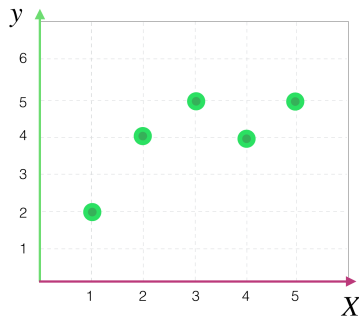
**First few PCs** would be **enough** to capture important information

# Principal Component Analysis



x	y
1	2
2	4
3	5
4	4
5	5

# Principal Component Analysis



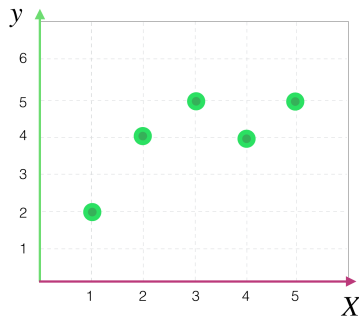
x	y
1	2
2	4
3	5
4	4
5	5

$$\bar{x} = 3$$

$$\bar{y} = 4$$



# Principal Component Analysis



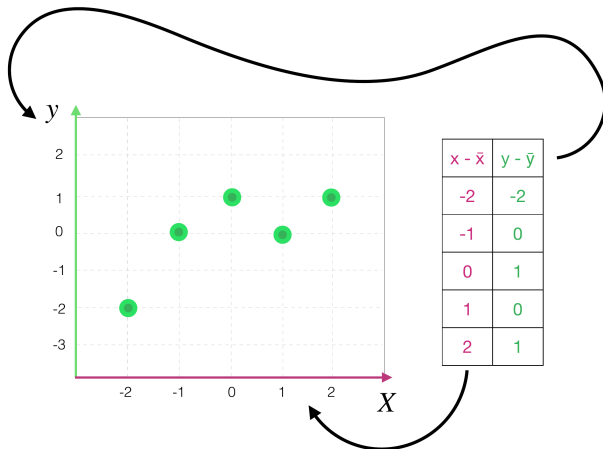
$x$	$y$
1	2
2	4
3	5
4	4
5	5

$$\bar{x} = 3$$

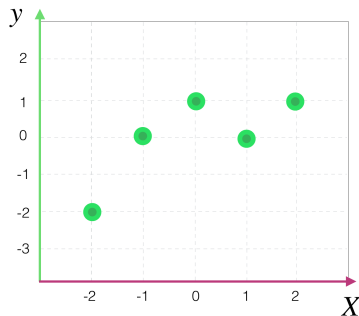
$$\bar{y} = 4$$

$x - \bar{x}$	$y - \bar{y}$
-2	-2
-1	0
0	1
1	0
2	1

# Principal Component Analysis



# Principal Component Analysis



$Z$

$x - \bar{x}$	$y - \bar{y}$
-2	-2
-1	0
0	1
1	0
2	1

# Principal Component Analysis

**Transpose** the matrix of coordinates

$Z$

-2	-2
-1	0
0	1
1	0
2	1

# Principal Component Analysis

**Transpose** the matrix of coordinates

$Z$

-2	-2
-1	0
0	1
1	0
2	1

What are the **dimensions**  
of the transposed matrix?

# Principal Component Analysis

**Transpose** the matrix of coordinates

$Z$

-2	-2
-1	0
0	1
1	0
2	1

$Z^T$



# Principal Component Analysis

**Transpose** the matrix of coordinates

$Z$

-2	-2
-1	0
0	1
1	0
2	1

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

# Principal Component Analysis

$$Z^T \times Z = S$$



# Principal Component Analysis

$$Z^T \times Z = S$$

-2	-1	0	1	2
-2	0	1	0	1

-2	-2
-1	0
0	1
1	0
2	1

?	?
?	?

# Principal Component Analysis

$$Z^T \times Z = S$$

-2	-1	0	1	2
-2	0	1	0	1

-2	-2
-1	0
0	1
1	0
2	1

10	6
6	6

# Principal Component Analysis

$$Z^T \times Z = S \times 4$$

-2	-1	0	1	2
-2	0	1	0	1

-2	-2
-1	0
0	1
1	0
2	1

2.5	1.5
1.5	1.5

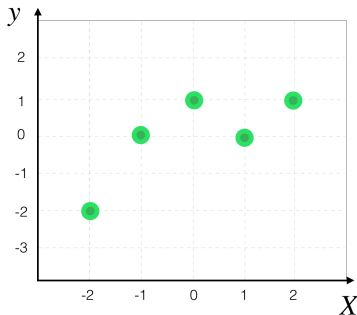
# Principal Component Analysis

$$\begin{array}{ccccc} & & Z^T & & \\ \hline & -2 & -1 & 0 & 1 & 2 \\ \hline & -2 & 0 & 1 & 0 & 1 \\ \hline \end{array} \times \begin{array}{cc} & Z \\ \hline -2 & -2 \\ -1 & 0 \\ 0 & 1 \\ 1 & 0 \\ 2 & 1 \\ \hline \end{array} = \begin{array}{cc} & S \\ \hline 2.5 & 1.5 \\ 1.5 & 1.5 \\ \hline \end{array} \times 4$$

**Covariance** matrix

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

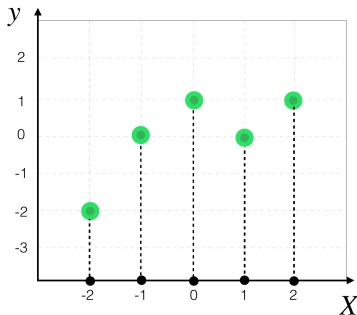

$$S$$

2.5	1.5
1.5	1.5

**Covariance** matrix

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

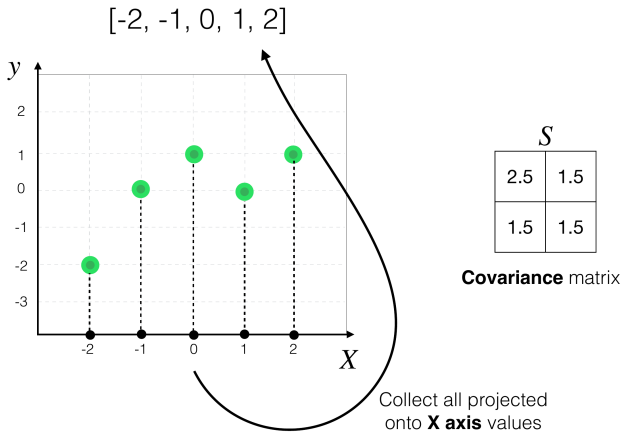

$$S$$

2.5	1.5
1.5	1.5

**Covariance** matrix

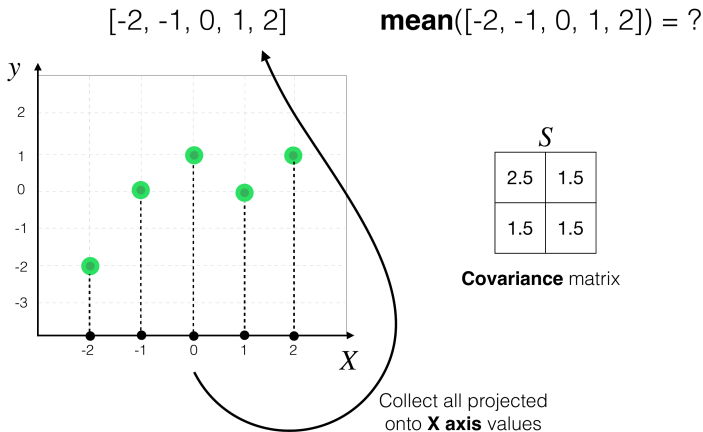
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



# Principal Component Analysis

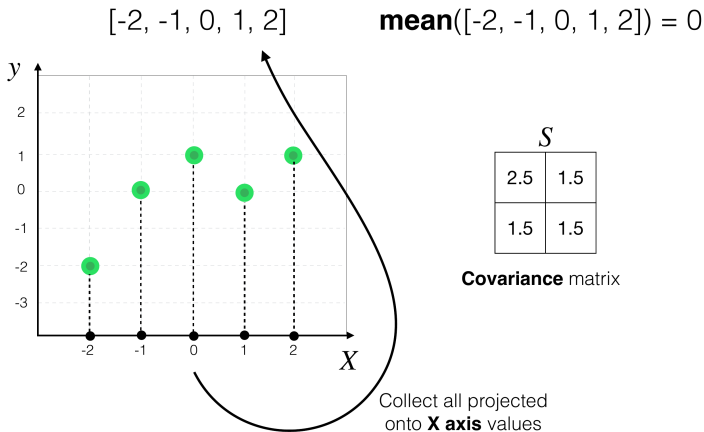
How to **interpret** values in **covariance** matrix?





# Principal Component Analysis

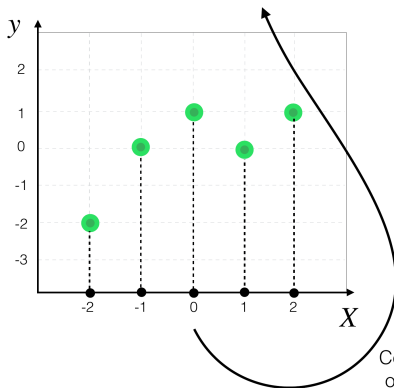
How to **interpret** values in **covariance** matrix?



# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $\bar{\mathbf{x}} = 0$


$$S$$

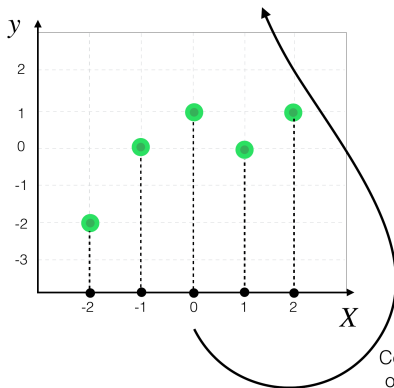
2.5	1.5
1.5	1.5

**Covariance** matrix

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

$$[-2, -1, 0, 1, 2] \quad \bar{\mathbf{x}} = 0 \quad \sigma = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$



$S$	
2.5	1.5
1.5	1.5

**Covariance** matrix

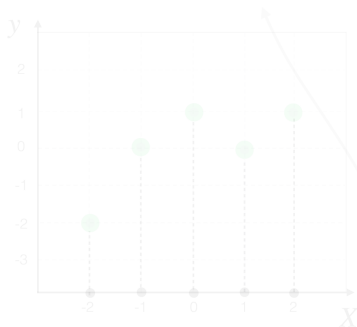
Collect all projected  
onto **X axis** values

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $\bar{x} = 0$

$$\sigma = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$



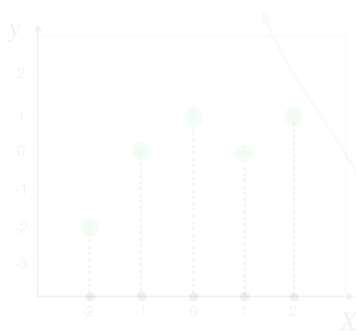
**Covariance** matrix

Collect all projected  
onto **X axis** values

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

**Variance**  $\rightarrow \sigma = \frac{\sum (x_i - \bar{x})^2}{n - 1}$

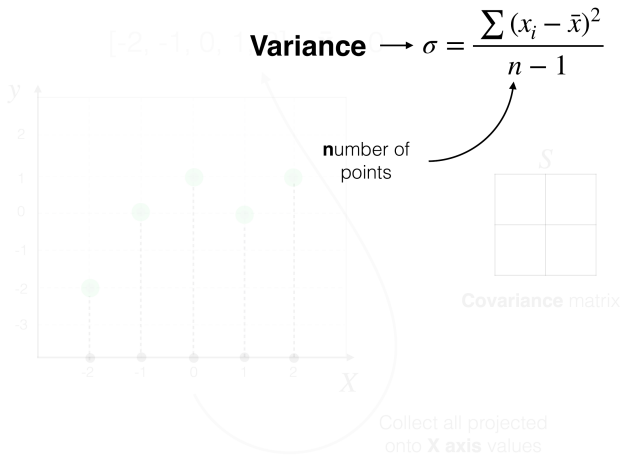


**Covariance** matrix

Collect all projected  
onto **X axis** values

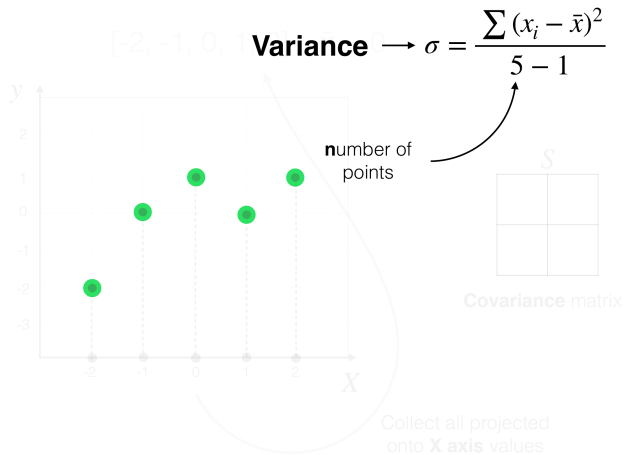
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



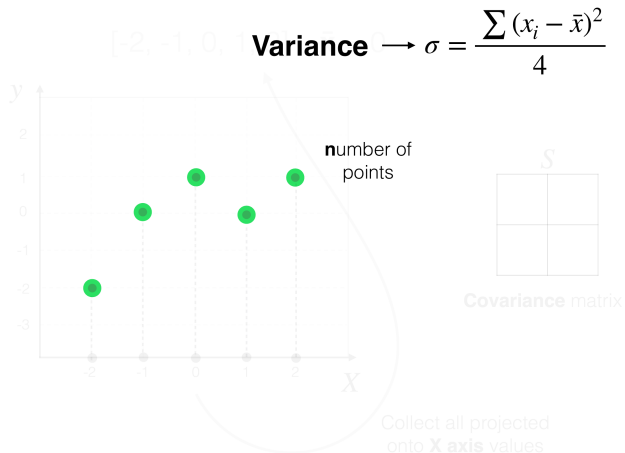
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



# Principal Component Analysis

How to **interpret** values in **covariance** matrix?





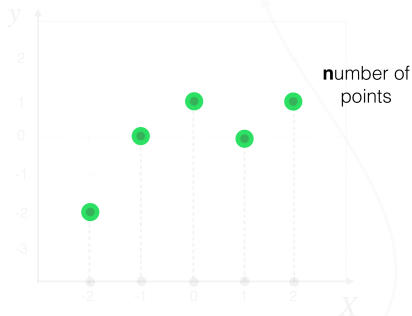
# Principal Component Analysis

How to **interpret** values in **covariance matrix**

**Variance**  $\rightarrow \sigma = \frac{\sum (x_i - \bar{x})^2}{4}$

mean of all  
points

$$\bar{x} = 0$$



**Covariance matrix**

Collect all projected  
onto **X axis** values

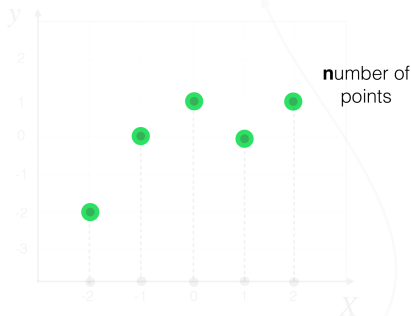
# Principal Component Analysis

How to **interpret** values in **covariance matrix**

mean of all points

$$\bar{\mathbf{x}} = 0$$

$$\text{Variance} \rightarrow \sigma = \frac{\sum (x_i - 0)^2}{4}$$



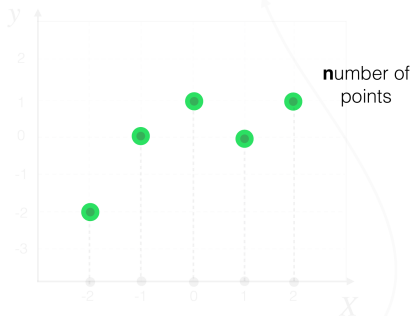
Covariance matrix

Collect all projected  
onto **X axis** values

# Principal Component Analysis

How to **interpret** values in **covariance**  $\bar{\mathbf{x}} = 0$  **mean of all points**

$$\text{Variance} \rightarrow \sigma = \frac{\sum (x_i)^2}{4}$$



Covariance matrix

Collect all projected  
onto **X axis** values

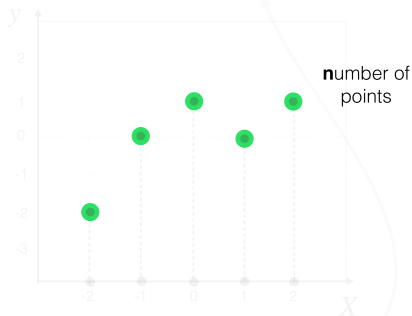
# Principal Component Analysis

How to interpret  $\bar{x}$  in covariance matrix  $S$

value of each point

$\bar{x} = 0$  mean of all points

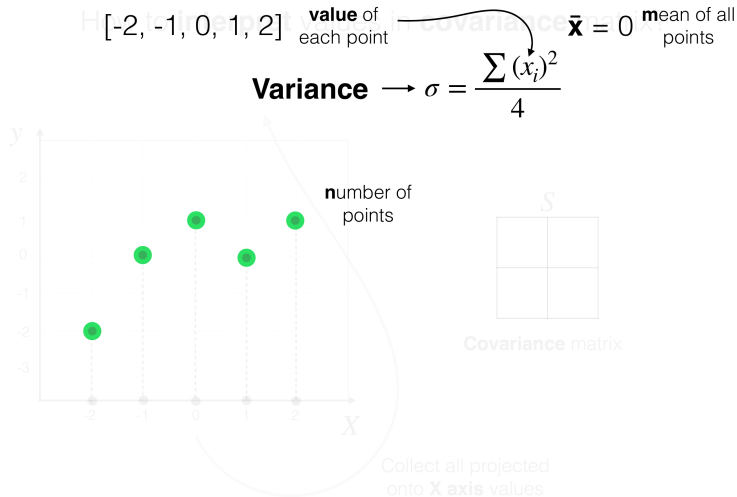
**Variance**  $\rightarrow \sigma = \frac{\sum (x_i)^2}{4}$



Covariance matrix

Collect all projected  
onto **X axis** values

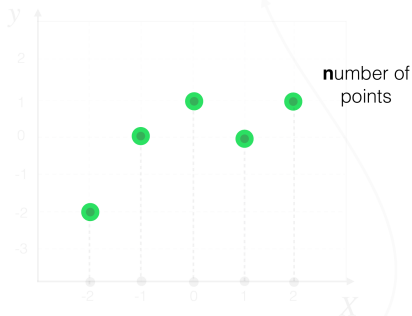
# Principal Component Analysis



# Principal Component Analysis

$[-2, -1, 0, 1, 2]$  value of each point  $\bar{\mathbf{x}} = 0$  mean of all points

**Variance**  $\rightarrow \sigma = \frac{(-2)^2 + (-1)^2 + (0)^2 + (1)^2 + (2)^2}{4}$



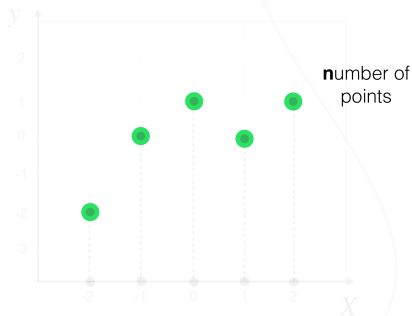
Covariance matrix

Collect all projected  
onto **X** axis values

# Principal Component Analysis

[-2, -1, 0, 1, 2] **value** of each point  $\bar{\mathbf{x}} = 0$  **mean** of all points

$$\text{Variance} \rightarrow \sigma = \frac{10}{4}$$



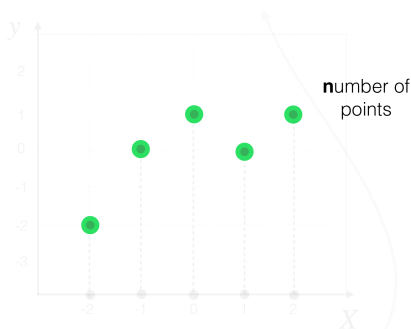
Covariance matrix

Collect all projected  
onto **X axis** values

# Principal Component Analysis

[-2, -1, 0, 1, 2] value of each point  $\bar{\mathbf{x}} = 0$  mean of all points

Variance  $\rightarrow \sigma = 2.5$



Covariance matrix

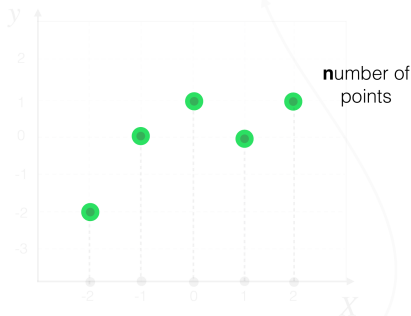
Collect all projected  
onto **X** axis values



# Principal Component Analysis

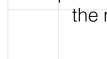
[-2, -1, 0, 1, 2] **value of each point**  $\bar{x} = 0$  **mean of all points**

**Variance**  $\rightarrow \sigma = 2.5$



$$\sigma = \frac{\sum (x_i - \bar{x})^2}{n - 1}$$

**Variance** is an expected value of the squared deviation from the mean



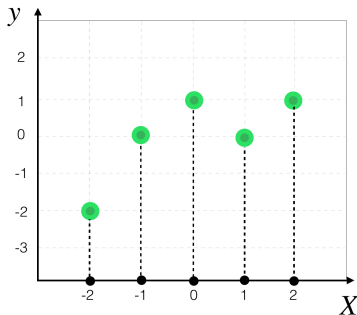
**Covariance matrix**

Collect all projected onto **X axis** values

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

$$[-2, -1, 0, 1, 2] \quad \bar{\mathbf{x}} = 0 \quad \sigma = 2.5$$


$$S$$

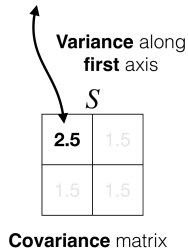
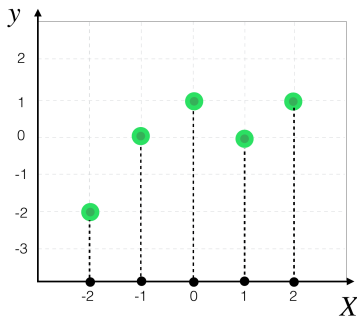
2.5	1.5
1.5	1.5

**Covariance** matrix

# Principal Component Analysis

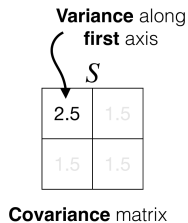
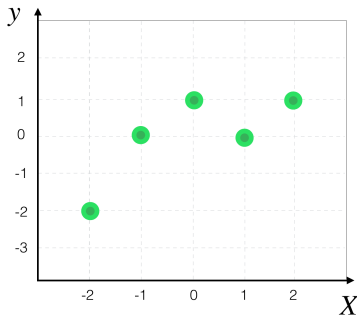
How to **interpret** values in **covariance** matrix?

$$[-2, -1, 0, 1, 2] \quad \bar{\mathbf{x}} = 0 \quad \sigma = 2.5$$



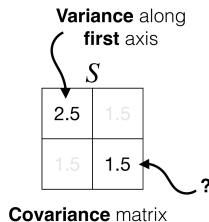
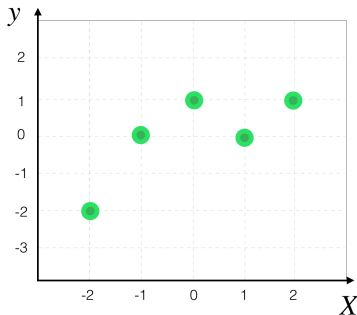
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



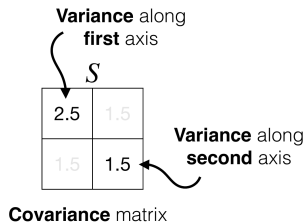
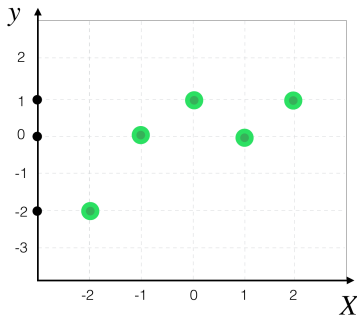
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



# Principal Component Analysis

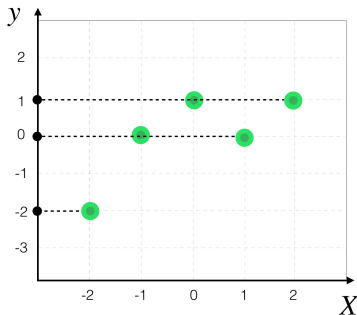
How to **interpret** values in **covariance** matrix?

**value** of  
each point

**number** of  
points

**mean** of all  
points

$$\sigma = \frac{\sum (y_i - \bar{y})^2}{n - 1}$$



Variance along  
first axis

$S$

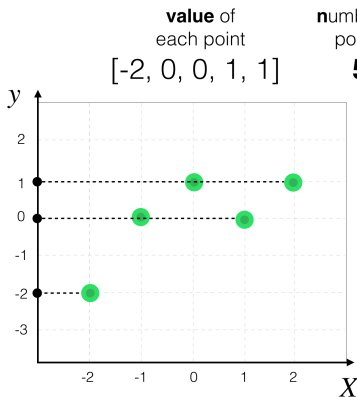
2.5	1.5
1.5	1.5

Variance along  
second axis

**Covariance** matrix

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



mean of all points  
 $\bar{y} = 0$

$$\sigma = \frac{\sum (y_i - \bar{y})^2}{n - 1}$$

Variance along first axis

$S$

2.5	1.5
1.5	1.5

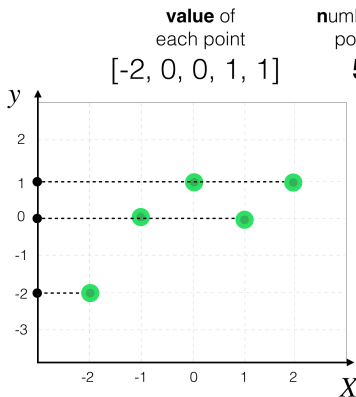
Variance along second axis

**Covariance** matrix

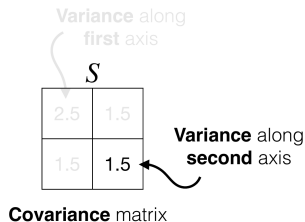


# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

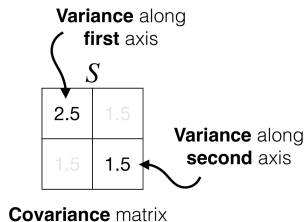
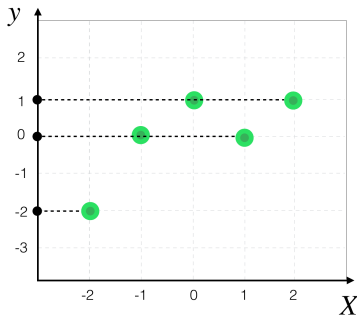


mean of all points  
 $\bar{y} = 0$

$$\sigma = \frac{6}{4}$$


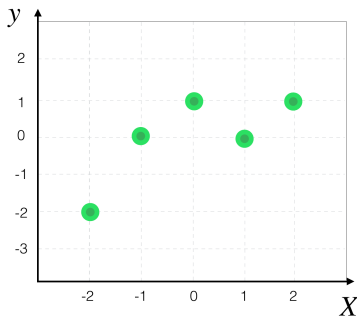
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

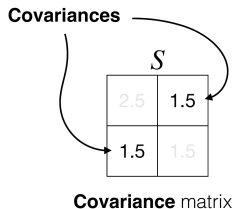
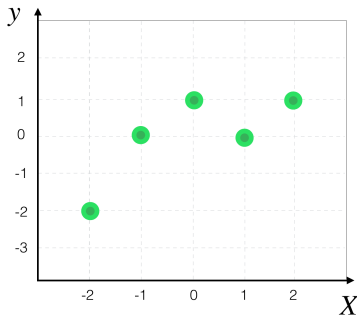

$$S$$

2.5	1.5
1.5	1.5

**Covariance** matrix

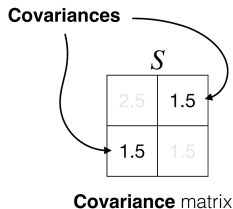
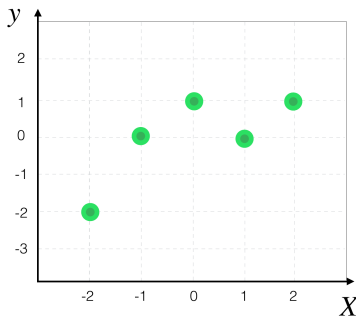
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



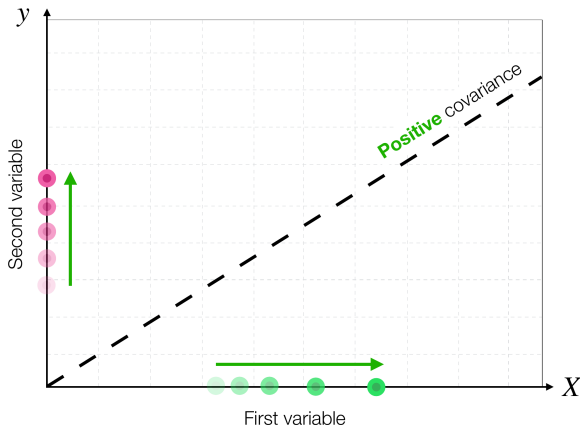
# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

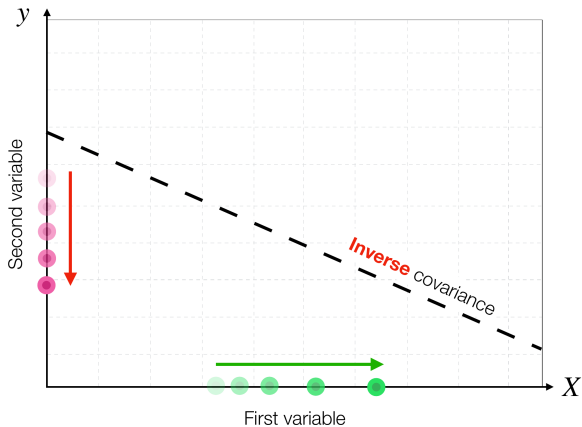


Covariance indicates how two variables are related. A **positive** covariance means the variables are **positively related**, while a **negative** covariance means the variables are **inversely related**.

# Principal Component Analysis

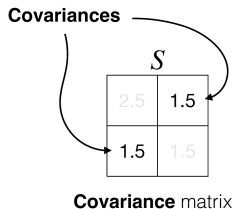
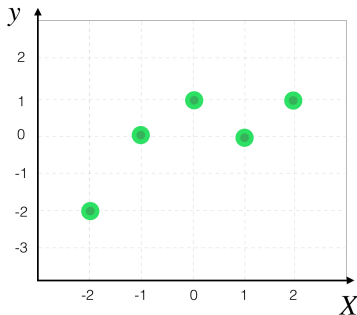


# Principal Component Analysis



# Principal Component Analysis

How to **interpret** values in **covariance** matrix?





# Principal Component Analysis

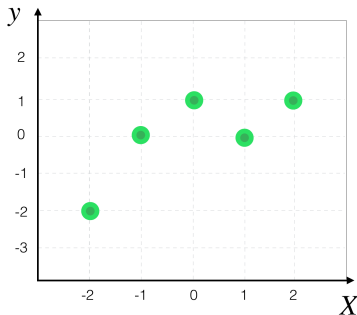
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[-2, 0, 0, 1, 1]$

$\bar{x} = 0$

$\bar{y} = 0$

$$\text{cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$



**Covariances**

$S$

2.5	1.5
1.5	1.5

**Covariance** matrix

# Principal Component Analysis

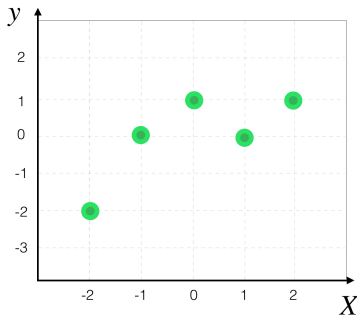
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[-2, 0, 0, 1, 1]$

$$\bar{x} = 0$$

$$\bar{y} = 0$$

$$\text{cov}(x, y) = \frac{6}{4}$$



**Covariances**

$S$

2.5	1.5
1.5	1.5

**Covariance** matrix

# Principal Component Analysis

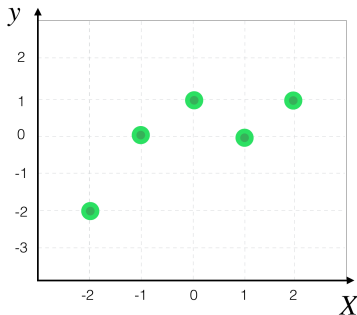
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[-2, 0, 0, 1, 1]$

$$\bar{x} = 0$$

$$\bar{y} = 0$$

$$\text{cov}(x, y) = 1.5$$



Covariances

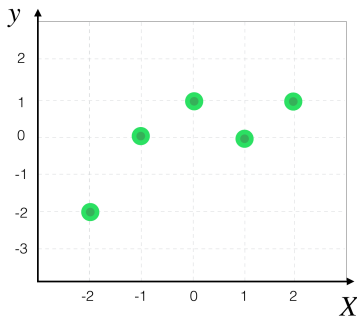
$S$

2.5	1.5
1.5	1.5

Covariance matrix

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?


$$S$$

2.5	1.5
1.5	1.5

**Covariance** matrix

# Principal Component Analysis

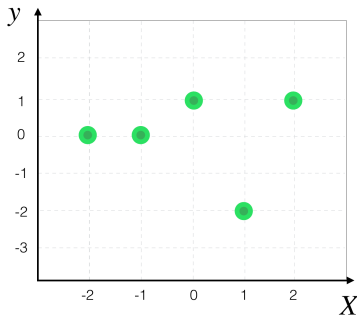
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[0, 0, 1, -2, 1]$

$\bar{\mathbf{x}} = ?$

$\bar{\mathbf{y}} = ?$

$$\text{cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$



$S$

2.5	?
?	1.5

**Covariance** matrix

# Principal Component Analysis

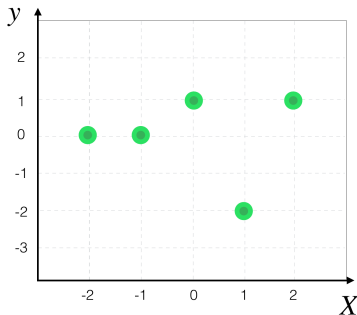
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[0, 0, 1, -2, 1]$

$\bar{x} = 0$

$\bar{y} = 0$

$$\text{cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n - 1}$$



$S$

2.5	?
?	1.5

**Covariance** matrix

# Principal Component Analysis

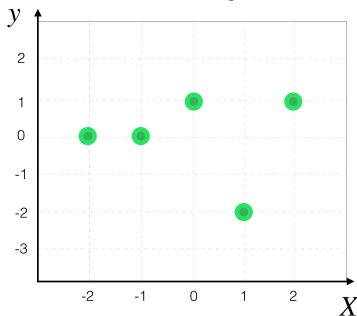
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[0, 0, 1, -2, 1]$

$$\bar{x} = 0$$

$$\bar{y} = 0$$

$$\text{cov}(x, y) = \frac{\sum (x_i)(y_i)}{4}$$



S	
2.5	?
?	1.5

**Covariance** matrix

# Principal Component Analysis

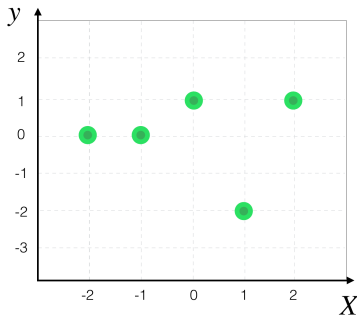
How to **interpret** values in **covariance** matrix?

[-2, -1, 0, 1, 2] [0, 0, 1, -2, 1]

$$\bar{x} = 0$$

$$\bar{y} = 0$$

$$\text{cov}(x, y) = \frac{(-2)(0) + (-1)(0) + (0)(1) + (1)(-2) + (2)(1)}{4}$$



$S$	
2.5	?
?	1.5

**Covariance** matrix



# Principal Component Analysis

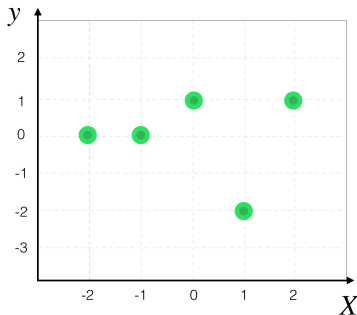
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[0, 0, 1, -2, 1]$

$$\bar{x} = 0$$

$$\bar{y} = 0$$

$$\text{cov}(x, y) = \frac{0}{4} = 0$$



S	
2.5	?
?	1.5

**Covariance** matrix

# Principal Component Analysis

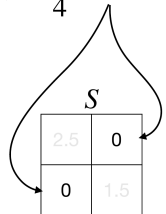
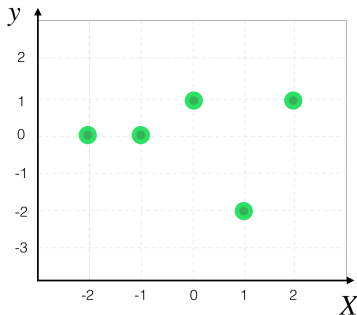
How to **interpret** values in **covariance** matrix?

$[-2, -1, 0, 1, 2]$   $[0, 0, 1, -2, 1]$

$$\bar{x} = 0$$

$$\bar{y} = 0$$

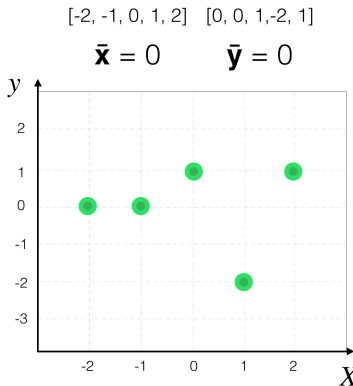
$$\text{cov}(x, y) = \frac{0}{4} = 0$$



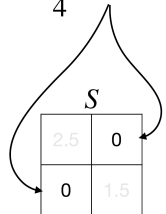
**Covariance** matrix

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



$$\text{cov}(x, y) = \frac{0}{4} = 0$$

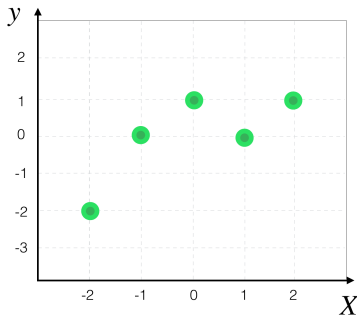


**Covariance** matrix

Covariance **0** means that there is **no relationship** between two variables. Knowing something about the value of one does not say anything about the value of the other.

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?

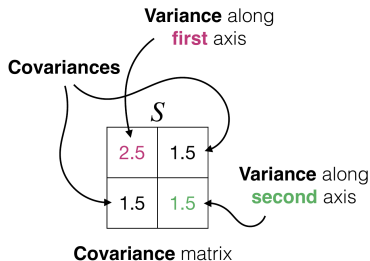
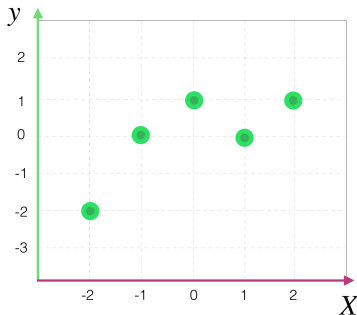

$$S$$

2.5	1.5
1.5	1.5

**Covariance** matrix

# Principal Component Analysis

How to **interpret** values in **covariance** matrix?



# Principal Component Analysis

$$Z^T \times Z = S \times 4$$

-2	-1	0	1	2
-2	0	1	0	1

-2	-2
-1	0
0	1
1	0
2	1

2.5	1.5
1.5	1.5

# Principal Component Analysis

$$S = PDP^T$$

2.5	1.5
1.5	1.5

# Principal Component Analysis

$$S = PDP^T$$

2.5	1.5
1.5	1.5

2.5	1.5
1.5	1.5

For an example: <https://www.scss.tcd.ie/Rozenn.Dahyot/CS1BA1/SolutionEigen.pdf>



# Principal Component Analysis

$$S = PD P^T$$

2.5	1.5
1.5	1.5

-2.9	0.24
-2.1	-0.33

-0.81	-0.58
0.58	-0.81

For an example: <https://www.scss.tcd.ie/Rozenn.Dahyot/CS1BA1/SolutionEigen.pdf>

# Principal Component Analysis

$$S = P D P^T$$

2.5	1.5
1.5	1.5

-0.81	0.58
-0.58	-0.81

3.58	0
0	0.41

-0.81	-0.58
0.58	-0.81

For an example: <https://www.scss.tcd.ie/Rozenn.Dahyot/CS1BA1/SolutionEigen.pdf>

## Eigendecomposition

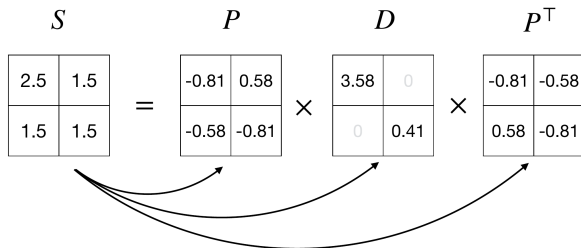
$$S = P D P^T$$

2.5	1.5
1.5	1.5

-0.81	0.58
-0.58	-0.81

3.58	0
0	0.41

-0.81	-0.58
0.58	-0.81



For an example: <https://www.scss.tcd.ie/Rozenn.Dahyot/CS1BA1/SolutionEigen.pdf>

## Eigendecomposition

Eigen**vectors**

$P$                        $D$                        $P^T$

$$\begin{array}{|c|c|} \hline 2.5 & 1.5 \\ \hline 1.5 & 1.5 \\ \hline \end{array} = \begin{array}{|c|c|} \hline -0.81 & 0.58 \\ \hline -0.58 & -0.81 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline 3.58 & 0 \\ \hline 0 & 0.41 \\ \hline \end{array} \times \begin{array}{|c|c|} \hline -0.81 & -0.58 \\ \hline 0.58 & -0.81 \\ \hline \end{array}$$

Eigen**values**

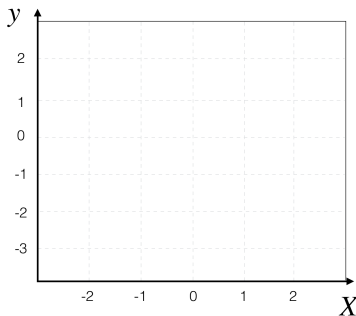
The diagram illustrates the eigendecomposition of a symmetric matrix  $S$ . The matrix  $S$  is shown as a 2x2 grid with values 2.5, 1.5, 1.5, and 1.5. It is equal to the product of three matrices:  $P$ ,  $D$ , and  $P^T$ . Matrix  $P$  is a 2x2 grid with values -0.81, 0.58, -0.58, and -0.81. Matrix  $D$  is a 2x2 grid with values 3.58, 0, 0, and 0.41. Matrix  $P^T$  is a 2x2 grid with values -0.81, -0.58, 0.58, and -0.81. An arrow points from the label 'Eigenvalues' to the diagonal matrix  $D$ . Another arrow points from the label 'Eigen vectors' to the matrix  $P$ , and a third arrow points from the same label to the matrix  $P^T$ .

For an example: <https://www.scss.tcd.ie/Rozenn.Dahyot/CS1BA1/SolutionEigen.pdf>

# Principal Component Analysis

Eigen**vectors**

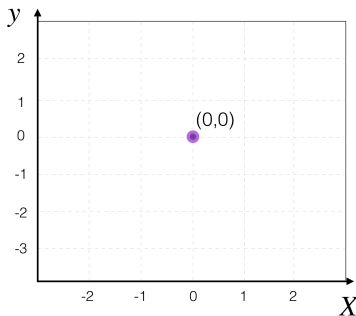
-0.81	0.58
-0.58	-0.81



# Principal Component Analysis

Eigen**vectors**

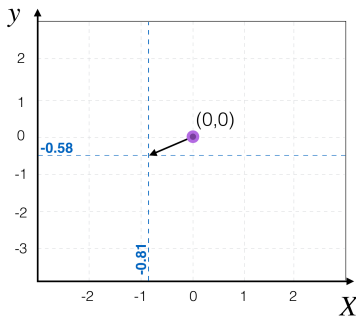
-0.81	0.58
-0.58	-0.81



# Principal Component Analysis

Eigen**vectors**

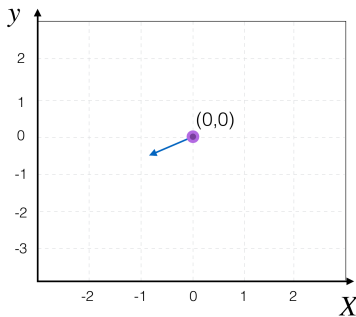
-0.81	0.58
-0.58	-0.81



# Principal Component Analysis

Eigen**vectors**

-0.81	0.58
-0.58	-0.81

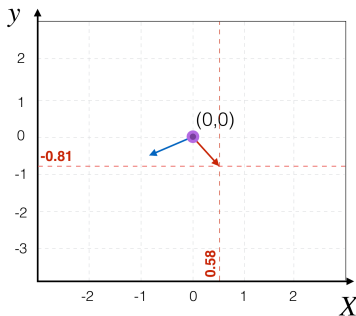




# Principal Component Analysis

Eigen**vectors**

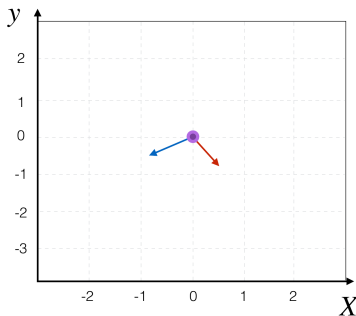
-0.81	0.58
-0.58	-0.81



# Principal Component Analysis

Eigen**vectors**

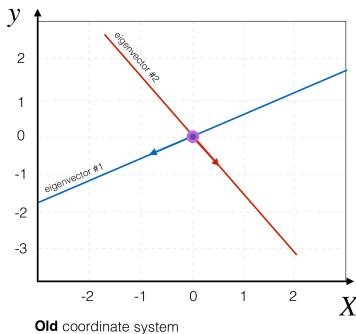
-0.81	0.58
-0.58	-0.81



# Principal Component Analysis

Eigenvectors

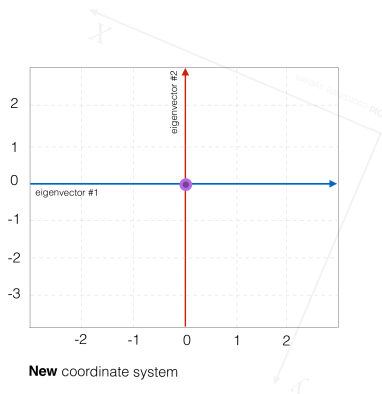
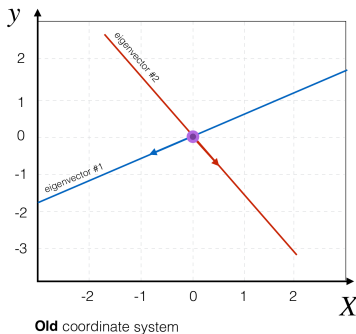
-0.81	0.58
-0.58	-0.81



# Principal Component Analysis

Eigenvectors

-0.81	0.58
-0.58	-0.81



# Principal Component Analysis

Eigen**vectors**

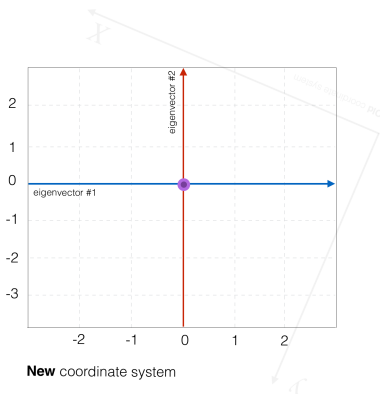
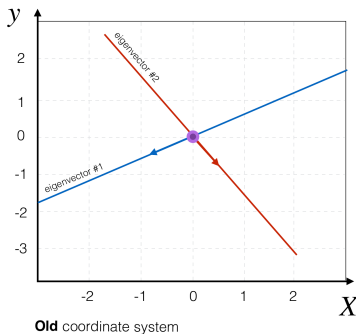
-0.81	0.58
-0.58	-0.81

transpose

=

Eigen**vectors**<sup>T</sup>

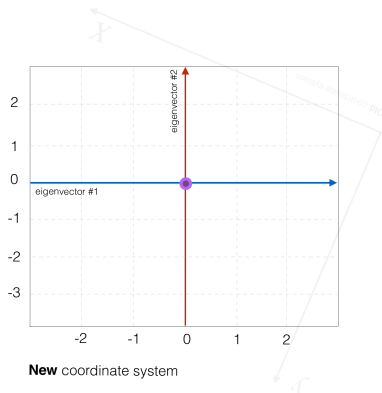
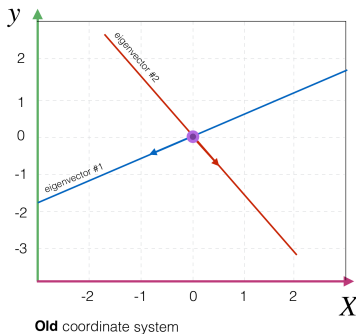
-0.81	-0.58
0.58	-0.81



# Principal Component Analysis

Eigen**vectors**<sup>T</sup>

-0.81	-0.58
0.58	-0.81



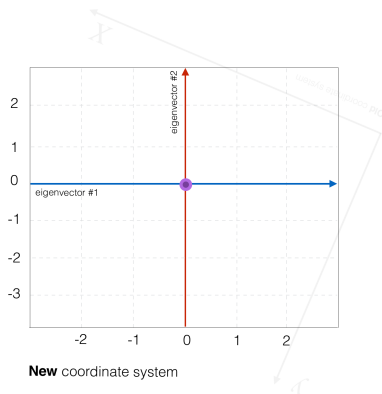
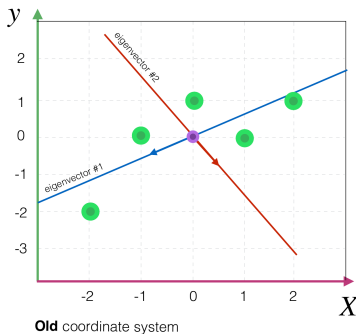
# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1



# Principal Component Analysis

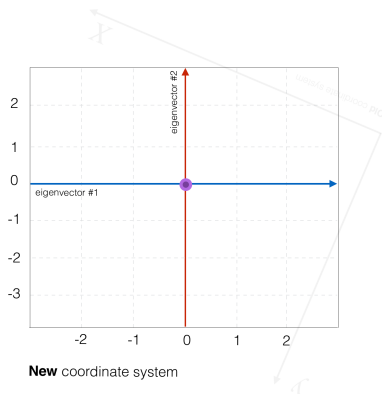
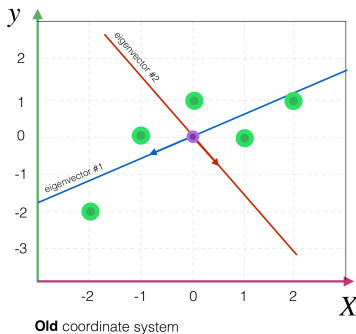
Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1





# Principal Component Analysis

Eigenvectors<sup>T</sup>

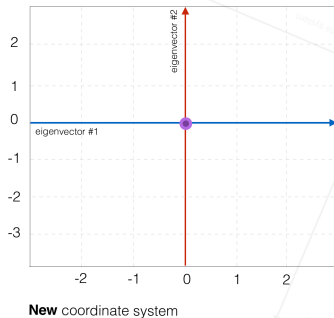
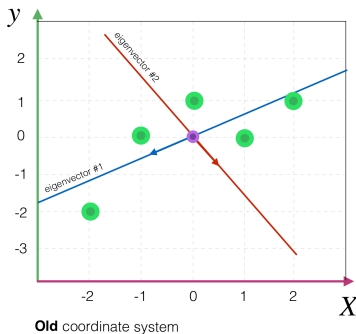
-0.81	-0.58
0.58	-0.81

×

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

=

# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

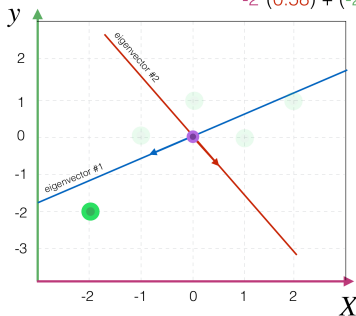
$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

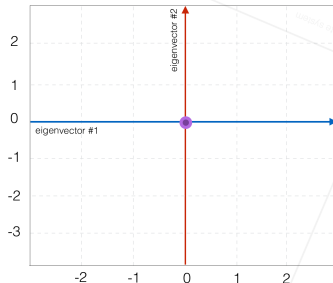
=


$$-2*(-0.81) + (-2)*(-0.58) = \mathbf{2.78}$$

$$-2*(0.58) + (-2)*(-0.81) = \mathbf{0.46}$$



Old coordinate system



New coordinate system

# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

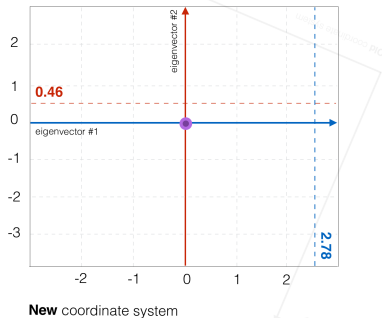
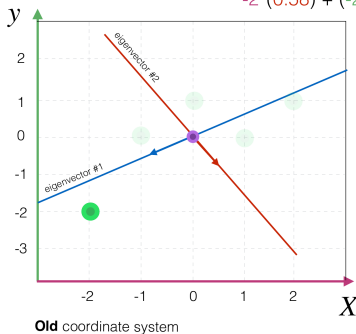
$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

=


$$-2*(-0.81) + (-2)*(-0.58) = 2.78$$

$$-2*(0.58) + (-2)*(-0.81) = 0.46$$



# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

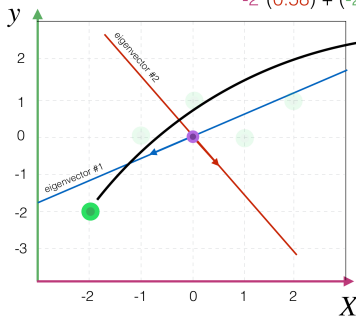
$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

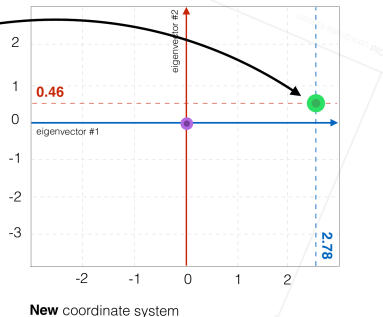
=


$$-2*(-0.81) + (-2)*(-0.58) = 2.78$$

$$-2*(0.58) + (-2)*(-0.81) = 0.46$$



Old coordinate system



New coordinate system

# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

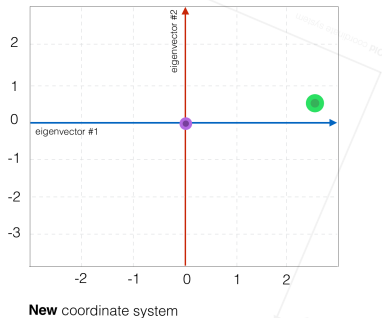
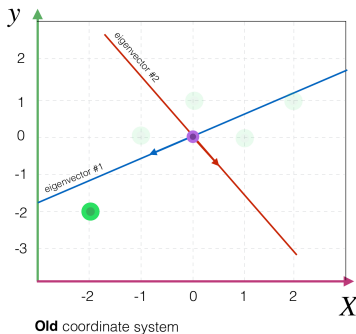
×

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46



# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

$Z^T$

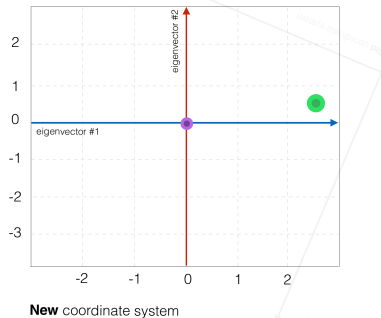
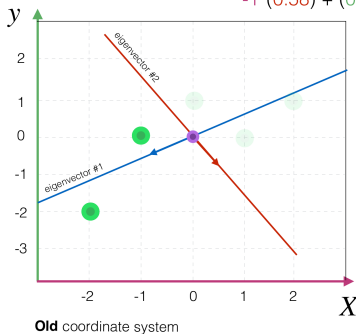
-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46

$$-1 * (-0.81) + (0) * (-0.58) = \mathbf{0.81}$$

$$-1 * (0.58) + (0) * (-0.81) = \mathbf{-0.58}$$



# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

$Z^T$

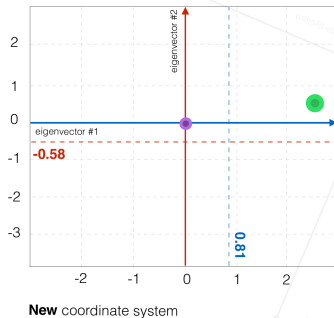
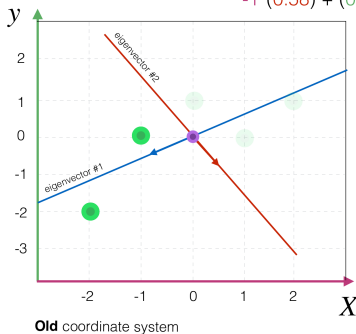
-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46

$$-1*(-0.81) + (0)*(-0.58) = \mathbf{0.81}$$

$$-1*(0.58) + (0)*(-0.81) = \mathbf{-0.58}$$



# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

$Z^T$

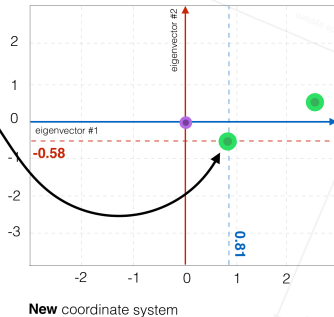
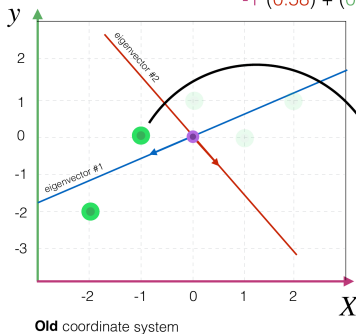
-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46
0.81	-0.58

$$-1 * (-0.81) + (0) * (-0.58) = 0.81$$

$$-1 * (0.58) + (0) * (-0.81) = -0.58$$





# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

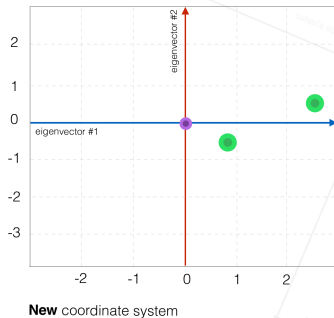
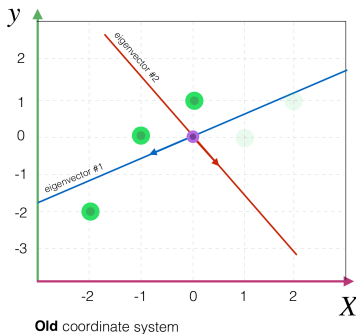
×

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46
0.81	-0.58



# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

$Z^T$

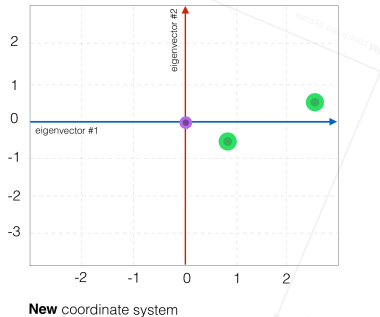
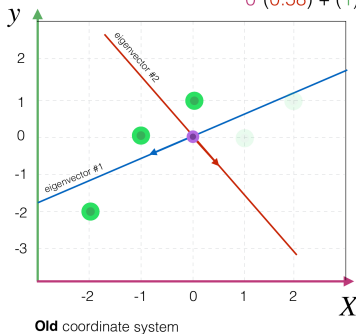
-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46
0.81	-0.58

$$0 * (-0.81) + (1) * (-0.58) = -0.58$$

$$0 * (0.58) + (1) * (-0.81) = -0.81$$



# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

×

$Z^T$

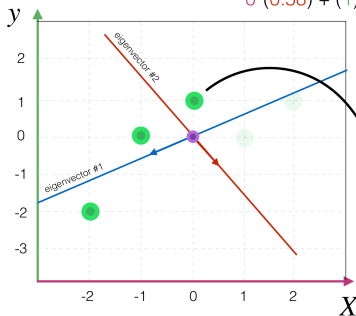
-2	-1	0	1	2
-2	0	1	0	1

=

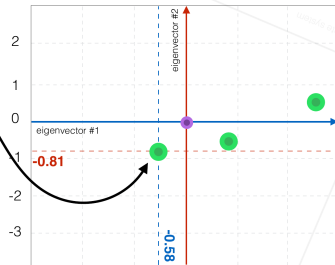
2.78	0.46
0.81	-0.58
-0.58	-0.81

$$0 * (-0.81) + (1) * (-0.58) = -0.58$$

$$0 * (0.58) + (1) * (-0.81) = -0.81$$



Old coordinate system



New coordinate system

# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

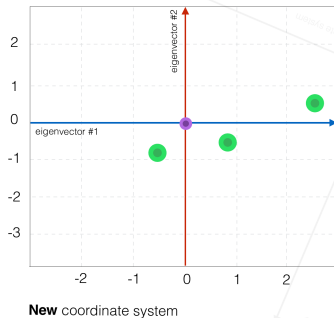
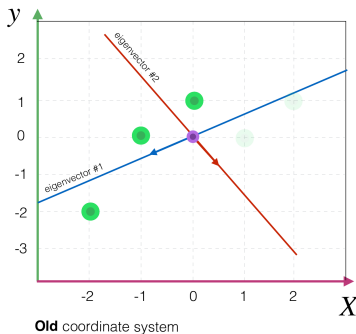
×

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46
0.81	-0.58
-0.58	-0.81



# Principal Component Analysis

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

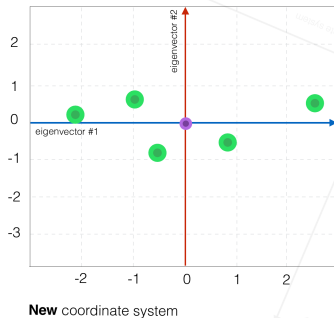
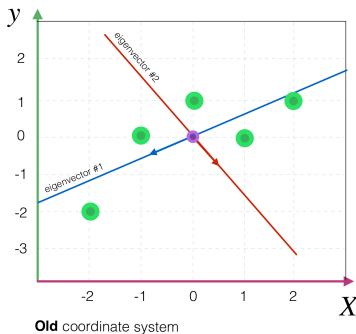
×

$Z^T$

-2	-1	0	1	2
-2	0	1	0	1

=

2.78	0.46
0.81	-0.58
-0.58	-0.81
-0.81	0.58
-2.2	0.35



# Principal Component Analysis

These eigenvectors are called **Principle Components**

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

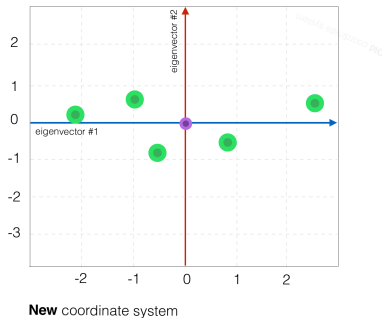
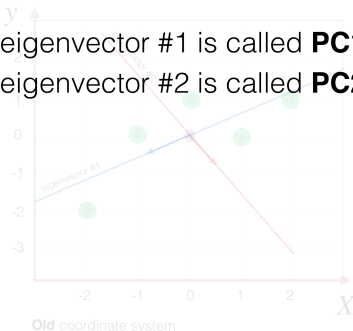
$Z^T$

-2	-1	0	1	2
2	0	1	0	1

=

2.78	0.46
0.81	-0.58
-0.58	-0.81
-0.81	0.58
-2.2	0.35

eigenvector #1 is called **PC1**  
eigenvector #2 is called **PC2**



# Principal Component Analysis

These eigenvectors are called **Principle Components**

Eigenvectors<sup>T</sup>

-0.81	-0.58
0.58	-0.81

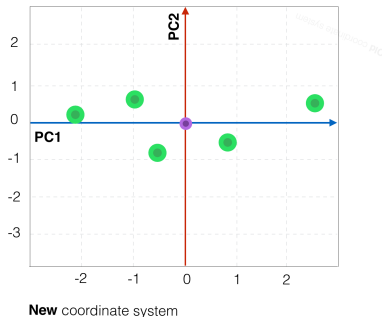
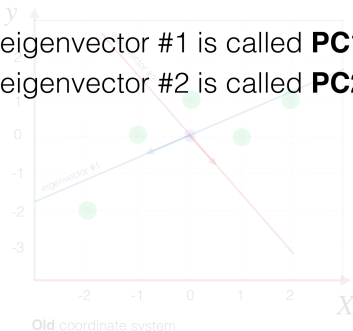
$Z^T$

-2	-1	0	1	2
2	0	1	0	1

=

2.78	0.46
0.81	-0.58
-0.58	-0.81
-0.81	0.58
-2.2	0.35

eigenvector #1 is called **PC1**  
eigenvector #2 is called **PC2**



# Principal Component Analysis

$$Z^T \times Z = S$$

-2	-1	0	1	2
-2	0	1	0	1

-2	-2
-1	0
0	1
1	0
2	1

2.5	1.5
1.5	1.5

Compute  
**covariance**  
matrix

$$S = P D P^T$$

2.5	1.5
1.5	1.5

-0.81	0.58
-0.58	-0.81

3.58	0
0	0.41

-0.81	-0.58
0.58	-0.81

Perform  
**eigendecomposition**

$$P^T \times Z^T =$$

-0.81	-0.58
0.58	-0.81

-2	-1	0	1	2
-2	0	1	0	1

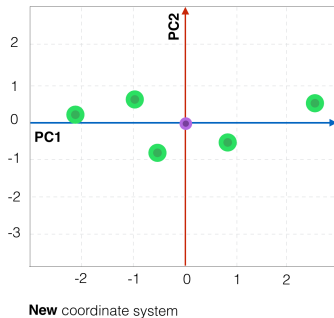
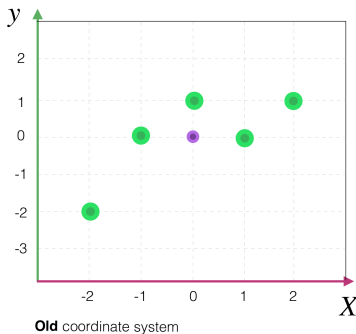
2.78	0.46
0.81	-0.58
-0.58	-0.81
-0.81	0.58
-2.2	0.35

Compute new  
**coordinates**



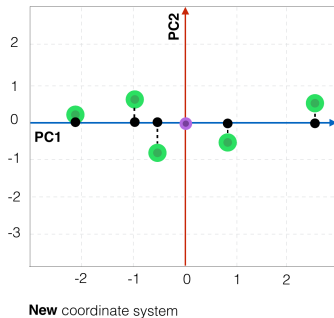
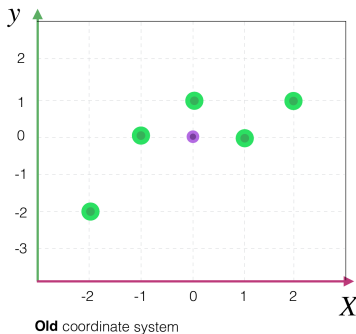
# Principal Component Analysis

Do you still remember what was it all about?



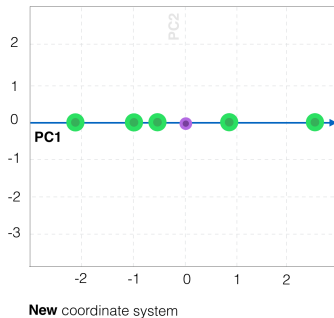
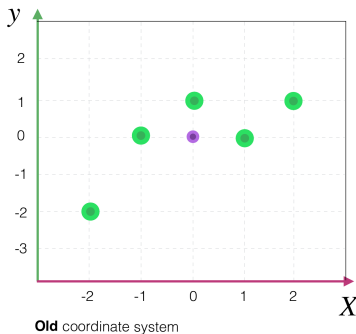
# Principal Component Analysis

Do you still remember what was it all about?  
We want to **reduce the dimensionality**!



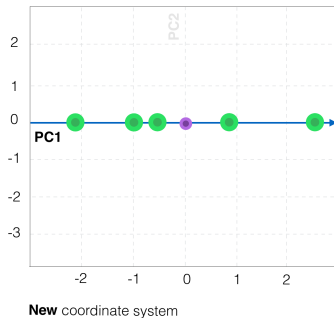
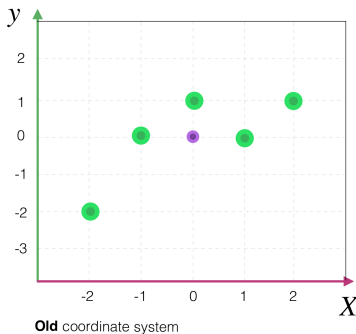
# Principal Component Analysis

Do you still remember what was it all about?  
We want to **reduce the dimensionality**!



# Principal Component Analysis

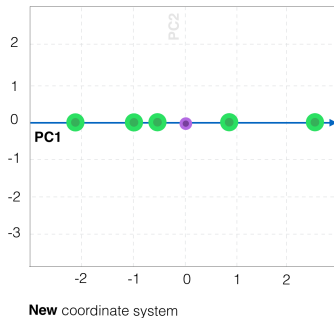
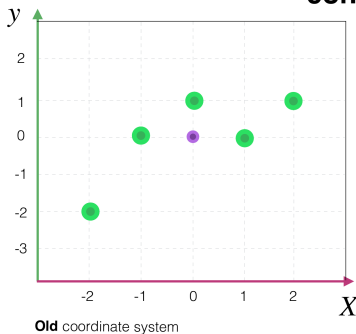
We can **ignore** the **second eigenvector** because it does not contain much information



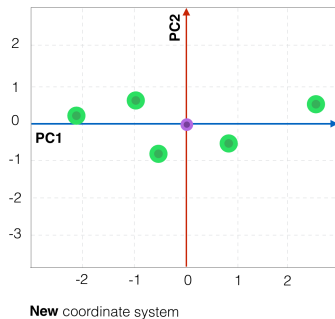
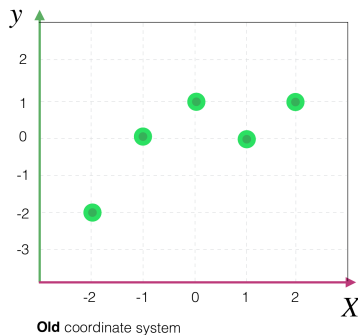
# Principal Component Analysis

We can **ignore** the **second eigenvector**  
because it does not contain much information

**How much** information **the second eigenvector**  
**contains?**



# Principal Component Analysis

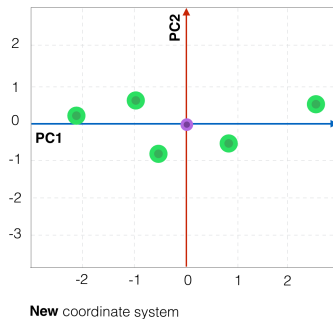
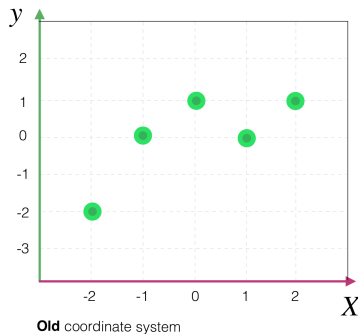


# Principal Component Analysis

$$D$$

3.58	0
0	0.42

Eigenvalues



# Principal Component Analysis

$$S$$

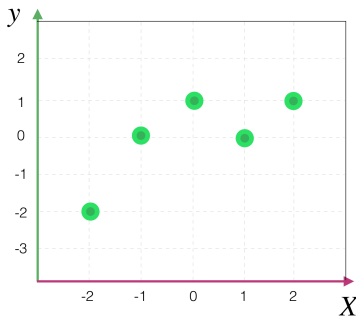
2.5	1.5
1.5	1.5

**Covariance** matrix

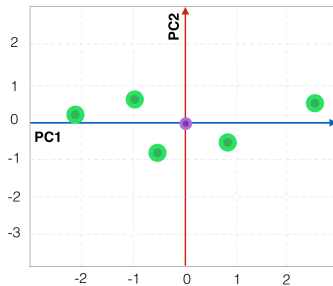
$$D$$

3.58	0
0	0.42

**Eigenvalues**



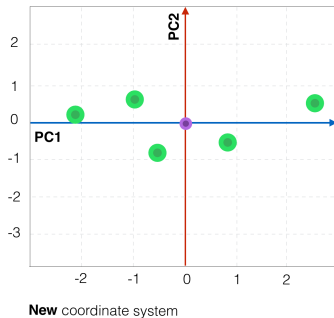
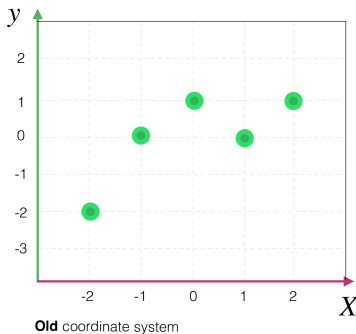
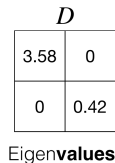
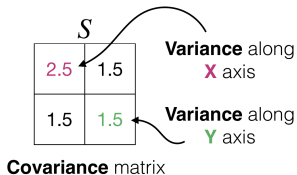
**Old** coordinate system



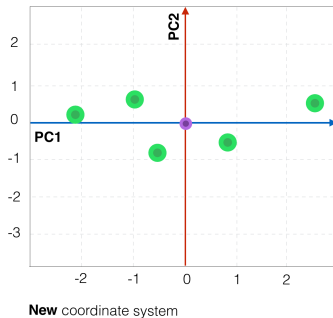
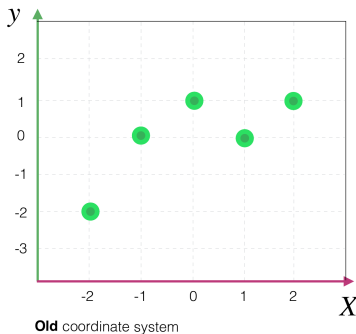
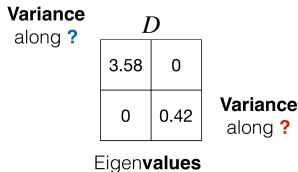
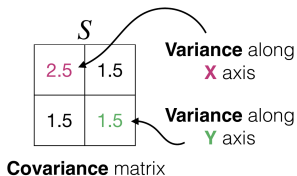
**New** coordinate system



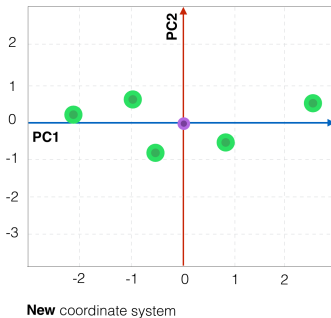
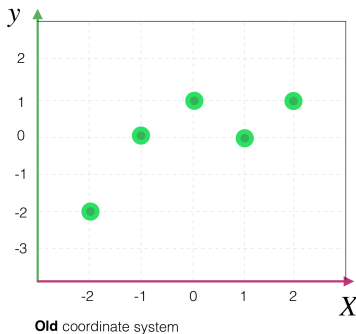
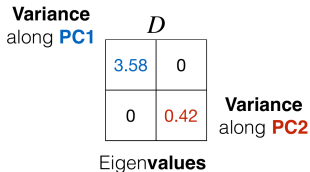
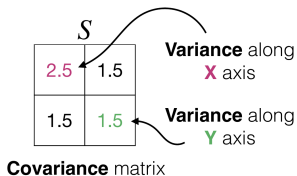
# Principal Component Analysis



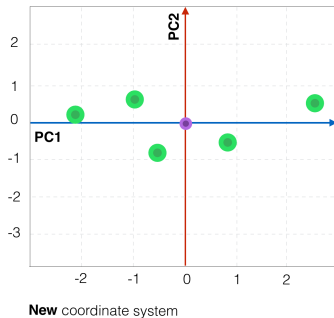
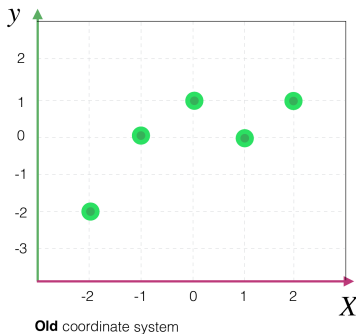
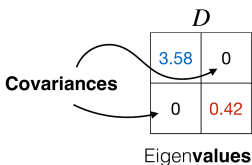
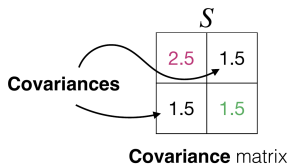
# Principal Component Analysis



# Principal Component Analysis



# Principal Component Analysis



# Principal Component Analysis

$$S$$

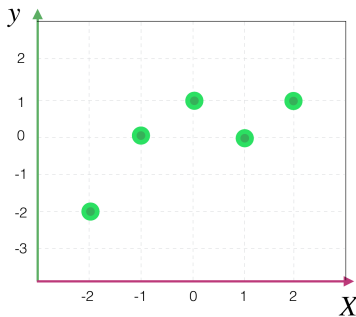
2.5	1.5
1.5	1.5

Covariance matrix

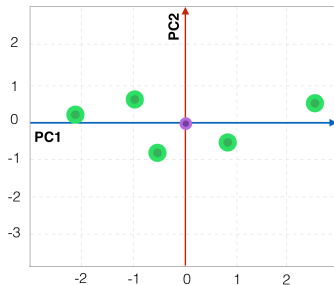
$$D$$

3.58	0
0	0.42

New covariance matrix



Old coordinate system



New coordinate system

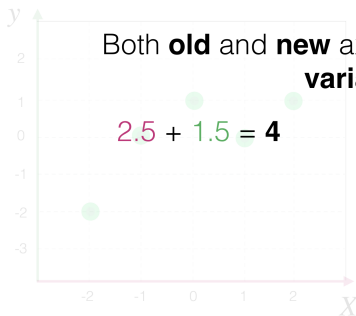
# Principal Component Analysis

$S$	
2.5	1.5
1.5	1.5

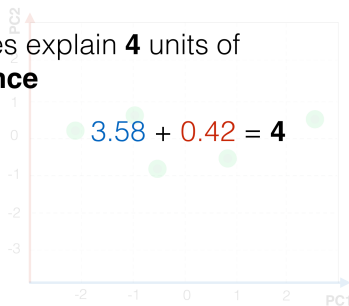
**Covariance** matrix

$D$	
3.58	0
0	0.42

New **covariance** matrix



**Old** coordinate system



**New** coordinate system

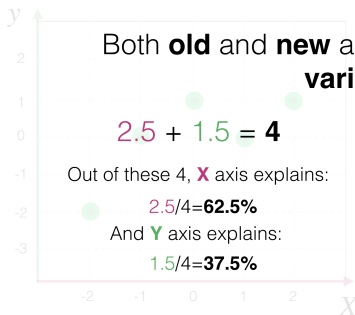
# Principal Component Analysis

$S$	
2.5	1.5
1.5	1.5

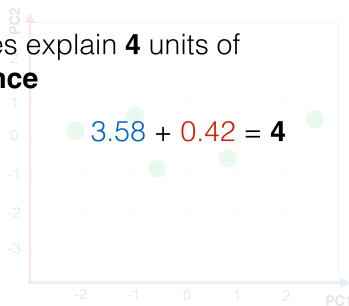
**Covariance** matrix

$D$	
3.58	0
0	0.42

New **covariance** matrix



**Old** coordinate system



**New** coordinate system

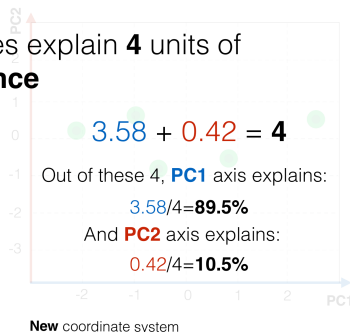
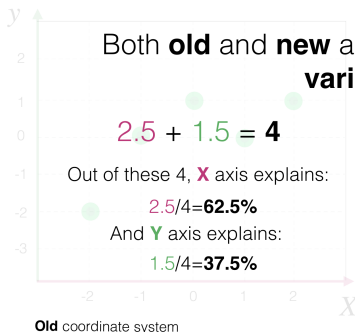
# Principal Component Analysis

$S$	
2.5	1.5
1.5	1.5

**Covariance** matrix

$D$	
3.58	0
0	0.42

New **covariance** matrix

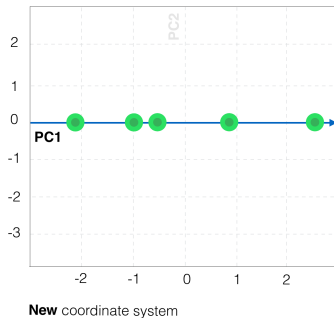
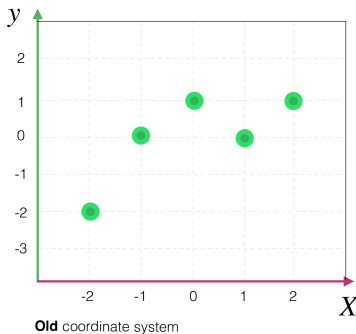




# Principal Component Analysis

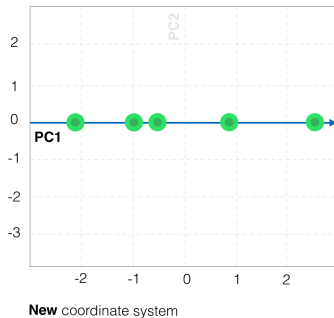
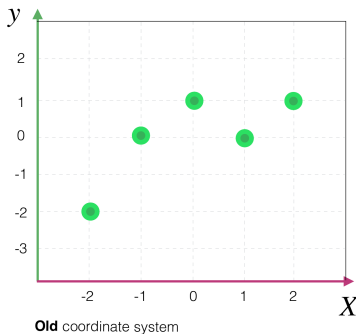
We can **ignore** the **second PC** because it does not contain much information

**How much** information **the second PC** contains?



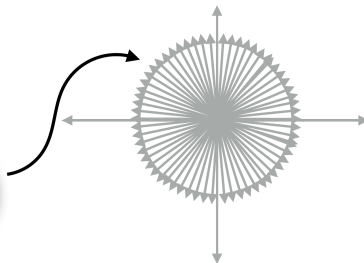
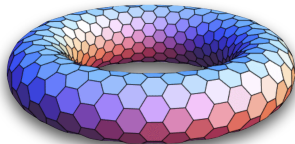
# Principal Component Analysis

We can **ignore** the **second PC** because it explains only **10.5%** of variation



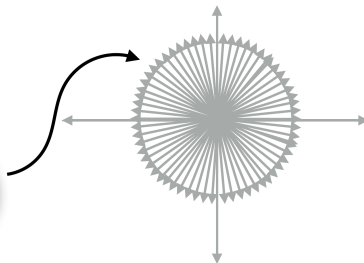
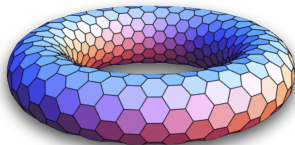
# Principal Component Analysis

**How many** PCs will be formed in **200D** space?



# Principal Component Analysis

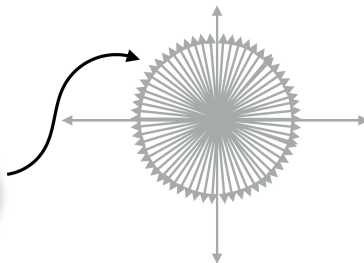
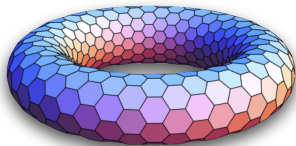
**How many** PCs will be formed in **200D** space?



No exceptions, **200 PCs**

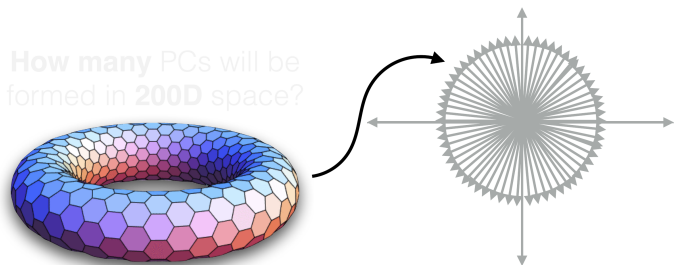
# Principal Component Analysis

How many PCs will be formed in **200D** space?



No exceptions, **200 PCs**

# Principal Component Analysis



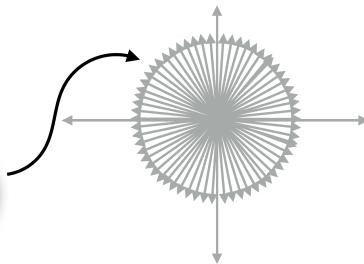
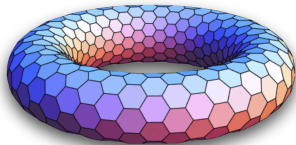
**How many** PCs should we keep?

No exceptions, 200 PCs

# Principal Component Analysis

**Variance explained** is a good criteria for choosing the total number of PCs to keep

How many PCs will be formed in 200D space?



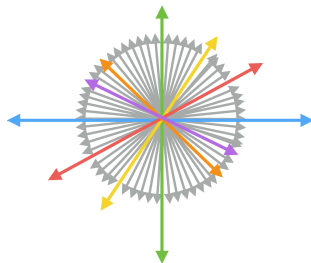
No exceptions, 200 PCs

**How many** PCs should we keep?

# Principal Component Analysis

**Variance explained** is a good criteria for choosing the total number of PCs to keep

You should keep as many PCs as it takes to explain **90% of total variance**



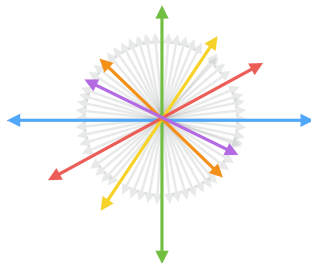
**How many** PCs should we keep?



# Principal Component Analysis

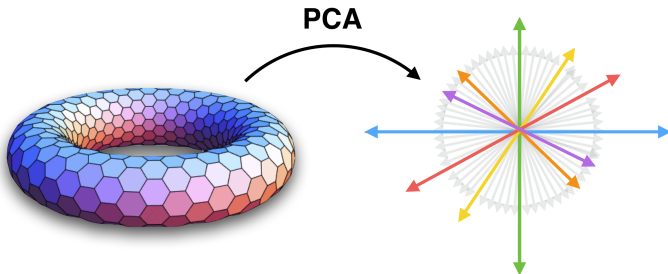
**Variance explained** is a good criteria for choosing the total number of PCs to keep

You should keep as many PCs as it takes to explain **90% of total variance**



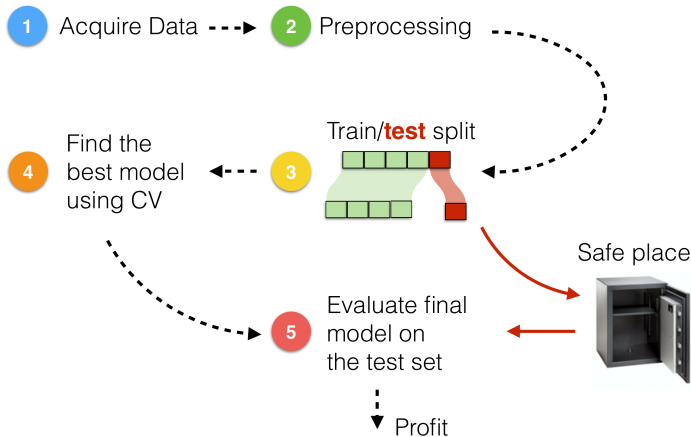
**How many** PCs should we keep?

## Principle Component Analysis (**PCA**)

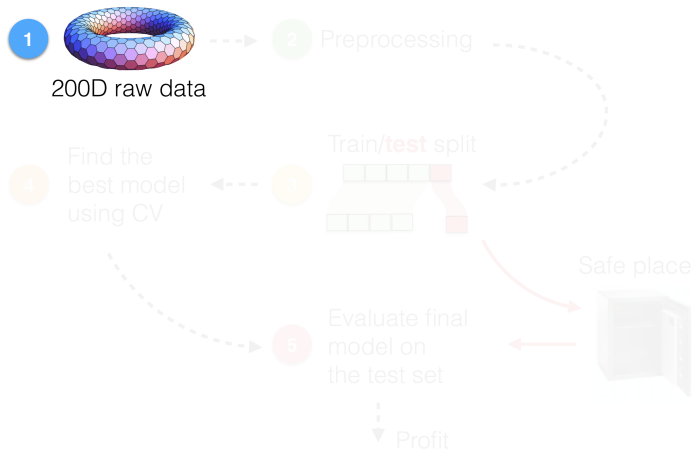


Can be used as part of supervised learning pipeline

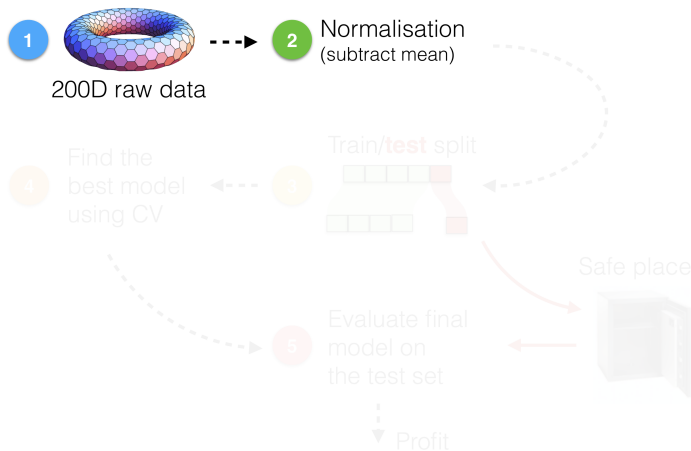
## Supervised Learning pipeline



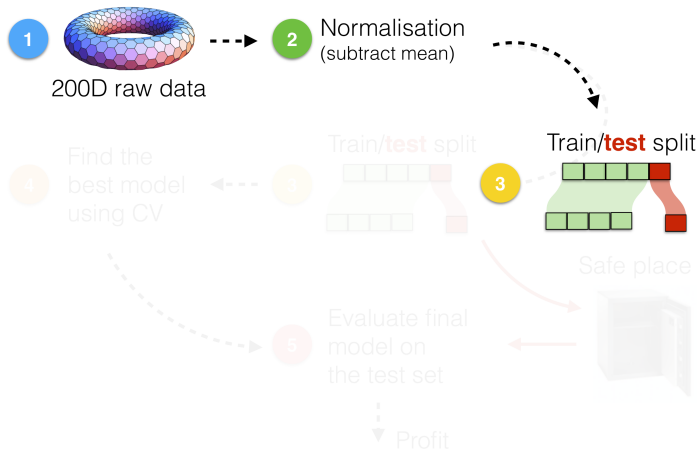
## Supervised Learning pipeline



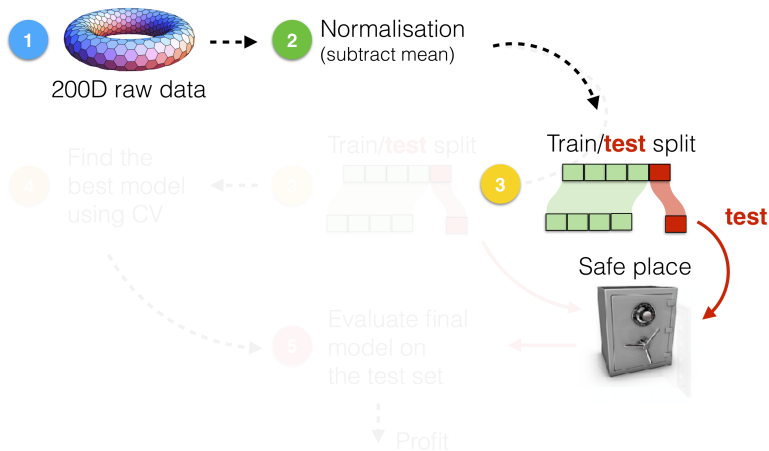
## Supervised Learning pipeline



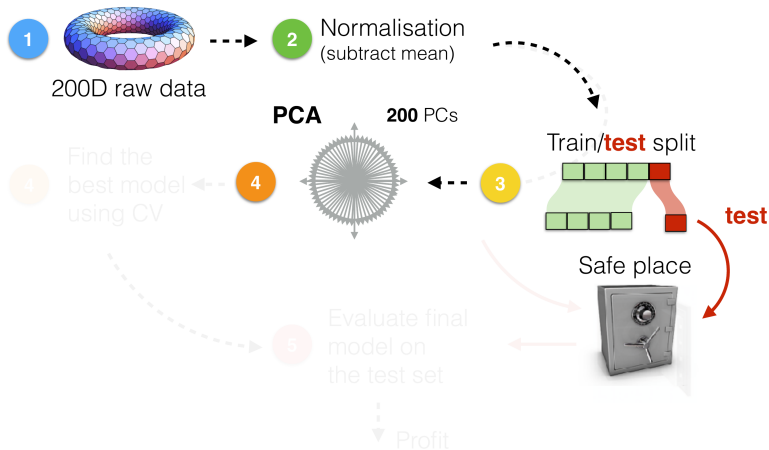
## Supervised Learning pipeline



## Supervised Learning pipeline

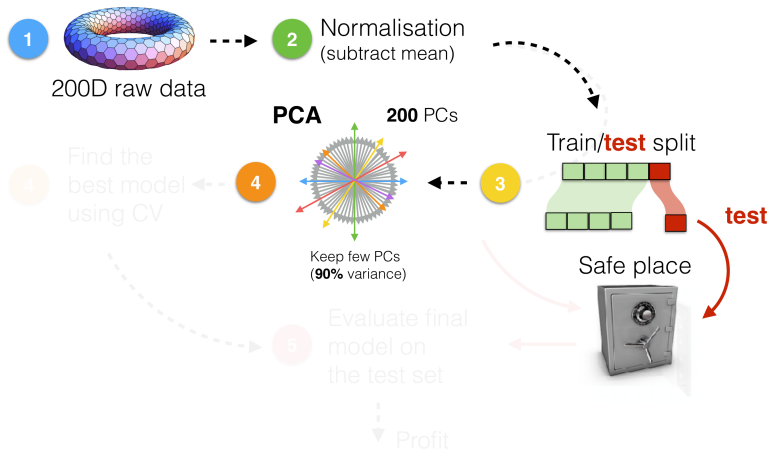


## Supervised Learning pipeline

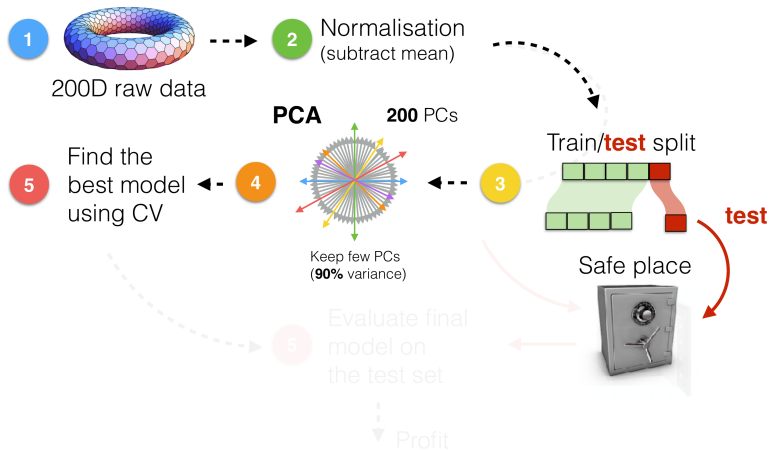




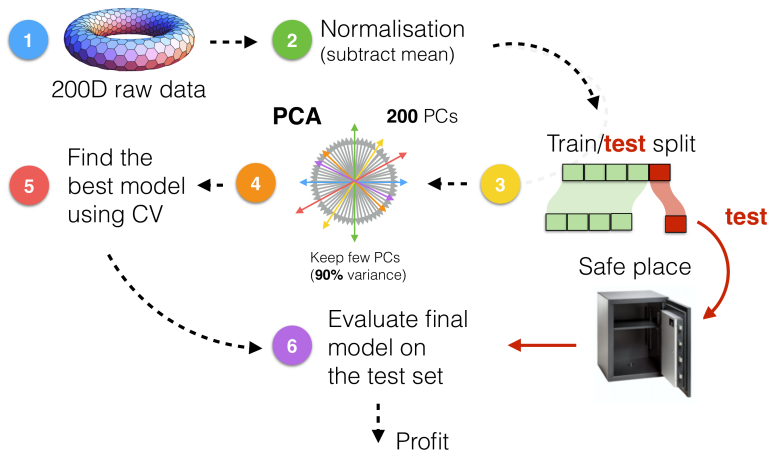
## Supervised Learning pipeline



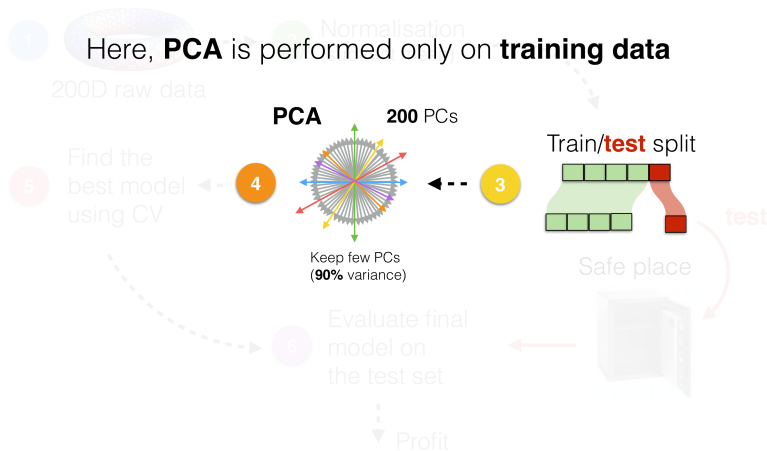
## Supervised Learning pipeline



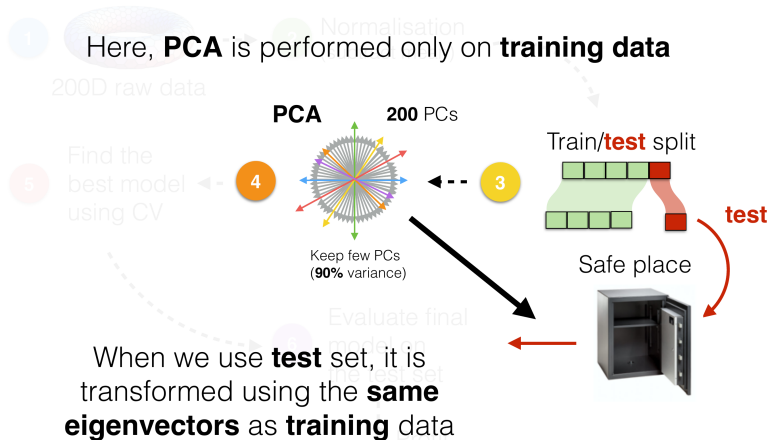
## Supervised Learning pipeline



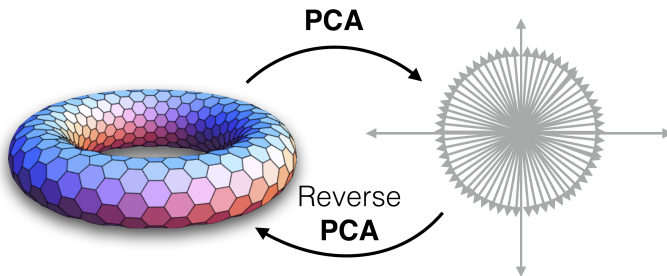
## Supervised Learning pipeline



## Supervised Learning pipeline



**PCA** has an “**undo**” button



You can **recover** the **original features** back!

**PCA** has an “**undo**” button

**Conventional** PCA

$$\text{eigenvectors} \times Z_{\text{original}}^{\top} = Z_{\text{transformed}}$$

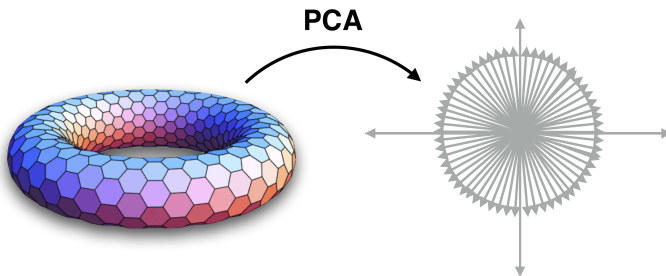
**Reversed** PCA

$$\text{eigenvectors}^{\top} \times Z_{\text{transformed}} = Z_{\text{original}}$$

You can **recover** the **original features** back!

# Principal Component Analysis

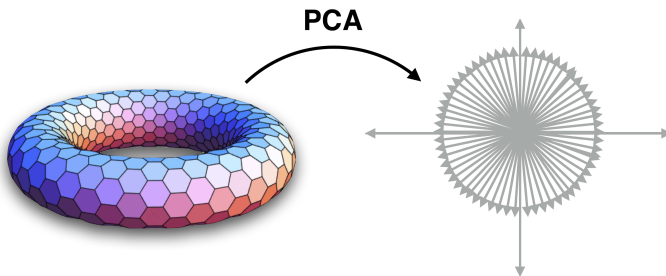
Principle components **are linear combinations of original features**





# Principal Component Analysis

Principle components **are linear combinations** of **original features**



So if you predict anything based on PCs, the **meaning** of original features **is not preserved** after the transformation

# Problem Setting of PCA

- In PCA, we are interested in finding projections  $\mathbf{x}'_n$  of data points  $\mathbf{x}_n$  that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.

# Problem Setting of PCA

- In PCA, we are interested in finding projections  $\mathbf{x}'_n$  of data points  $\mathbf{x}_n$  that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times D}$ , where  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $\mathbb{E}(\mathbf{X}) = 0$  and the sample covariance of  $\mathbf{X}$  is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

# Problem Setting of PCA

- In PCA, we are interested in finding projections  $\mathbf{x}'_n$  of data points  $\mathbf{x}_n$  that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times D}$ , where  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $\mathbb{E}(\mathbf{X}) = 0$  and the sample covariance of  $\mathbf{X}$  is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

- We assume there exists a low-dimensional compressed representation (code) of  $\mathbf{x}_n$

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$  is the projection matrix.

# Problem Setting of PCA

- In PCA, we are interested in finding projections  $\mathbf{x}'_n$  of data points  $\mathbf{x}_n$  that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times D}$ , where  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $\mathbb{E}(\mathbf{X}) = 0$  and the sample covariance of  $\mathbf{X}$  is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

- We assume there exists a low-dimensional compressed representation (code) of  $\mathbf{x}_n$

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$  is the projection matrix.

- We further assume that the columns of  $\mathbf{B}$  are orthonormal, so that  $\mathbf{b}_i^T \mathbf{b}_j = 0$  iff  $i \neq j$  and  $\mathbf{b}_i^T \mathbf{b}_i = 1$

# Problem Setting of PCA

- In PCA, we are interested in finding projections  $\mathbf{x}'_n$  of data points  $\mathbf{x}_n$  that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times D}$ , where  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $\mathbb{E}(\mathbf{X}) = 0$  and the sample covariance of  $\mathbf{X}$  is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

- We assume there exists a low-dimensional compressed representation (code) of  $\mathbf{x}_n$

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$  is the projection matrix.

- We further assume that the columns of  $\mathbf{B}$  are orthonormal, so that  $\mathbf{b}_i^T \mathbf{b}_j = 0$  iff  $i \neq j$  and  $\mathbf{b}_i^T \mathbf{b}_i = 1$
- We seek an  $M$ -dimensional subspace  $U \subseteq \mathbb{R}^D$ ,  $\dim(U) = M < D$  onto which we project the data

# Problem Setting of PCA

- In PCA, we are interested in finding projections  $\mathbf{x}'_n$  of data points  $\mathbf{x}_n$  that are as similar to the original data points as possible, but which have a significantly lower intrinsic dimensionality.
- Suppose we have a dataset  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times D}$ , where  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $\mathbb{E}(\mathbf{X}) = 0$  and the sample covariance of  $\mathbf{X}$  is

$$S = \frac{1}{N} \mathbf{X}^T \mathbf{X} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T$$

- We assume there exists a low-dimensional compressed representation (code) of  $\mathbf{x}_n$

$$\mathbf{z}_n = \mathbf{B}^T \mathbf{x}_n \in \mathbb{R}^M,$$

where  $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_M] \in \mathbb{R}^{D \times M}$  is the projection matrix.

- We further assume that the columns of  $\mathbf{B}$  are orthonormal, so that  $\mathbf{b}_i^T \mathbf{b}_j = 0$  iff  $i \neq j$  and  $\mathbf{b}_i^T \mathbf{b}_i = 1$
- We seek an  $M$ -dimensional subspace  $U \subseteq \mathbb{R}^D$ ,  $\dim(U) = M < D$  onto which we project the data
- We denote the projected data as  $\mathbf{x}'_n \in U$  and their coordinates by  $\mathbf{z}_n$

# Problem Setting of PCA

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**



# Problem Setting of PCA

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation

# Problem Setting of PCA

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.

# Problem Setting of PCA

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.
- We start by seeking a single vector  $\mathbf{b}_1 \in \mathbb{R}^D$  that maximizes variance of the projected data

$$\text{Var}(z_1) = \frac{1}{N} \sum_{n=1}^N z_{1n}^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^T \mathbf{x}_n)^2 =$$

# Problem Setting of PCA

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.
- We start by seeking a single vector  $\mathbf{b}_1 \in \mathbb{R}^D$  that maximizes variance of the projected data

$$\begin{aligned}\mathbb{V}ar(z_1) &= \frac{1}{N} \sum_{n=1}^N z_{1n}^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^T \mathbf{x}_n)^2 = \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{b}_1 = \mathbf{b}_1^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1 = \mathbf{b}_1^T S \mathbf{b}_1,\end{aligned}$$

where  $S$  is the sample covariance matrix.

# Problem Setting of PCA

- We can describe the information contained in the data by looking at the spread of the data and measure it with **variance**
- In this setting, retaining most information after data compression is equivalent to capturing the largest amount of variance in the low-dimensional representation
- We maximize the variance of the low-dimensional code using a sequential approach.
- We start by seeking a single vector  $\mathbf{b}_1 \in \mathbb{R}^D$  that maximizes variance of the projected data

$$\begin{aligned}\mathbb{V}ar(z_1) &= \frac{1}{N} \sum_{n=1}^N z_{1n}^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_1^T \mathbf{x}_n)^2 = \\ &= \frac{1}{N} \sum_{n=1}^N \mathbf{b}_1^T \mathbf{x}_n \mathbf{x}_n^T \mathbf{b}_1 = \mathbf{b}_1^T \left( \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^T \right) \mathbf{b}_1 = \mathbf{b}_1^T S \mathbf{b}_1,\end{aligned}$$

where  $S$  is the sample covariance matrix.

- Note that if we would not put a unit vector constraint on  $\mathbf{b}_1$ , the variance would increase for longer-length  $\mathbf{b}_1$ .

We need to solve the following constrained optimization problem

$$\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$$

$$\text{subject to } \|\mathbf{b}_1\|^2 = 1$$

We need to solve the following constrained optimization problem

$$\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$$

$$\text{subject to } \|\mathbf{b}_1\|^2 = 1$$

The Lagrangian will be

$$\Lambda(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^T S \mathbf{b}_1 + \lambda_1(1 - \mathbf{b}_1^T \mathbf{b}_1)$$

# Mathematical Derivations

We need to solve the following constrained optimization problem

$$\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$$

$$\text{subject to } \|\mathbf{b}_1\|^2 = 1$$

The Lagrangian will be

$$\Lambda(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^T S \mathbf{b}_1 + \lambda_1(1 - \mathbf{b}_1^T \mathbf{b}_1)$$

Finding the partial derivatives w.r.t.  $\mathbf{b}_1$  and  $\lambda_1$

$$\frac{\partial \Lambda}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^T S - 2\lambda_1 \mathbf{b}_1^T, \quad \frac{\partial \Lambda}{\partial \lambda_1} = 1 - \mathbf{b}_1^T \mathbf{b}_1$$



# Mathematical Derivations

We need to solve the following constrained optimization problem

$$\max_{\mathbf{b}_1} \mathbf{b}_1^T S \mathbf{b}_1$$

$$\text{subject to } \|\mathbf{b}_1\|^2 = 1$$

The Lagrangian will be

$$\Lambda(\mathbf{b}_1, \lambda_1) = \mathbf{b}_1^T S \mathbf{b}_1 + \lambda_1(1 - \mathbf{b}_1^T \mathbf{b}_1)$$

Finding the partial derivatives w.r.t.  $\mathbf{b}_1$  and  $\lambda_1$

$$\frac{\partial \Lambda}{\partial \mathbf{b}_1} = 2\mathbf{b}_1^T S - 2\lambda_1 \mathbf{b}_1^T, \quad \frac{\partial \Lambda}{\partial \lambda_1} = 1 - \mathbf{b}_1^T \mathbf{b}_1$$

Setting these partial derivatives to 0 gives us

$$S \mathbf{b}_1 = \lambda_1 \mathbf{b}_1$$

$$\mathbf{b}_1^T \mathbf{b}_1 = 1$$

Does this look familiar?

$$S\mathbf{b}_1 = \lambda_1\mathbf{b}_1$$
$$\mathbf{b}_1^T\mathbf{b}_1 = 1$$

$\mathbf{b}_1$  is an eigenvector of the data covariance matrix  $S$ , and the Lagrange multiplier  $\lambda_1$  plays the role of the corresponding eigenvalue.

$$S\mathbf{b}_1 = \lambda_1\mathbf{b}_1$$
$$\mathbf{b}_1^T\mathbf{b}_1 = 1$$

$\mathbf{b}_1$  is an eigenvector of the data covariance matrix  $S$ , and the Lagrange multiplier  $\lambda_1$  plays the role of the corresponding eigenvalue.

The objective can be re-written as

$$\mathbb{V}ar(z_1) = \mathbf{b}_1^T S \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^T \mathbf{b}_1 = \lambda_1,$$

the variance of the data projected onto a one-dimensional subspace equals the eigenvalue that is associated with the basis vector  $\mathbf{b}_1$  that spans this subspace.

$$\begin{aligned}S\mathbf{b}_1 &= \lambda_1\mathbf{b}_1 \\ \mathbf{b}_1^T\mathbf{b}_1 &= 1\end{aligned}$$

$\mathbf{b}_1$  is an eigenvector of the data covariance matrix  $S$ , and the Lagrange multiplier  $\lambda_1$  plays the role of the corresponding eigenvalue.

The objective can be re-written as

$$\mathbb{V}ar(z_1) = \mathbf{b}_1^T S \mathbf{b}_1 = \lambda_1 \mathbf{b}_1^T \mathbf{b}_1 = \lambda_1,$$

the variance of the data projected onto a one-dimensional subspace equals the eigenvalue that is associated with the basis vector  $\mathbf{b}_1$  that spans this subspace. To obtain the projection of  $\mathbf{x}_n$  on the obtained subspace, we can use the following formula

$$\mathbf{x}'_n = \mathbf{b}_1 z_{1n} = \mathbf{b}_1 \mathbf{b}_1^T \mathbf{x}_n \in \mathbb{R}^D$$

# Mathematical Derivations

- Assume we have found the first  $m - 1$  principal components as the  $m - 1$  eigenvectors of  $S$  that are associated with the largest  $m - 1$  eigenvalues.

# Mathematical Derivations

- Assume we have found the first  $m - 1$  principal components as the  $m - 1$  eigenvectors of  $S$  that are associated with the largest  $m - 1$  eigenvalues.
- The  $m$ -th principal component can be found by subtracting the effect of the first  $m - 1$  principal components  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  from the data, by trying to find principal components that compress the remaining information.

$$\hat{\mathbf{X}} = \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X} = \mathbf{X} - \mathbf{B}_{m-1} \mathbf{X},$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  and  $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$  is the projection matrix onto the subspace spanned by  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$ .

# Mathematical Derivations

- Assume we have found the first  $m - 1$  principal components as the  $m - 1$  eigenvectors of  $S$  that are associated with the largest  $m - 1$  eigenvalues.
- The  $m$ -th principal component can be found by subtracting the effect of the first  $m - 1$  principal components  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  from the data, by trying to find principal components that compress the remaining information.

$$\hat{\mathbf{X}} = \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X} = \mathbf{X} - \mathbf{B}_{m-1} \mathbf{X},$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  and  $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$  is the projection matrix onto the subspace spanned by  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$ .

- The  $m$ -th PC can be found by maximizing the variance

$$\text{Var}(z_m) = \frac{1}{N} \sum_{n=1}^N z_{mn}^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_m^T \hat{\mathbf{x}}_n)^2 = \mathbf{b}_m^T \hat{S} \mathbf{b}_m$$

$$\text{subject to } \|\mathbf{b}_m\|^2 = 1$$

where  $\hat{S}$  is the covariance matrix of  $\hat{\mathbf{X}}$

# Mathematical Derivations

- Assume we have found the first  $m - 1$  principal components as the  $m - 1$  eigenvectors of  $S$  that are associated with the largest  $m - 1$  eigenvalues.
- The  $m$ -th principal component can be found by subtracting the effect of the first  $m - 1$  principal components  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$  from the data, by trying to find principal components that compress the remaining information.

$$\hat{\mathbf{X}} = \mathbf{X} - \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T \mathbf{X} = \mathbf{X} - \mathbf{B}_{m-1} \mathbf{X},$$

where  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_N] \in \mathbb{R}^{D \times N}$  and  $\mathbf{B}_{m-1} = \sum_{i=1}^{m-1} \mathbf{b}_i \mathbf{b}_i^T$  is the projection matrix onto the subspace spanned by  $\mathbf{b}_1, \dots, \mathbf{b}_{m-1}$ .

- The  $m$ -th PC can be found by maximizing the variance

$$\text{Var}(z_m) = \frac{1}{N} \sum_{n=1}^N z_{mn}^2 = \frac{1}{N} \sum_{n=1}^N (\mathbf{b}_m^T \hat{\mathbf{x}}_n)^2 = \mathbf{b}_m^T \hat{S} \mathbf{b}_m$$

$$\text{subject to } \|\mathbf{b}_m\|^2 = 1$$

where  $\hat{S}$  is the covariance matrix of  $\hat{\mathbf{X}}$

- $\mathbf{b}_m$  is also an eigenvector of  $S$  and  $\text{Var}(z_m) = \mathbf{b}_m^T S \mathbf{b}_m = \lambda_m$  is the  $m$ -th largest eigenvalue of  $S$



# Summary of PCA

- We can reduce the dimensionality of our data, by using the  $M$  eigenvectors of the covariance matrix  $S$  associated with the  $M$  largest eigenvalues.

# Summary of PCA

- We can reduce the dimensionality of our data, by using the  $M$  eigenvectors of the covariance matrix  $S$  associated with the  $M$  largest eigenvalues.
- The projection matrix  $\mathbf{B}$  consists of the eigenvectors of  $S$
- The amount of variance PCA captured with the first  $M$  principal components is

$$\sum_{i=1}^M \lambda_m$$

# Summary of PCA

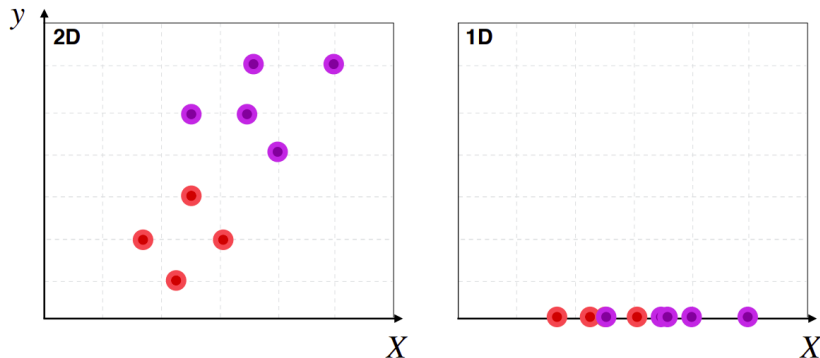
- We can reduce the dimensionality of our data, by using the  $M$  eigenvectors of the covariance matrix  $S$  associated with the  $M$  largest eigenvalues.
- The projection matrix  $\mathbf{B}$  consists of the eigenvectors of  $S$
- The amount of variance PCA captured with the first  $M$  principal components is

$$\sum_{i=1}^M \lambda_m$$

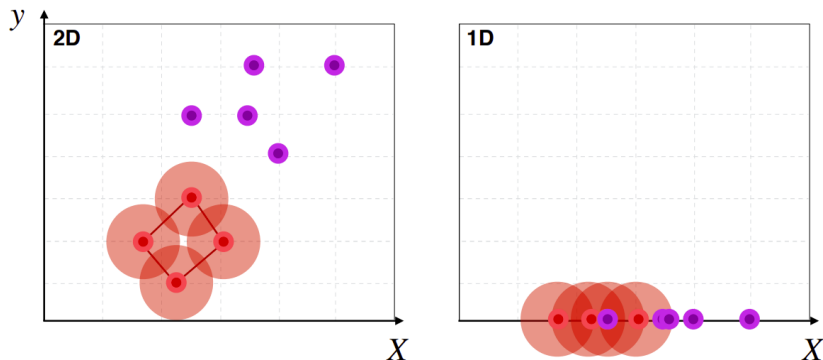
- The variance lost by data compression via PCA is

$$\sum_{i=M+1}^D \lambda_m$$

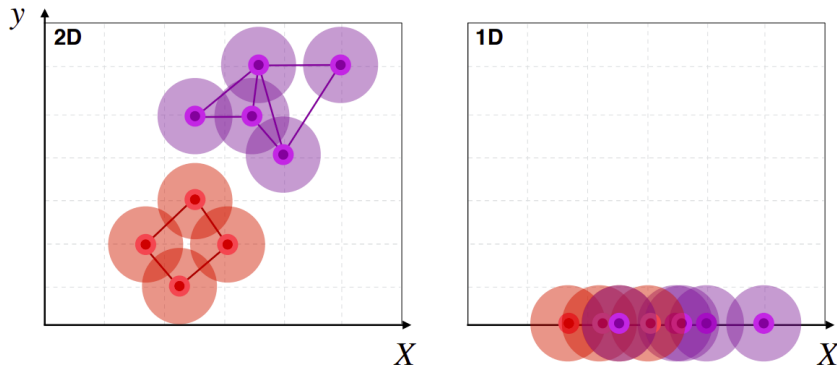
# Intuition behind t-SNE



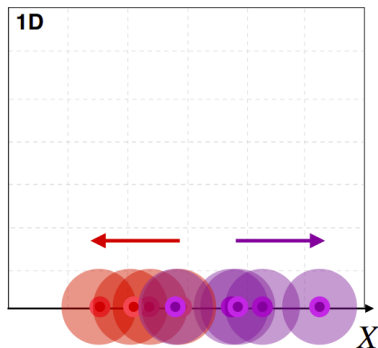
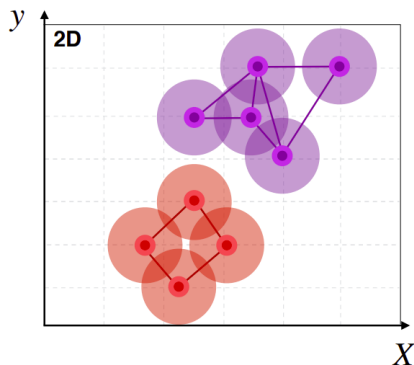
# Intuition behind t-SNE



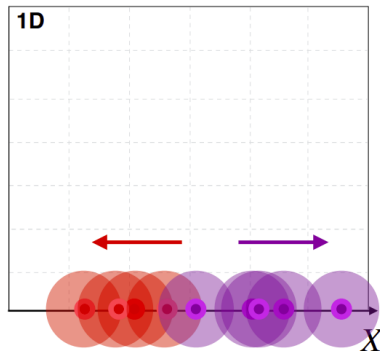
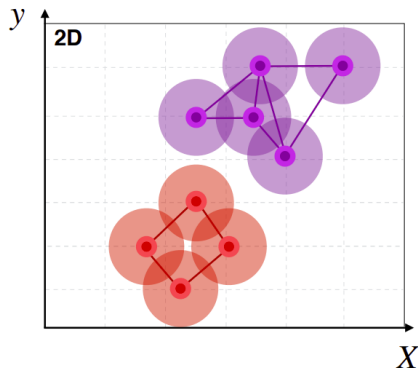
# Intuition behind t-SNE



# Intuition behind t-SNE

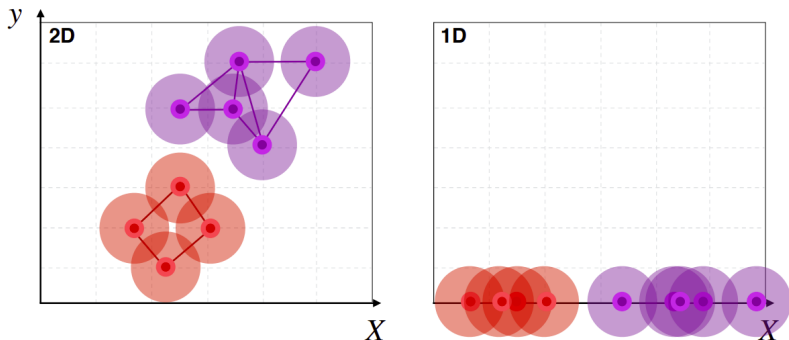


# Intuition behind t-SNE



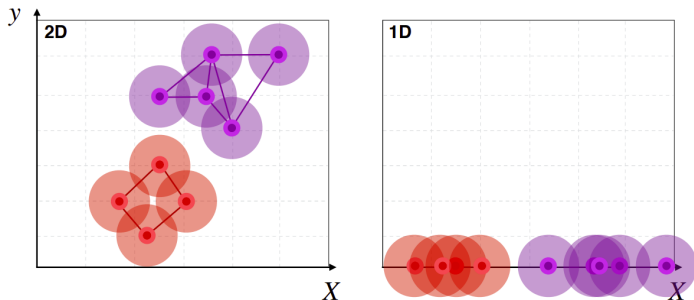


# Intuition behind t-SNE



# Intuition behind t-SNE

**t-SNE** iteratively tries to make distances in low-dimensional space to be similar to distances in high-dimensional space



# What have we learned today?

- ✓ Dimensionality Reduction
- ✓ Principal Component Analysis
- ✓ t-Distributed Stochastic Neighbor Embedding (t-SNE)