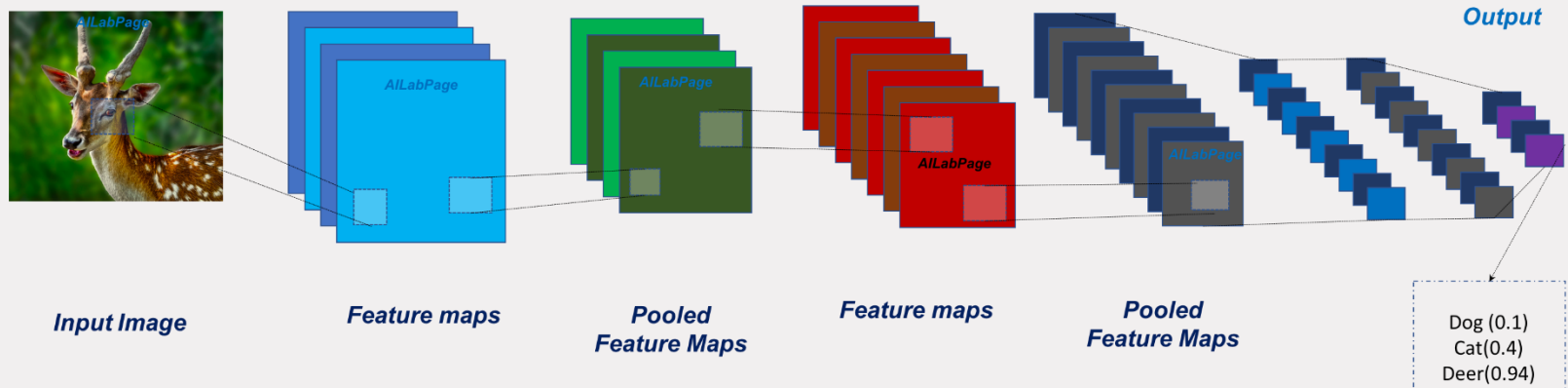


# Unsupervised machine learning (8DC00)

Cian Scannell

# Previous lecture

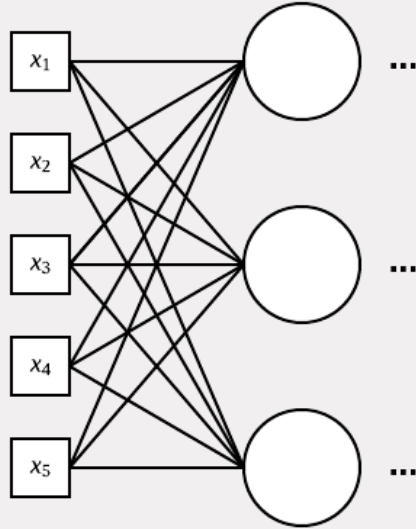
- Convolutional neural networks (CNN)



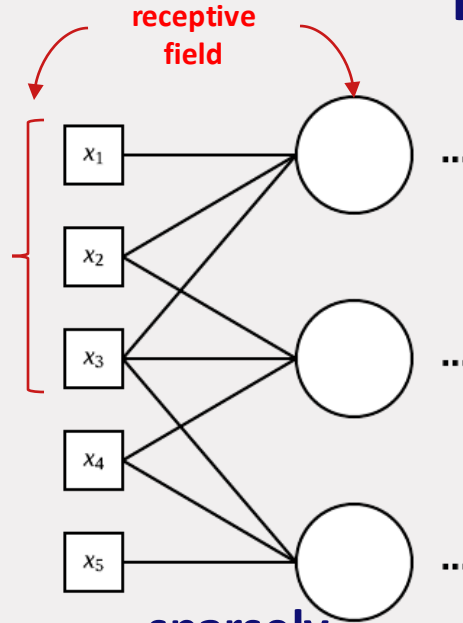
- By now: students can design a (simple) CNN

Image from: <https://vinodsblog.com/>

## Previous lecture

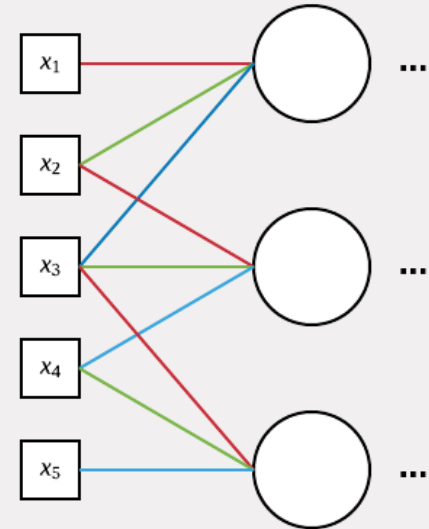


**“regular” NN**  
15 weights



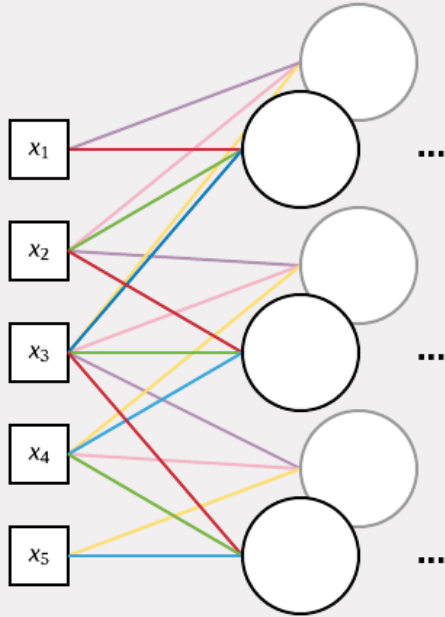
**sparsely  
connected NN**  
9 weights

## Reducing # of weights



**shared weights**  
3 weights

## Previous lecture



two sets shared weights  
6 weights

$$\begin{bmatrix} a_{1,1} & a_{1,2} & a_{1,3} \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 \end{bmatrix} * \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \end{bmatrix}$$

$$\begin{bmatrix} a_{2,1} & a_{2,2} & a_{2,3} \end{bmatrix} = \begin{bmatrix} x_1 & x_2 & x_3 & x_4 & x_5 \end{bmatrix} * \begin{bmatrix} w_{2,1} & w_{2,2} & w_{2,3} \end{bmatrix}$$

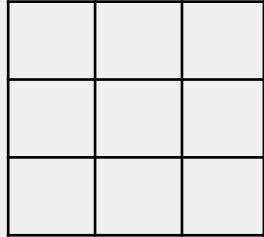
$\begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} \end{bmatrix}$ , and  $\begin{bmatrix} w_{2,1} & w_{2,2} & w_{2,3} \end{bmatrix}$   
are **convolution kernels**. They extract  
features.

# Previous lecture: Kernel size

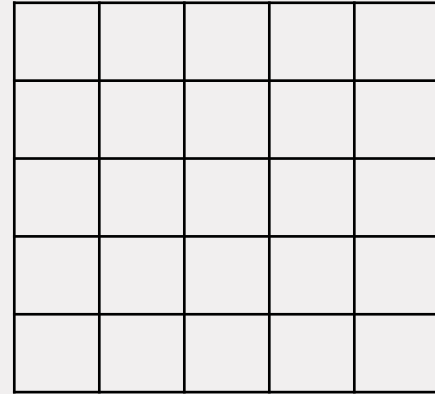
1 x 1



3 x 3

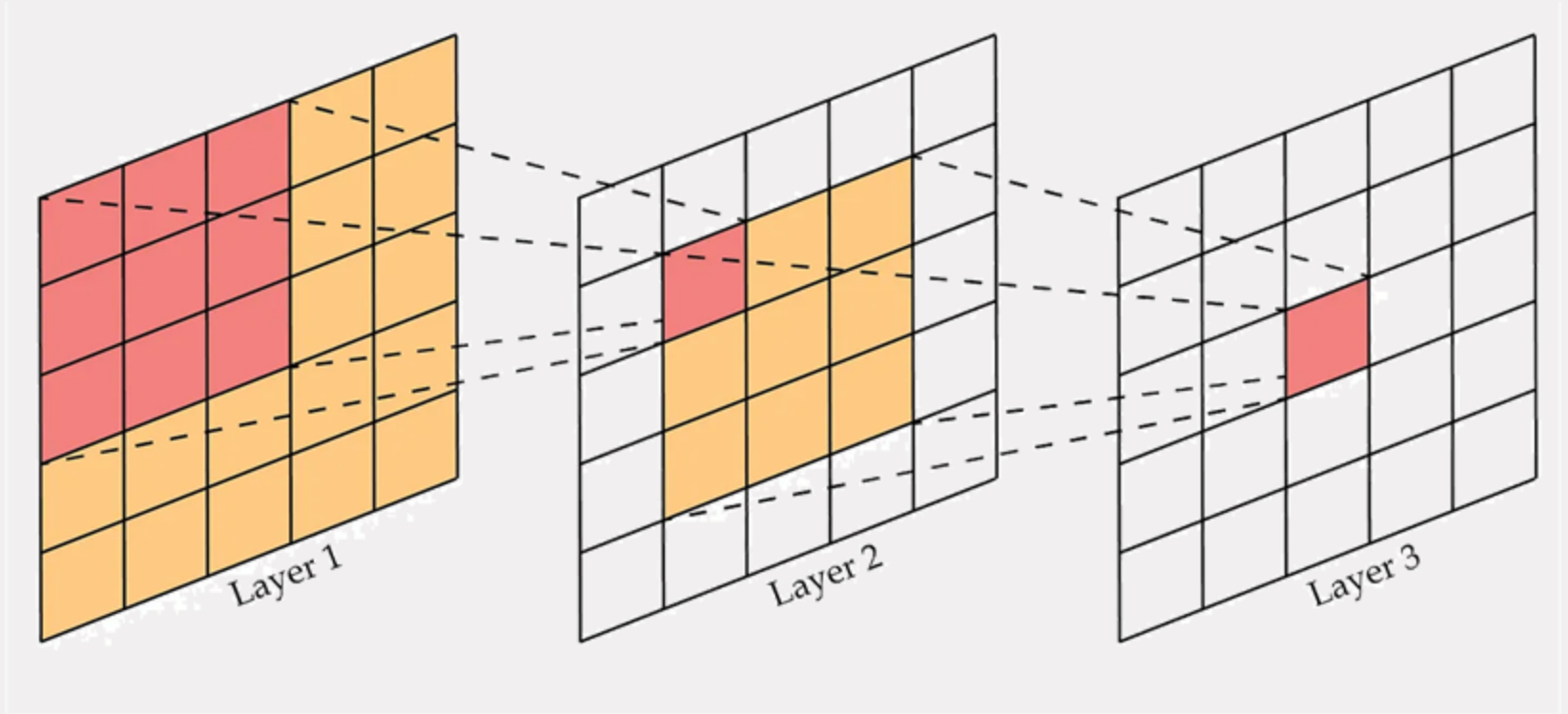


5 x 5



- More weights → more information
- More computations / memory
- Receptive field

# Receptive field



# Previous lecture: Max-pooling

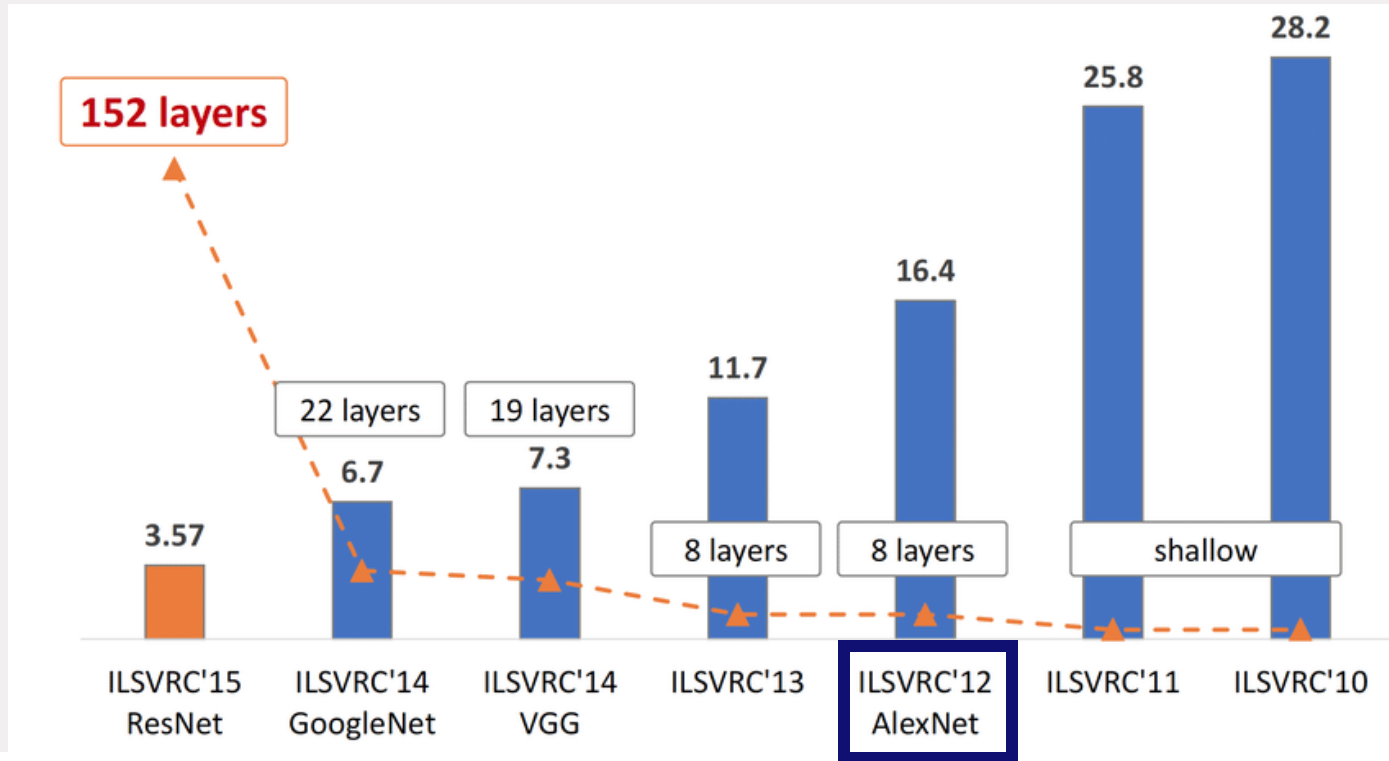
- Reduce size of feature space
- Maximum of features
- Typical kernel size =  $2 \times 2$

*Addition to last lecture:*

Stride = step size of the sliding window

- Typically 1 for convolutions
- Typically equal to kernel size for pooling operations

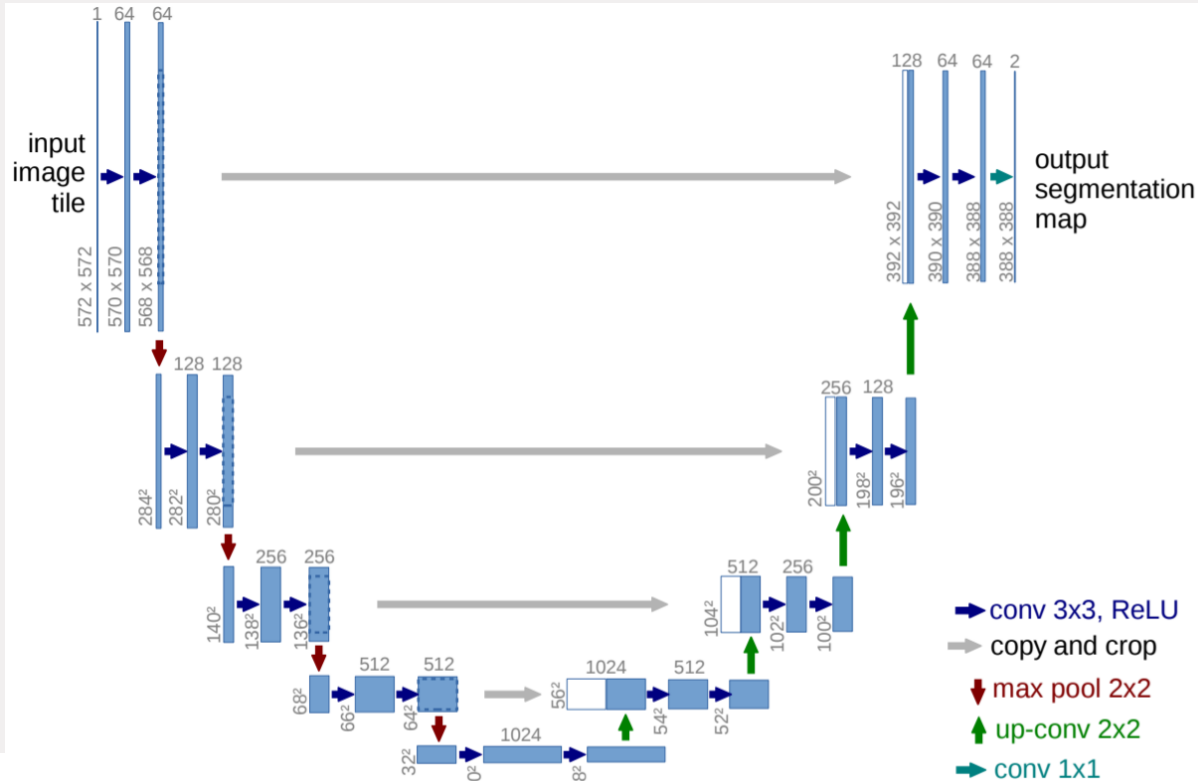
# Error rates ImageNet (top 5 accuracy)





# U-Net: Segmentation

Ronneberger et al., MICCAI 2015



# Solving machine learning problems

- Type of problem?
- Preprocessing steps?
- What kind of 'labels' do you need?
- What should the output look like?
- What should be the activation function in the final layer?
- What model do you choose?
- How do you evaluate performance?
- Other considerations?

# Previous examples required labeled data

learning from labeled data = **supervised training**

# Learning outcomes

- Student can describe the difference between **supervised** and **unsupervised learning** and name advantages of both methods
- Student can apply **K-means** to find **clusters** in data
- Student can explain **Principal Component Analysis** and motivate **dimensionality reduction**
- Student can explain the concept of an **Autoencoder** and motivate why abstract features (latent variables) can be used for a secondary task.

# Learning strategies

## Supervised

Learning from examples (=training data) that are labeled with their desired outputs. The goal is to learn general rules that maps inputs to outputs.

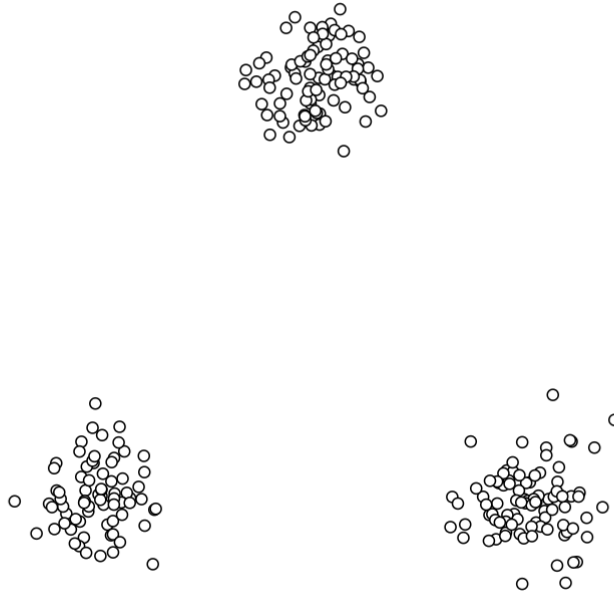
## Unsupervised

Learning from examples without labels. The goals are:

- Learning the entire probability distribution that generated a dataset
- Finding structure in data
- Reducing dimensionality → feature learning

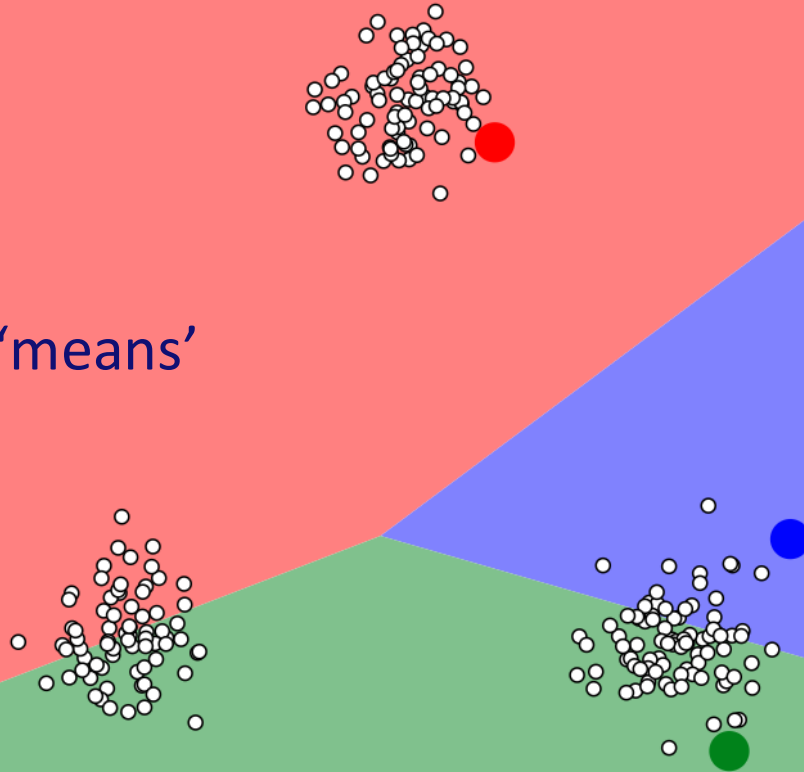
# Finding clusters using K-means

<https://www.naftaliharris.com/blog/visualizing-k-means-clustering/>



Initialization:

Choose K initial 'means'



# K-means – evaluate clustering performance

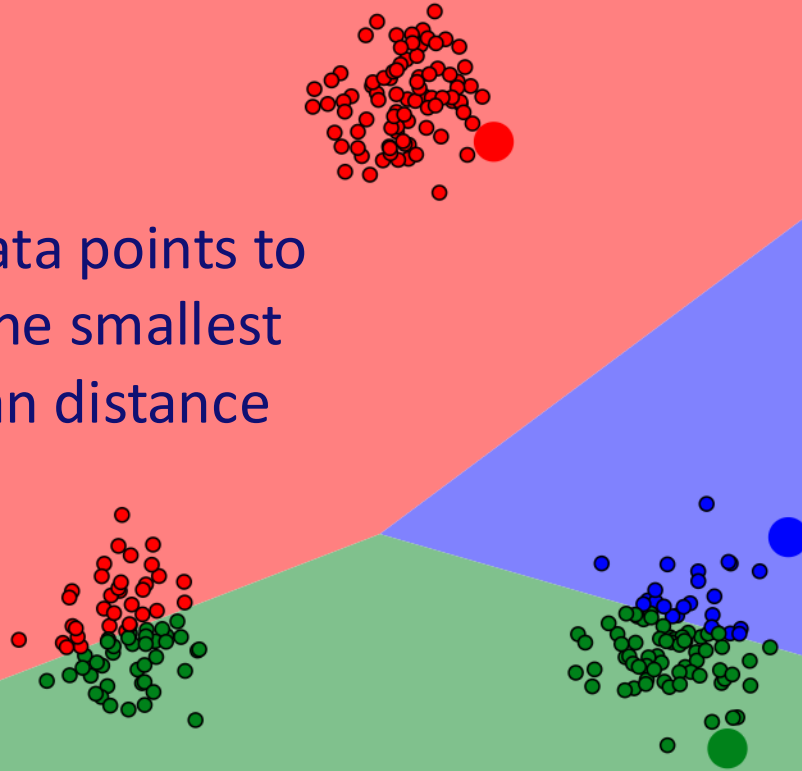
Average squared Euclidean distance between each point and the closest cluster:

$$J(W) = \frac{1}{N} \sum_i ||\min_k (W_k - x_i) ||_2^2$$

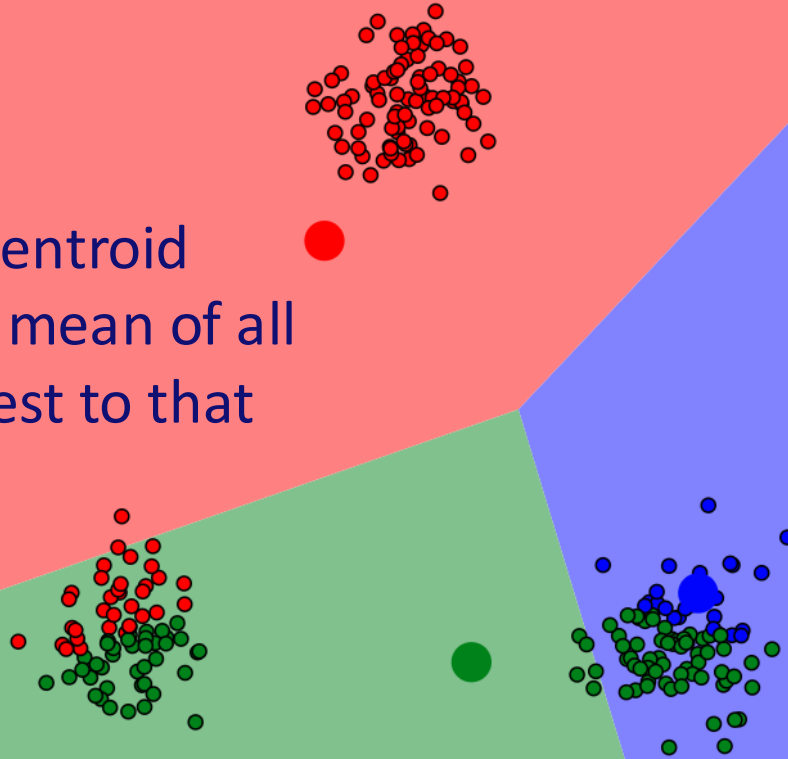
$x_i$  are the points,  $W$  are the cluster centroids



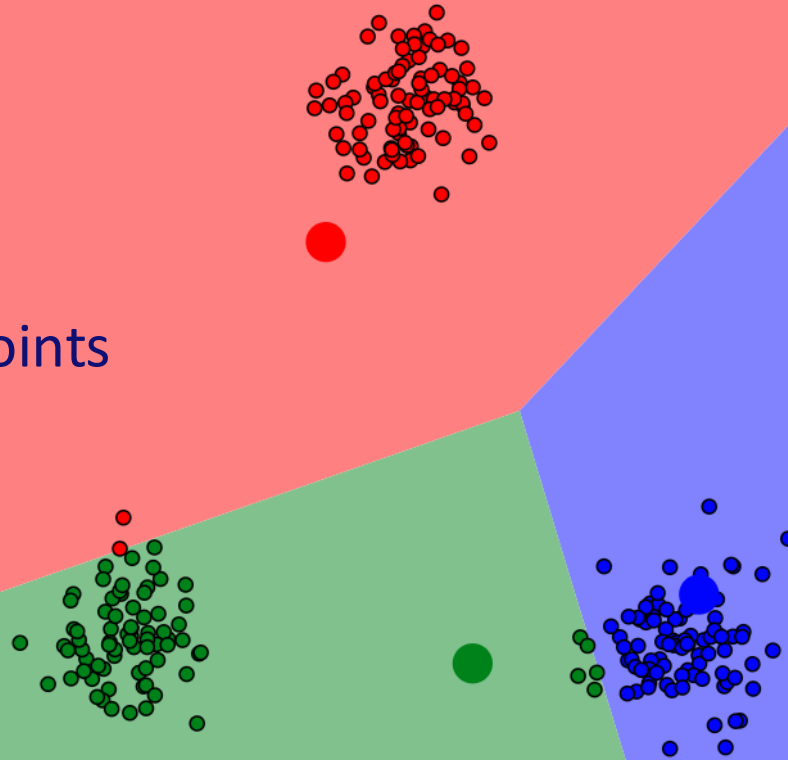
Step 1: assign data points to centroids with the smallest squared Euclidian distance



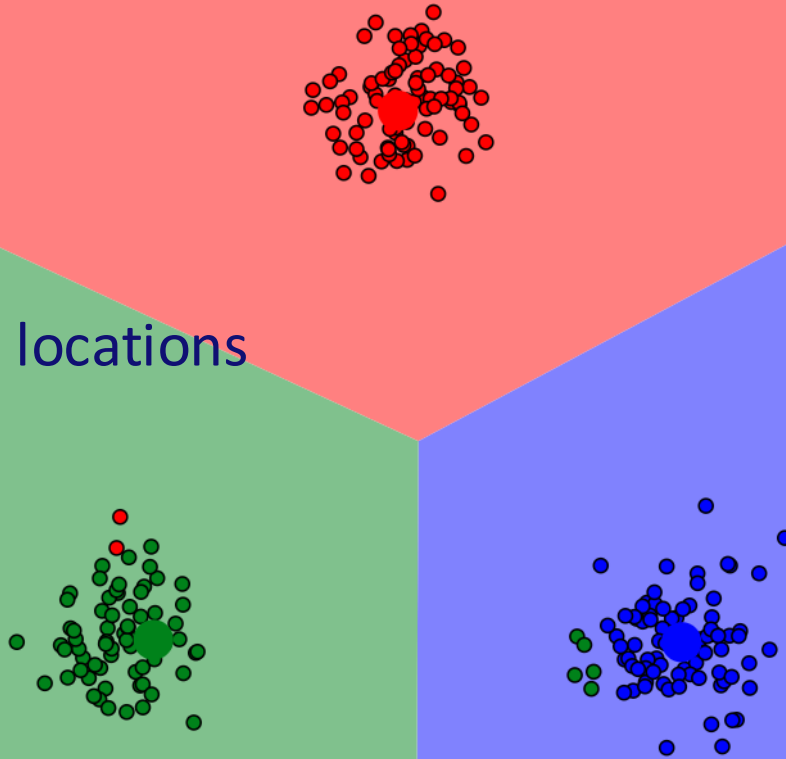
Step 2: update centroid locations by the mean of all data points closest to that centroid



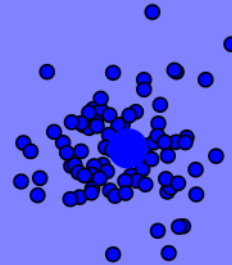
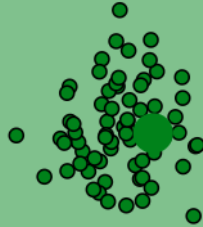
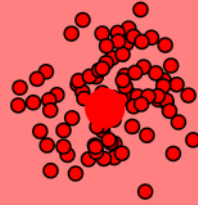
Repeat step 1:  
Reassign data points



Repeat step 2:  
Update centroid locations

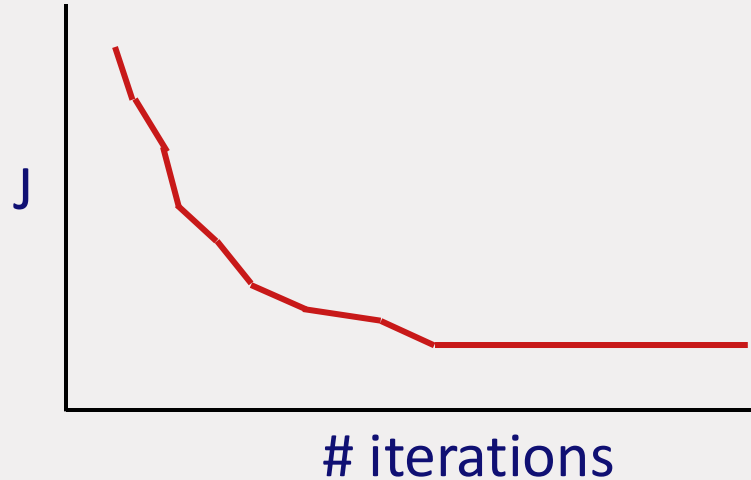


Repeat step 1:  
Reassign data points



# When do we stop?

- When the error  $J$  does not decrease anymore
- After  $n$  iterations

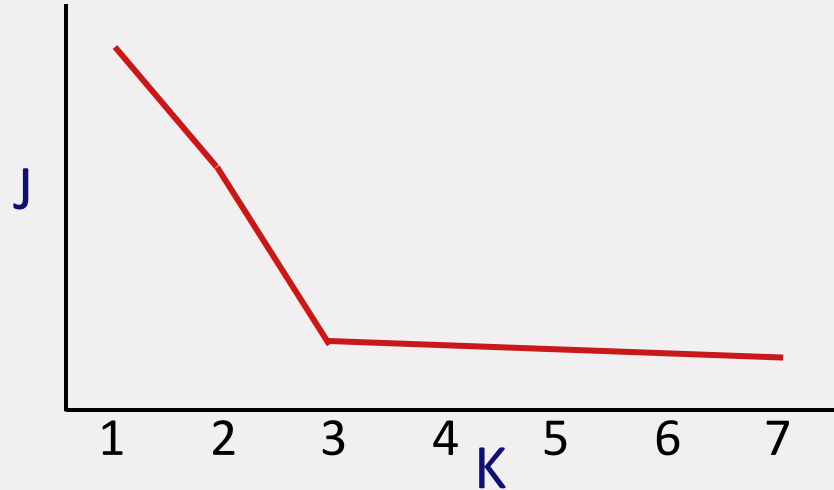


# How do we choose initial centroid locations?

- Random
- Farthest points
- Manual?
  - Supervised
  - Difficult for high-dimensional data

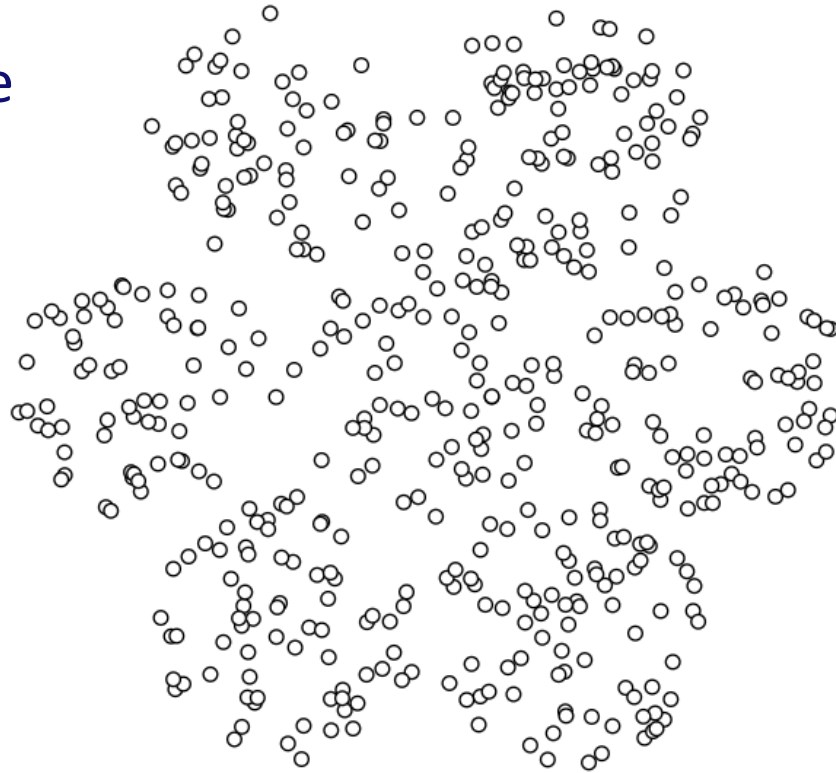
# How do we choose K?

- $K$  (=number of means) is a hyperparameter

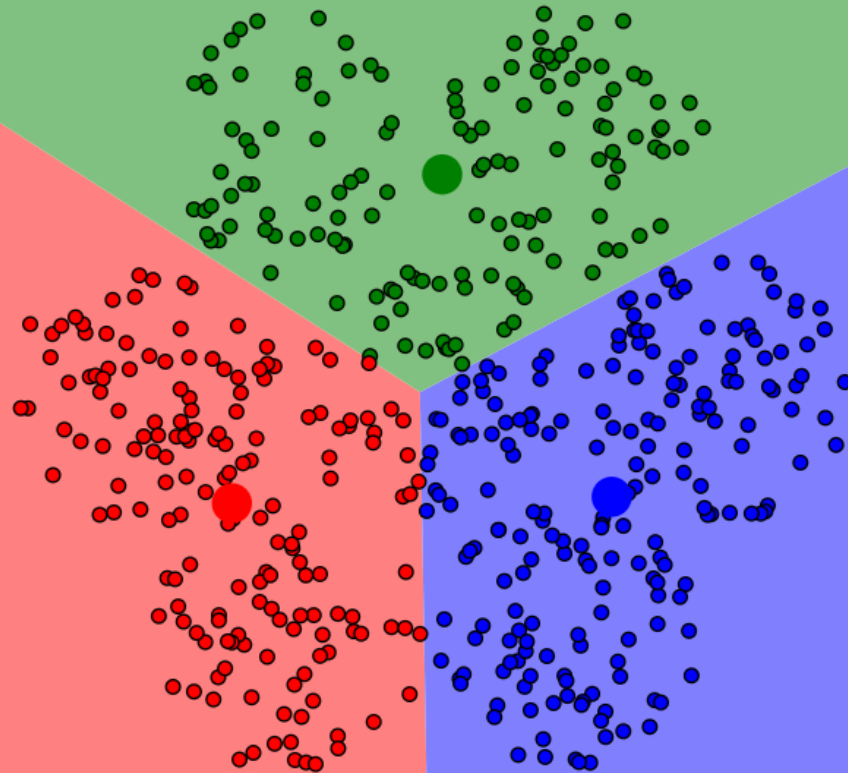




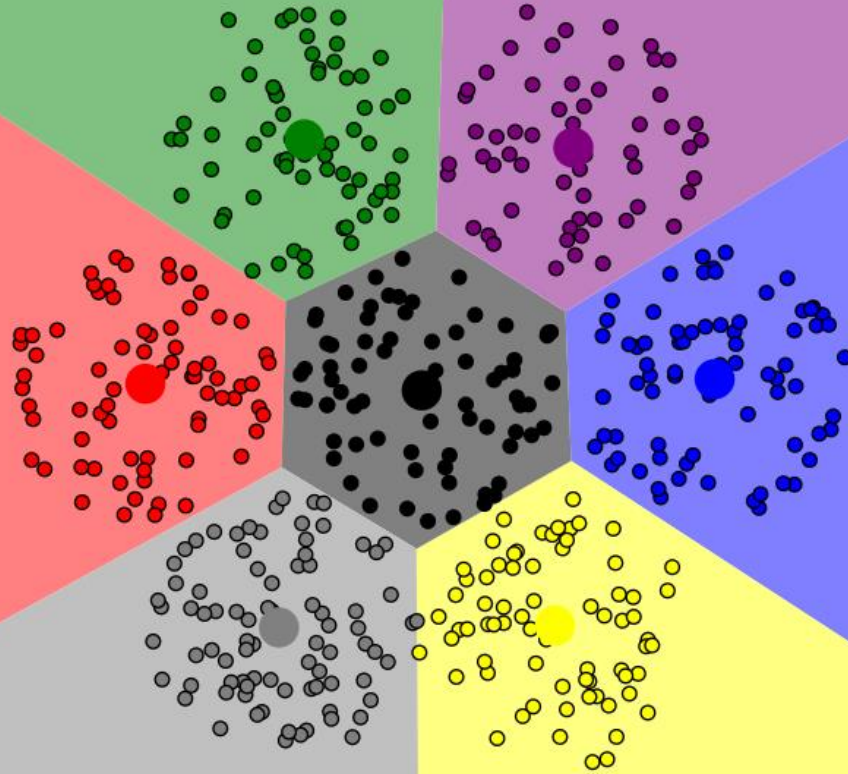
## Another example



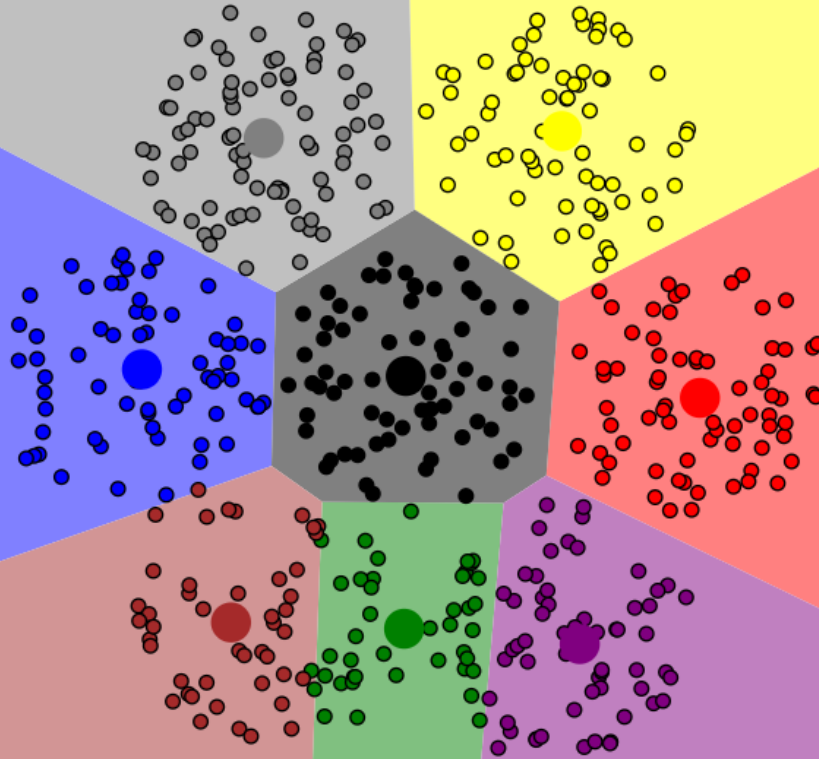
$K = 3$



$K = 7$



$K = 8$



# Questions so far?

# Principal Component Analysis (PCA)

**Goal:** Finding the principal components that describe our data.

= finding the directions in which the data shows most variation

Useful for dimensionality reduction

- E.g. find low-dimensional classification boundaries

Results in better generalization!

# Principal Component Analysis

## Example data set

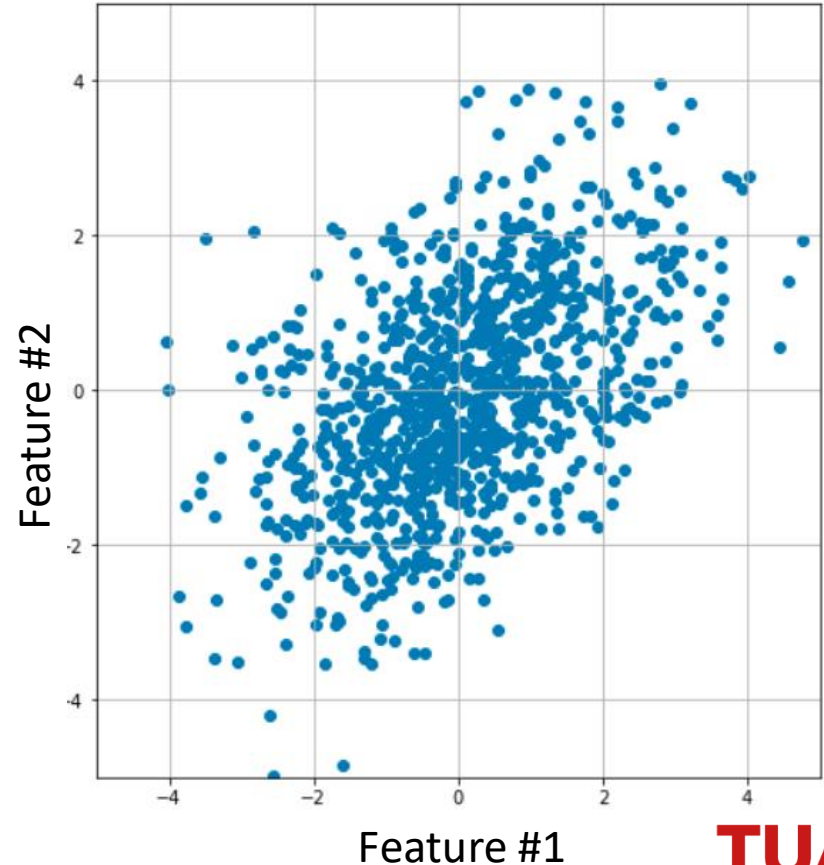
$M$ -by-2 matrix  $X$  containing  $M$  points

Sampled from 2D Gaussian distribution

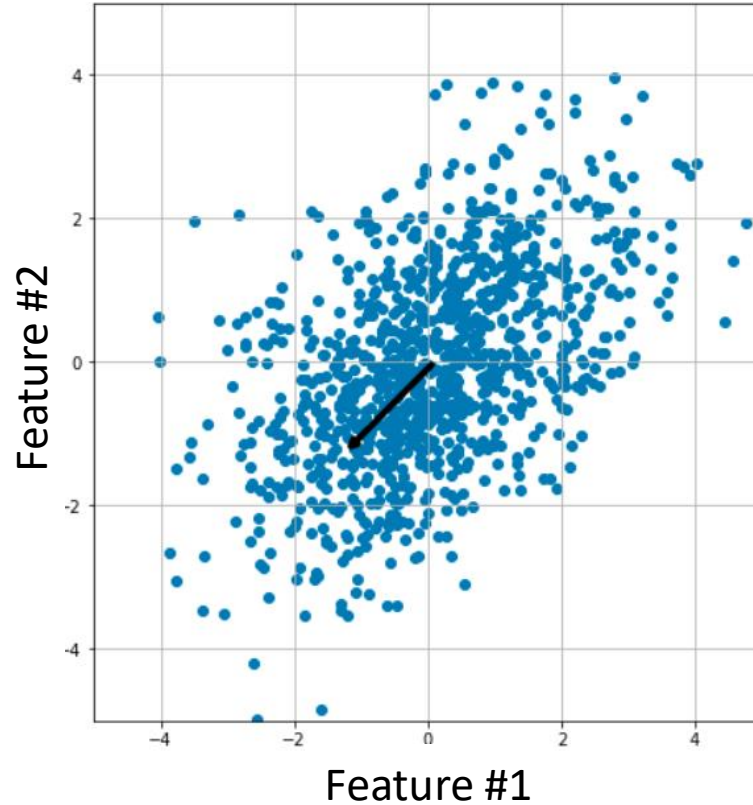
$$\mu_1 = 0$$

$$\mu_2 = 0$$

$$\Sigma = \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}$$

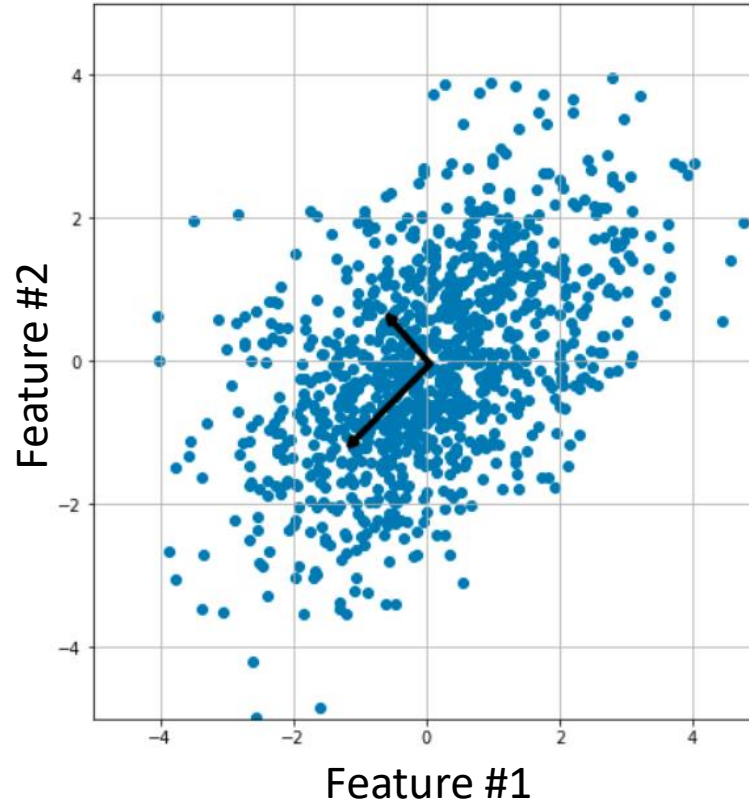


# Principal Component 1





# Principal component are orthogonal to one another



# Finding the principal components

- Center data by subtracting mean of each variable

$$\hat{X} = X - \bar{X}$$

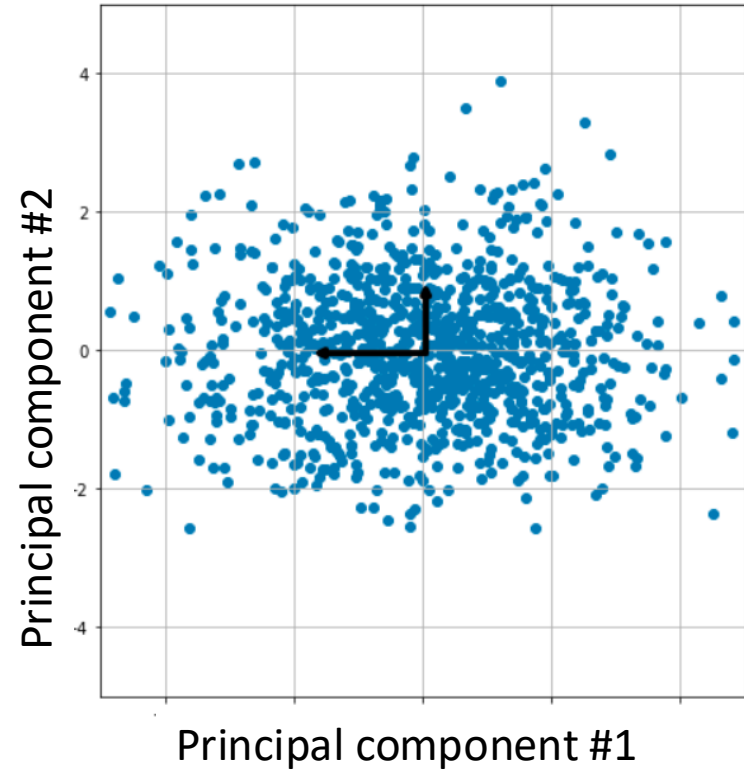
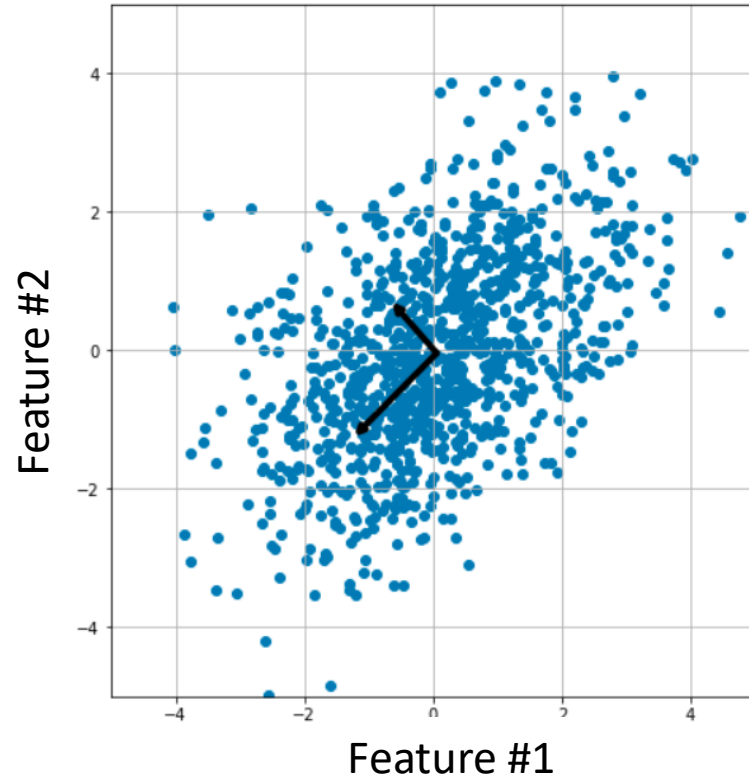
- Calculate covariance matrix

$$\Sigma = \frac{1}{M-1} X^T X$$

- Singular value decomposition (SVD) to find a matrix  $U$  that contains eigenvectors, ordered by largest to smallest variance

→ **principal components**

- Multiply  $\hat{X}$  with  $U$  to obtain  $X_{pca}$

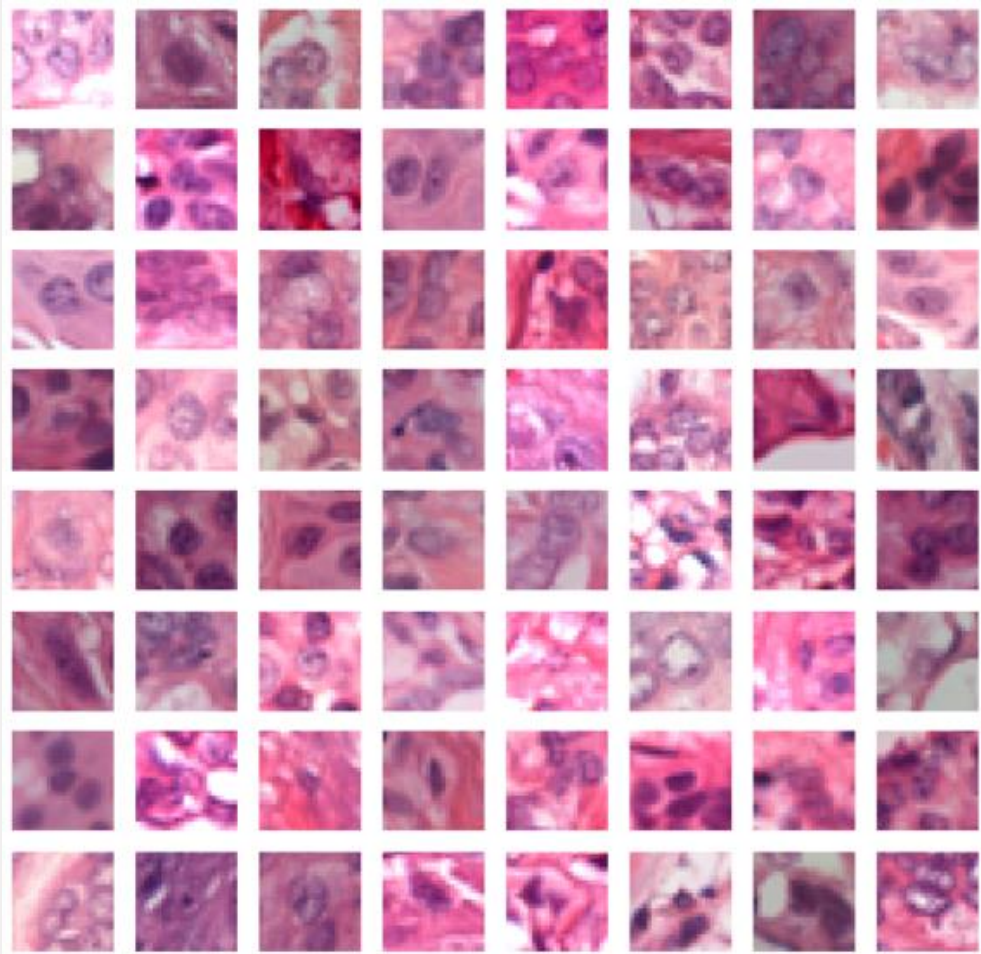


# Dimensionality reduction

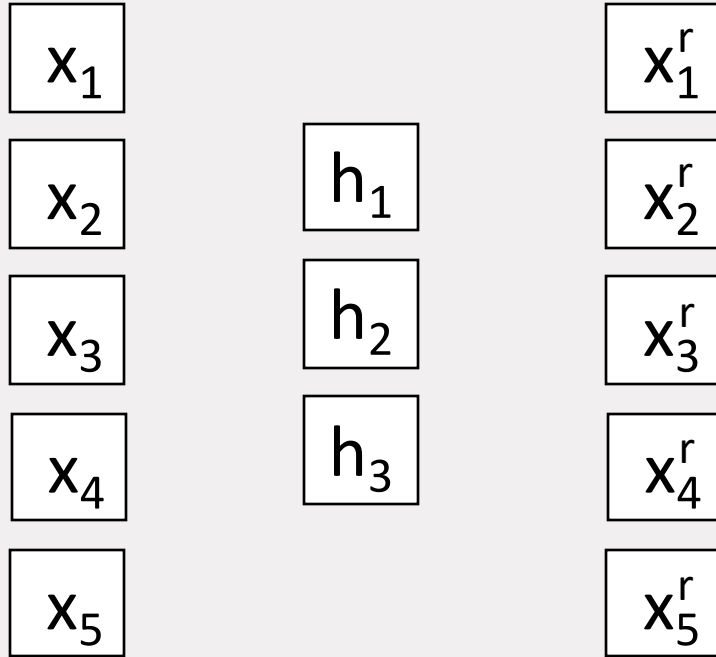
Instead of using all eigenvectors from  $\mathbf{U}$  we can select a set of  $n$  principal components.

For example, we can select the eigenvectors that contain 95% of the variance.

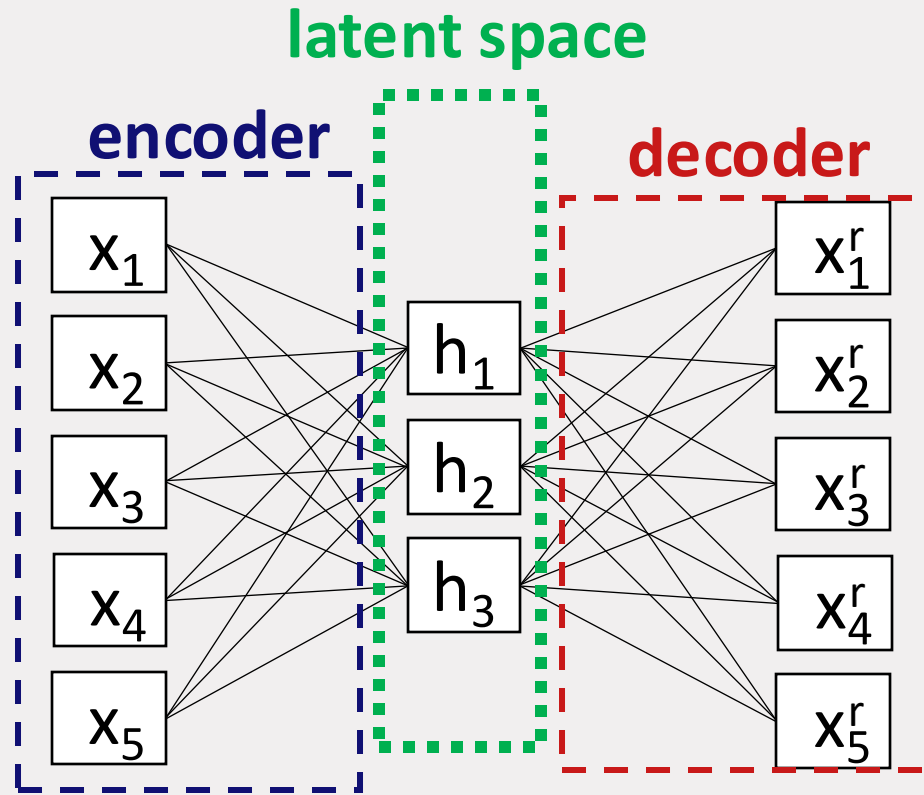
**More info → PCA demo!**

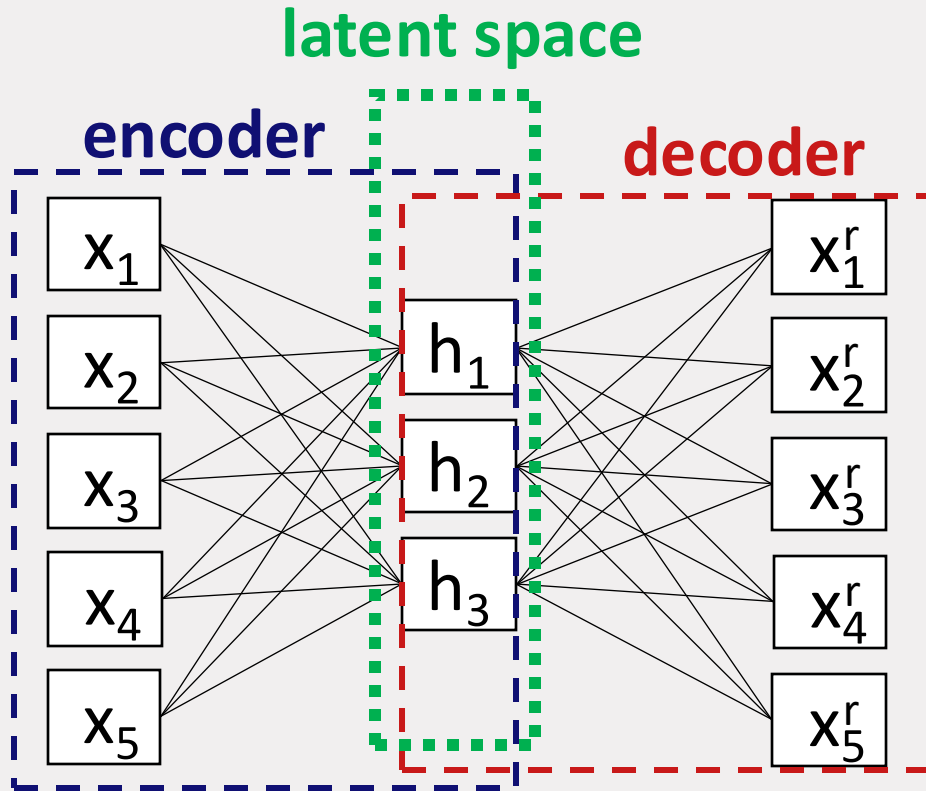


# Autoencoder



Goal = reconstruct input  $x_i$ ,  
using a restricted number of  
latent variables  $h_i$





Encoder:

$$h = f(x)$$

Decoder:

$$x^r = g(h)$$

Penalize dissimilarity

$$L(x, g(f(x)))$$

# Autoencoder

- Encoder/decoder can be simple or complex
- For example: a deep convolutional neural network

## Applications

- Dimension reduction!
- Latent variables can be used for secondary objective, e.g. classification
- Denoising (by adding noise to the input and reconstructing the original)
- Generative models – generating new (image) data



# 'Supervised' learning terminology

**Supervised  
methods**

*Weakly  
Supervised*

*Semi-  
supervised*

**Unsupervised  
methods**

*Self-  
supervised*

# Some remarks on semi-supervised learning

- Fewer labeled data needed
- For many medical applications data is still limited, e.g. because disease is rare
- Use knowledge from a related task
- Humans also learn in a semi-supervised fashion

# Summary

- Supervised versus unsupervised
- Finding structures (e.g. K-means)
- Dimension reduction (e.g. PCA, autoencoders)
- Semi-supervised learning

