# Unsupervised learning

Jordi Villà i Freixa

Universitat de Vic - Universitat Central de Catalunya
Study Abroad

*jordi.villa@uvic.cat*

course 2023-2024

# Índex

# Preliminary note

The material in these slides is strongly based on [1]. When other materials are used, they are cited accordingly.

Mathematical notation follows as good as it can a good practices proposal from the Beijing Academy of Artificial Intelligence.

# What to expect?

In this session we will discuss:

- Unsupervised learning
- The Expectation-Maximization Algorithm
- Principal component analysis
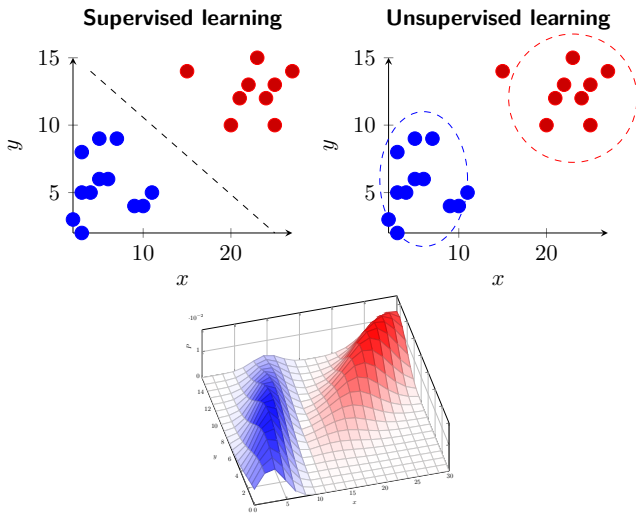
# Unsupervised learning



Figure 1: Supervised-unsupervised learning, and actual pdf behind the data

# Generating data

## Exercise 1

Generate bivariate data Generate data from a collection of three different bivariate distributions with arbitrary means and covariances.

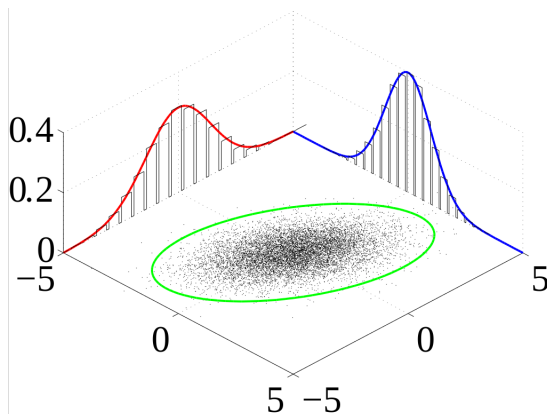# Joint probability distribution



Figure 2: Multivariate normal sample

# Maximum likelihood and Expectation-Maximization algorithm

We already saw how to use the K-means algorithm, based on the *distances* of the different data points. Now we will use another method that is based on assuming a given distribution of the data.

Maximum likelihood estimation (MLE) is an approach to density estimation for a dataset by searching accross probability distributions and their parameters. However, it needs that *latents* (or hidden) variables not to exist in the training dataset. It requires the dataset to be complete, or fully observed: that all variables that are relevant to the problem are present.

The **expectation-maximization (EM)** algorithm is an approach for performing MLE when latent variables are present. It is common to use it in estimating the density when data is missing, such as clustering in the **Gaussian Mixture Model**.

# The EM algorithm

Let us assume a *mixture model*: unspecified combinadtion of multiple probability distribution functions. In our case, we will deal with the *Gaussian mixture model (GMM)*, where the parameters to estimate are the mean and the standard deviation for each Gaussian (normal) distribution in the mixture.
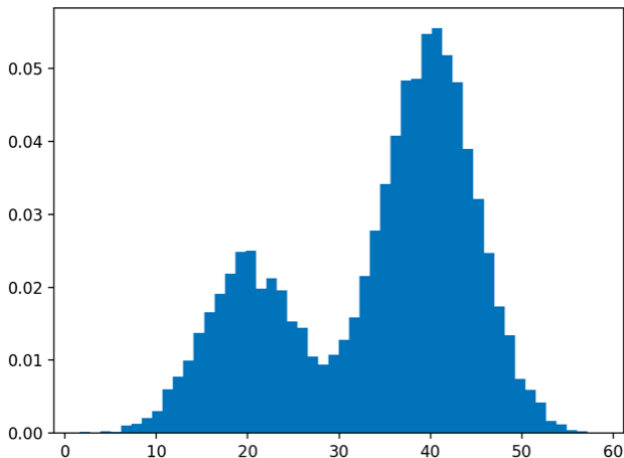
1. E-step: Estimate the missing variables in the dataset: the expected value for each latent variable.
2. M-step: Maximize the parameters of the model in the presence of the data, using MLE.

# A one dimensional example for EM

Adapted from here

```python
# example of a bimodal constructed from two gaussian p
from numpy import hstack
from numpy.random import normal
from matplotlib import pyplot
# generate a sample
X1 = normal(loc=20, scale=5, size=3000)
X2 = normal(loc=40, scale=5, size=7000)
X = hstack((X1, X2))
# plot the histogram
pyplot.hist(X, bins=50, density=True)
pyplot.show()
```

# A one dimensional example for EM

# A one dimensional example for EM

```python
from numpy import hstack
from numpy.random import normal
from sklearn.mixture import GaussianMixture
# generate a sample
X1 = normal(loc=20, scale=5, size=3000)
X2 = normal(loc=40, scale=5, size=7000)
X = hstack((X1, X2))
# reshape into a table with one column
X = X.reshape((len(X), 1))
# fit model
model = GaussianMixture(n_components=2, init_params='random
model.fit(X)
yhat = model.predict(X) # predict latent values
print(yhat[:100])       # check latent value for first few
print(yhat[-100:])      # check latent value for last few p
```

# The maths behind

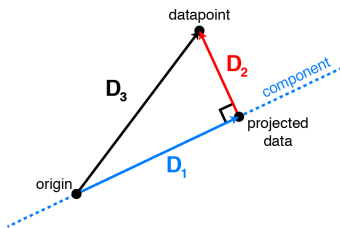$$\mu_k = \frac{\sum_{i=1}^{n} p_i^{(t)}(k)\mathbf{x}_i}{\sum_{i=1}^{n} p_i^{(t)}(k)}$$

$$\mathrm{Cov}_k = \frac{\sum_{i=1}^{n} p_i^{(t)}(k)(\mathbf{x}_i - \mu_k)(\mathbf{x}_i - \mu_k)^T}{\sum_{i=1}^{n} p_i^{(t)}(k)}$$

where we estimate the pdf by using the previous iteration expectation and covariace:

$$p_i^{(t)}(k) \propto w_k^{(t-1)} \phi_k(\mathbf{x}_i | \mu_k^{(t-1)}, \mathrm{Cov}_k^{(t-1)})$$

# PCA

1. Get some data
2. Substract the mean from each data dimension. Better standardize: $z = \frac{x - \bar{x}}{\sigma}$. Calculate the covariance matrix $XX^T$.
3. Diagonalize and obtain unit eigenvectors $c^T X X^T c - \lambda c^T c$ via SVD: $XX^T = UD^2U$
4. Choosing components and forming a feature vector
5. Recast the data along the chosen PC axes. $\mathbf{x}_i^{\mathrm{PC}} = U_k U_k^T \mathbf{x}_i$



Figure 3: Maximizing variance in principal component space is equivalent to minimizing least-squares reconstruction error. Taken from this post.

# EX2

## Exercise 2

EX2. PCA example Get the *Iris* dataset from the UCI repository.

1. Using the **kdeplot** method of **seaborn**, generate a figure for the kernel density plots of the variable *Petal.length* for each of the three species of *Iris sp.* in the dataset.

2. Show a scatter plot of the variable Petal.length vs the variable Sepal.length and comment on the correlation you see.

3. In order to analyze the possible correlatoons between the different data in the file, perform a PCA and obtain the principal component matrix $U$ as well as the diagonal matrix $D^2$. What can you say in terms of the variance?

4. Project the data on the first component and plot the kernel density estimate of the PCA-combined data.

📄 Dirk P. Kroese, Zdravko Botev, Thomas Taimre, and Radislav: Vaisman.
*Data Science and Machine Learning: Mathematical and Statistical Methods*.
Machine Learning & Pattern Recognition. Chapman & Hall/CRC, 2020.