

Introduction

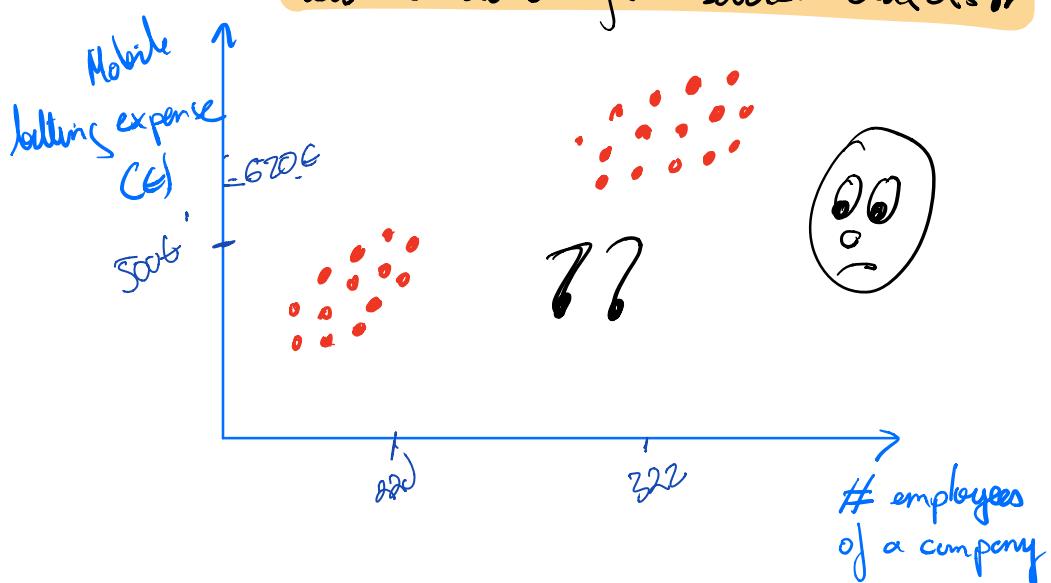
From previous lessons, we know:

- ① K-means don't work well with **non-circular clusters**
- ② Gaussian models will allow to model our data with a **Gaussian distribution** for each cluster

$\left\{ \begin{array}{l} \text{means } (\mu) \Leftrightarrow - \text{center} \\ \text{covariance} \Leftrightarrow - \text{direction and shape} \\ (\Sigma) \end{array} \right.$

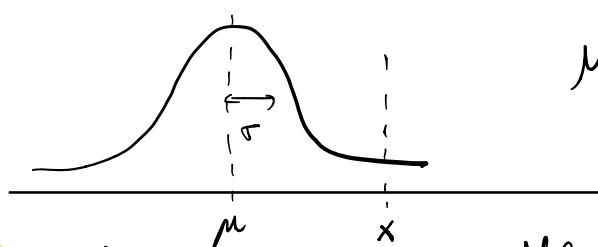
\Rightarrow We have more flexible to capture clusters with **different shapes** and less sensitive to **rescale of variables**

\Rightarrow Until now, we know to model 1 cluster. But, **how to do it for several clusters?**



Recall:

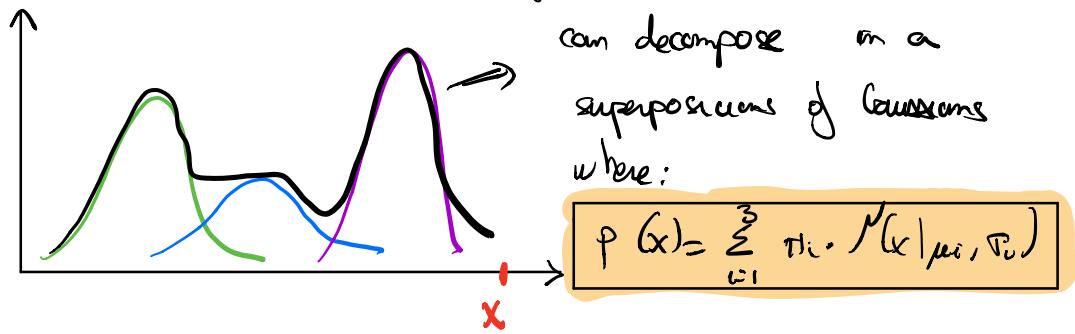
The Gaussian distribution in 1D is:



$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

The probability of observing x given $N(\mu, \sigma^2)$ is $p(x) = N(x|\mu, \sigma^2)$
or likelihood

When the dataset "is organized" in clusters, we can consider that our data follows a complex function distribution and we

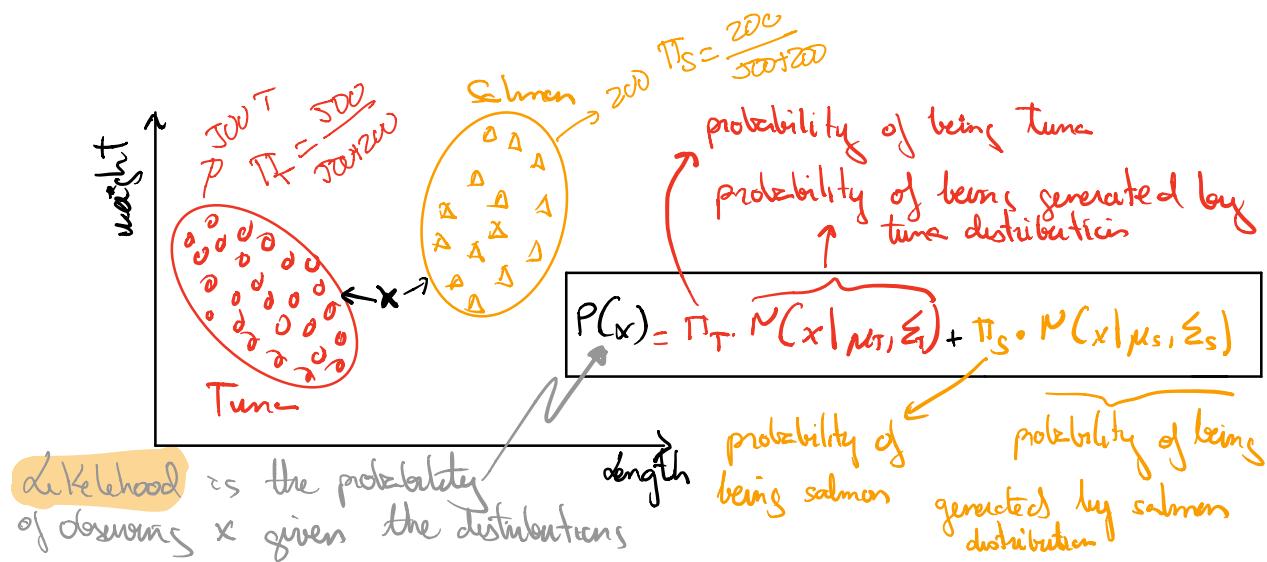


where π_k is prior probability that any datapoint is assigned to cluster k without observing the particular datapoint

where $N(x|\mu_k, \sigma_k^2)$ is the likelihood or the probability of observing x given a gaussian distribution K

An example:

Let's imagine a dataset with **salmon** and **tuna**. For a new datapoint x , the probability of observing x is:



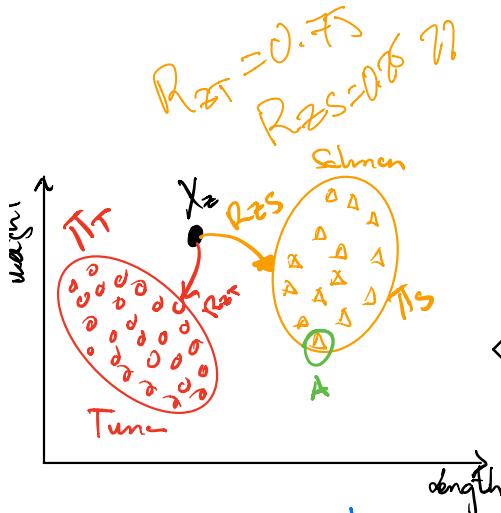
For the whole dataset \vec{X} of K datapoints, where x_i is a datapoint with dimension D , the total likelihood function is:

$$P(\vec{X} | \pi, \mu, \Sigma) = \prod_{i=1}^K \left(\sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k) \right)$$

As it's too computational cost, we calculate the log-likelihood:

$$\ln(P(\vec{X} | \pi, \mu, \Sigma)) = \sum_{i=1}^K \left[\left(\sum_{k=1}^K \pi_k \cdot N(x_i | \mu_k, \Sigma_k) \right) \right]$$

The log-likelihood is the function to optimize in the training stage.



Data points are not hard assigned to each cluster. Every datapoint, x_i , has a responsibility or probability of being generated by each gaussian component.

We calculate the derivative and we obtain :

$$\mu_k = \frac{1}{N_k} \sum_{i=1}^N r_{ik} \cdot x_i , \quad N_k = \sum_{i=1}^N r_{ik}$$

$$\Sigma_k = \frac{1}{N_k} \cdot \sum_{i=1}^N r_{ik} \cdot (x_i - \mu_k) (x_i - \mu_k)^T$$

$$\pi_k = \frac{N_k}{N}$$

where r_{ik} is the responsibility of datapoint i respect to cluster k ; i.e r_{ik} is the posterior probability that a datapoint x_i is generated by component or cluster k :

$$p(k, x_i) = \frac{p(k) \cdot p(x_i | k)}{\sum_{l=1}^k p(l) \cdot p(x_i | l)} = \frac{\pi_k \cdot N(x_i | \mu_k, \Sigma_k)}{\sum_{l=1}^k \pi_l \cdot N(x_i | \mu_l, \Sigma_l)}$$

Therefore, μ_k , Σ_k and π_k are weighted parameters according to datapoint's responsibilities

Therefore, our goal is to calculate the gaussian distributions, i.e. μ_k , Σ_k , π_k that give a maximum value of $P(\vec{x})$ for the whole dataset where : $\ln(p(\vec{x} | \pi, \mu, \Sigma)) = \sum_{i=1}^N \left[\sum_{k=1}^K \pi_k \cdot p(x_i | \mu_k, \Sigma_k) \right]$

is the log-likelihood of dataset \vec{x}

$$P(S, A) = \frac{P(S) \cdot P(A|S)}{P(S) \cdot P(A|S) + P(T) \cdot P(A|T)}$$

But, what does it mean?

$$P(T, A) = P(T) \cdot P(A|T)$$
$$\underbrace{P(T)}_{P(T) + P(\bar{T})} \cdot \underbrace{P(A|T)}_{P(A|T) + P(\bar{A}|T)}$$

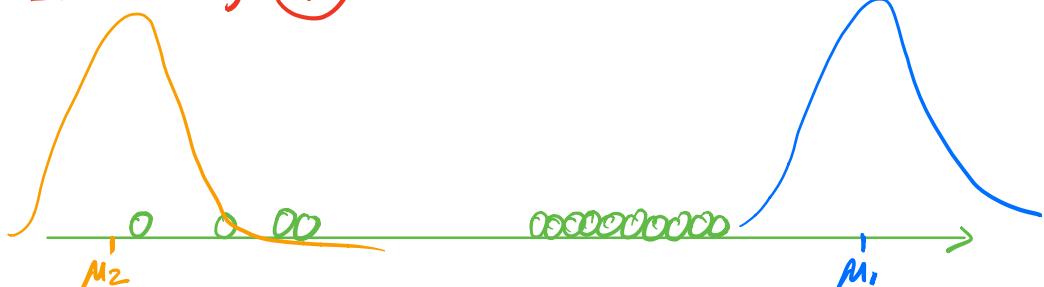
Imagine thus 1-D datapoints distribution

$$R_{P,T}$$

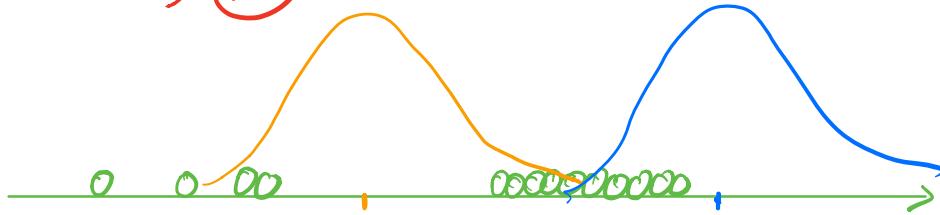


Which is the best Gaussian distribution??

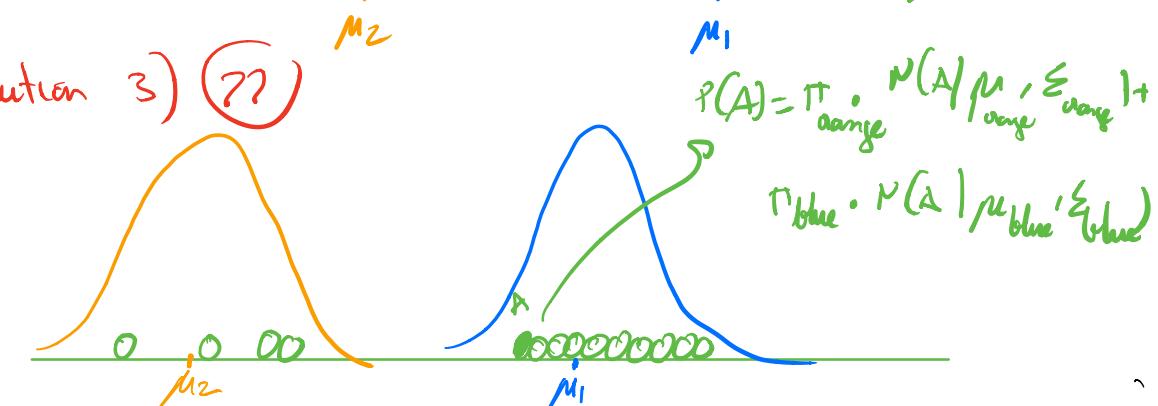
Solution 1) ??



Solution 2) ??



Solution 3) ??



⇒ Gaussian distributions of solution 3) maximizes the likelihood of the dataset because it maximizes $P(x) = \sum_{k=1}^2 \pi_k N(x | \mu_k, \sigma_k^2)$ $x \in \text{dataset}$

The Expectation - Maximization (EM) algorithm

The issue: When we have a new dataset \vec{X} , we would to calculate μ_k , Σ_k and π_k for every cluster k . However, we need to know $N(\vec{x}, \mu_k, \Sigma_k)$ of every cluster or component to calculate R_{ik} and to calculate $N(\vec{x}, \mu_k, \Sigma_k)$ we need to know μ_k , Σ_k and π_k .

We are in the EGG-CHICKEN problem



The solution:

Stage 0: We initialize randomly (or smarter) μ_k, Σ_k, π_k for every k

Stage 1: We calculate the membership or responsibility of each datapoint x_i (estimation) to every gaussian component

For each x_i , which its the probability to be generated by each component $k \Rightarrow R_{ik}$ for all k

Stage 2: We re-calculate or update each Gaussian based on the datapoints' responsibilities, R_{ik} :

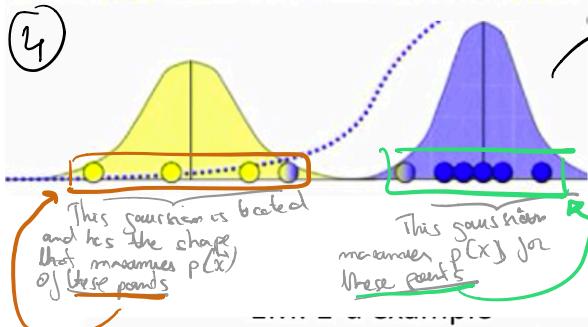
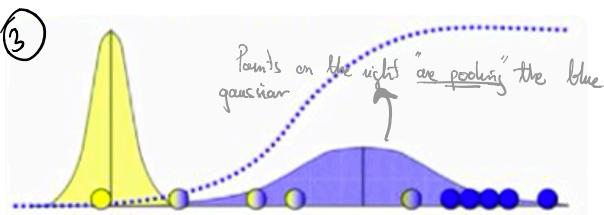
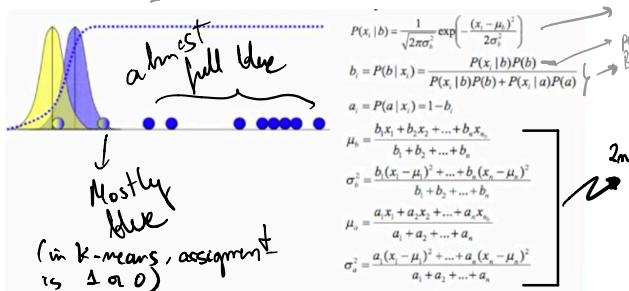
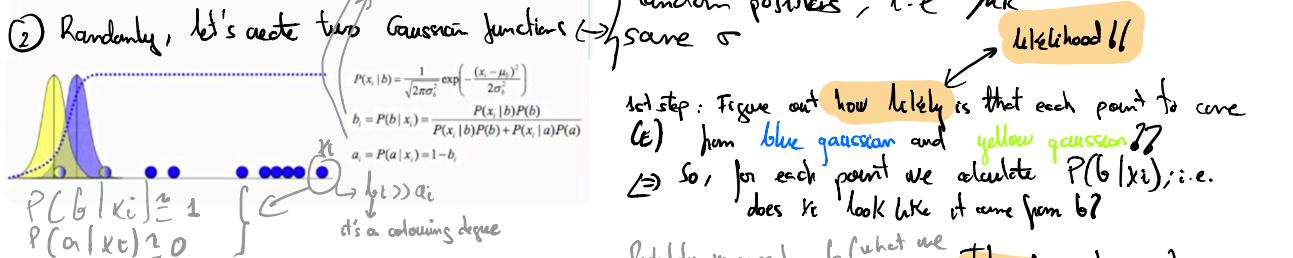
For each k , calculate μ_k, Σ_k and π_k

Example 10:-

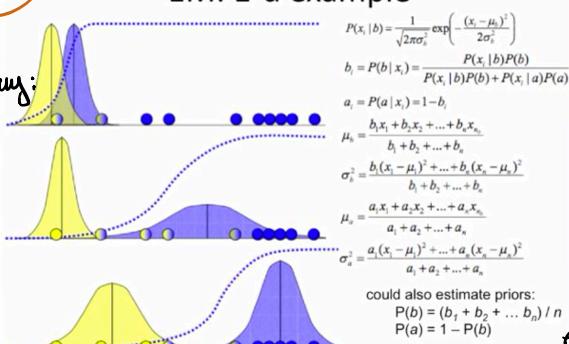
① Let's suppose this is dataset:

○ ○ ○○ ○○○○○

probability of x_i respects to
respectability of b respects to "a"



Summary:



random positions; i.e. μ_k
same σ

likelihood \mathcal{L}

1st step: Figure out how likely is that each point to come from blue gaussian and yellow gaussian \Rightarrow
 \Leftrightarrow So, for each point we calculate $P(b | x_i)$; i.e. does it look like it came from b ?

Probability x_i comes from b (what we know)

prior probability given distribution

Thanks to posterior probability, we can calculate $P(b | k_i)$

2nd step: Calculate the new μ_b , μ_y with the datapoints' fraction assigned to each cluster

Following this process, we guarantee the optimisation of the function of likelihood for the whole dataset \bar{x} ; i.e. likelihood takes a maximum value

From the final picture, we see that we achieved a maximum value for $P(\bar{x} | \pi, \mu, \varepsilon)$ due to gaussians are located at centers with \leq to maximize the value likelihood value for the whole dataset.

$$\ln(P(\bar{x} | \pi, \mu, \varepsilon)) = \sum_{i=1}^n \left[\left(\sum_{k=1}^K \pi_k P(x_i | \mu_k, \varepsilon_k) \right) \right]$$

↳ probability of observing \bar{x} given the distributions

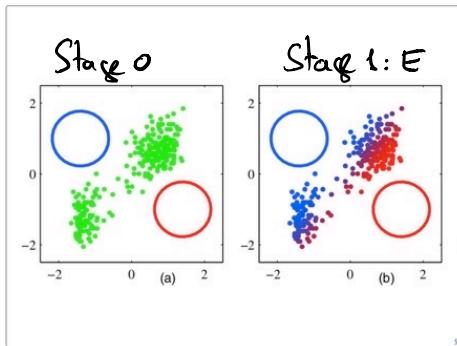
Stage 0: Invent gaussian distributions for blue and yellow.

the distributions

Stage 1: Estimate for each x_i their probability to come from blue a yellow distributions:
 $P(b | x_i) \rightarrow$ by posterior prob

Stage 2: Recalculate new distributions according to the probability of each k_i to come from each distribution

Example 2-D:



Here we see EM in action. We start with two random Gaussians and color the points by how much responsibility each component takes. The blue points are mostly claimed by the blue component, the red points are mostly claimed by the red, and the purple points in the middle have the responsibility divided equally between the components.

source: Machine Learning and Pattern Recognition,
Christopher Bishop.

