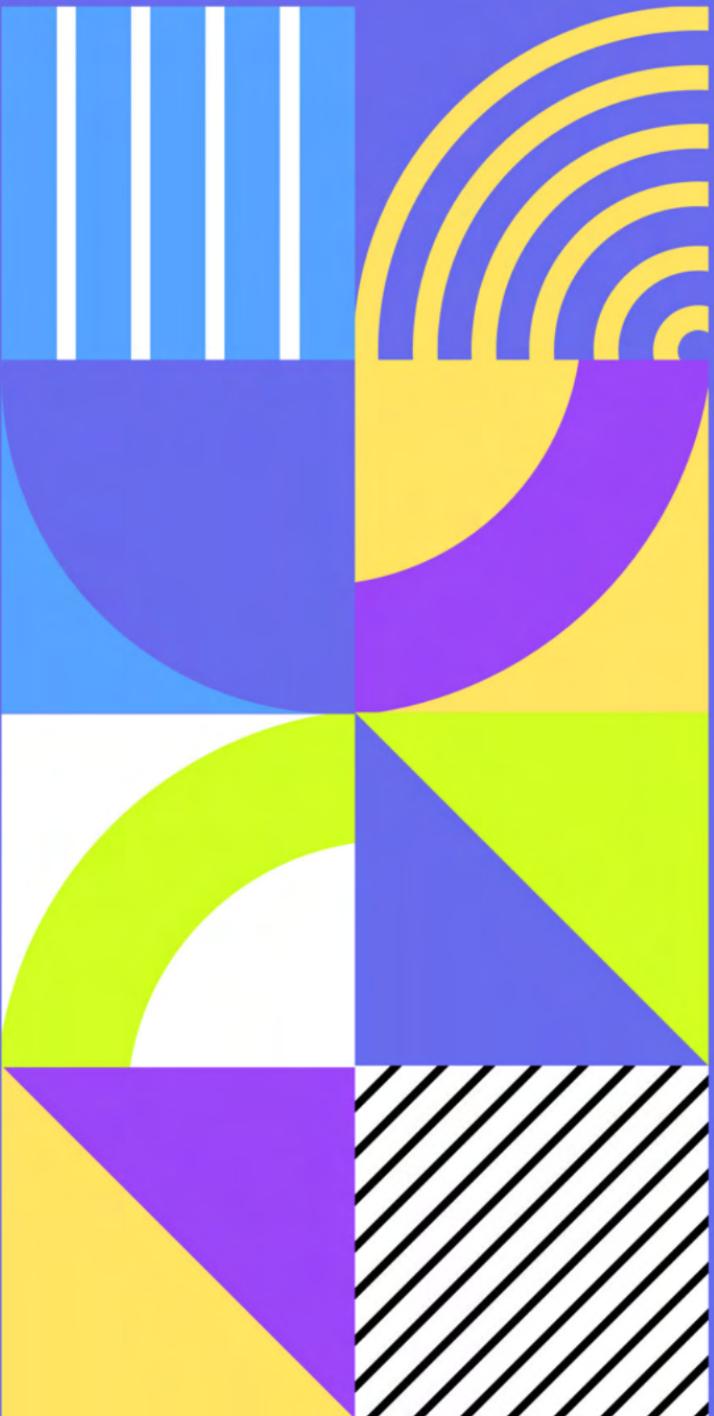


JNOTE.

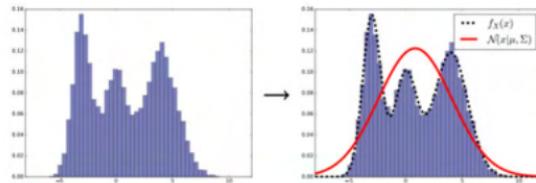


# Gaussian Mixture Models

We have seen that we can solve classification problems by building generative models that describe the distribution of samples

The Gaussian classifier is an example that assumes that class-conditional distributions are Gaussian

In many cases, however, the assumption can be quite inaccurate



single Gaussian fall short because:

- Real-world data rarely follows a perfect bell curve
- Many datasets have multiple peaks (multimodal distributions)
- Some distributions have asymmetric shapes or heavy tails
- Data clusters might appear in separate regions

So when a single Gaussian isn't sufficient we need more flexible models. We have several options:

- ① Domain-specific distribution: In some fields, we might know that data follows specific distributions (e.g. Poisson for count data)
- ② Gaussian Mixture Models (GMMs): A more general solution that can approximate virtually any smooth distribution by combining multiple Gaussian components.

GMMs can model complex multimodal distributions, it's possible to control the approximation precision by adding more components, GMMs maintain mathematical tractability while handling complex shape

Of course, since we are estimating the density from data, we require a sufficient amount of data to obtain good estimates

The use of GMMs is not restricted to classification

GMMs can be employed also in other tasks that require estimating a population density (e.g. data generation)

As we will see, they also allow to solve different kind of problems

For example, GMMs provide an alternative to K-means for clustering

It can be used also for open set classification

- We have already encountered an example of GMN in the Gaussian classifier model
- In this classifier we have modeled each class with its own Gaussian distributions:  $f_{x|C}(x|c) = N(x|\mu_c, \Sigma_c)$
  - Assigned prior probabilities to each class  $P(C=c) = \pi_c$
  - Calculate the overall data distribution as a weighted sum of these class-conditional distributions.

This overall distribution is what we now call a Gaussian Mixture Model:

$$f_x(x) = \sum_{c=1}^K \pi_c N(x|\mu_c, \Sigma_c)$$

Digression on it:

First, remember that in a classification setting, we have K different classes. For each class c, we have:

A prior probability  $\pi_c$  that represents how likely any random data point belongs to class c  $P(C=c) = \pi_c, c = 1, \dots, K$

A class-conditional density  $N(x | \mu_c, \Sigma_c)$  that describes the distribution of data points within that class

Now, to find the probability density of observing any particular data point x (regardless of which class it comes from), we need to consider all possible classes. This uses a fundamental concept in probability called the "law of total probability." For each potential class c:

We multiply the probability that a point comes from class c ( $\pi_c$ )

By the probability density of observing x given that it's from class c ( $N(x | \mu_c, \Sigma_c)$ )  $N(x | \mu_c, \Sigma_c)$

Then we add up these products for all classes

This gives us the overall probability density for any point x, without knowing which class it belongs to:

$$f_x(x) = \sum_{c=1}^K \pi_c N(x | \mu_c, \Sigma_c)$$

Think of it like if we randomly pick a data point from our entire dataset, this formula tells us how likely we are to observe a specific value x, considering all possible classes that could have generated it.

Each component represent a different source that could have generated our data point.

## Note on Gaussian Classification -

POTREBBE SERVIRE IN UN MODELLA DI CLASSIFICAZIONE A INSERIRE UNUSO (TUTTI LE POSSIBILI CLASSI SONO NOTE & IL CAMPO DEVE APPARTENERE AD UNA DI QUESTE). APPLICATO SOLO PER TRAMITE STA MODALITÀ

$$P(C=c | X_t=x_t) = \frac{f_{X,C}(x_t, c)}{\sum_{c' \in C} f_{X,C}(x_t, c')}$$

$$\sum_{c' \in C} f_{X,C}(x_t, c') = f_X(x_t)$$

Dove il denominatore sono tutte le possibili classi  $\{c\}$  un valore costante che normalizza la probabilità a posteriori garantendo che sommato tutte le  $f_{X,C}$  il valore sia 1.

Di per sé è la densità di probabilità multivariata dove  $X=x_t$  valore fisso e  $C$  una delle  $K$  possibilità

## Formal definition of GMs

More generally a Gaussian Mixture model represents a probability distribution as a weighted sum of multiple Gaussian components:

$$f_X(z) = \sum_{c=1}^K w_c N(z|\mu_c, \Sigma_c)$$

Where:

- $w_c$  are the mixture weights (analogous to the class priors in classification)
- $\mu_c$  are the means of each Gaussian component.
- $\Sigma_c$  are the covariance matrix for each component
- $K$  is the number of Gaussian components in our mixture.

**Visual intuition:** Imagine trying to model a mountain range. Each individual mountain can be represented by a Gaussian "bump," and the entire range is the combination of all these bumps, each with its own height, width, and position.

To fully specify a GM, we need three sets of parameters:

- Component means  $M = [\mu_1, \mu_2, \dots, \mu_K]$ , protuberance
  - These determine the center of each Gaussian "bump"
  - Each  $\mu_i$  is a vector with the same dimensionality of our data.
- Component Covariance  $S = [\Sigma_1, \Sigma_2, \dots, \Sigma_K]$ 
  - This determine the shape, orientation, and spread of each Gaussian.
  - Each  $\Sigma_i$  is a positive-definite matrix of size  $o \times o$  ( $o$  is the data dimensionality)
- Mixture weights:  $w = [w_1, w_2, \dots, w_K]$ 
  - This determine the relative contribution of each component to the overall mixture
  - Large weights mean that component has more influence in the final distribution.

$$X \sim GM(M, S, w)$$

$$f_X(z) = \sum_{c=1}^K w_c N(z|\mu_c, \Sigma_c)$$

The distribution parameters are the component means

$$M = [\mu_1 \dots \mu_K]$$

the component covariances

$$\mathcal{S} = [\Sigma_1 \dots \Sigma_K]$$

and the weights

$$\mathbf{w} = [w_1 \dots w_K]$$

Remember that, for  $f_X$  to be a density, we need that its integral is equal to 1. Integrating w.r.t.  $x$  we have:

$$\int f_X(x) = \int \sum_{c=1}^K w_c \mathcal{N}(x; \mu_c, \Sigma_c) dx \sum_{c=1}^K w_c \int \mathcal{N}(x; \mu_c, \Sigma_c) dx = \sum_{c=1}^K w_c = 1$$

if the weights  
are the prior(s)  
then sum of

i.e. the weights must sum to 1

When working with a dataset  $D = [x_1, x_2, \dots, x_n]$ , GMMs provide us with a powerful framework for modeling the underlying distribution. We make a fundamental assumption: all samples have been independently generated by the same GMM distribution.

More formally, we treat the data points as independent and identically distributed (i.i.d.) random variables, where each  $X_i$  follows the same distribution  $X \sim \text{GMM}(M, S, w)$ .

This is important to understand in context:

$$X_i \sim X \sim LN(\mu, \sigma^2) \quad \text{sum of diff. rest Gaussian component}$$

- In density estimation, we're simply trying to model the overall distribution of our data
  - Unlike supervised classification, we're not assigning predefined class labels
  - The dataset D is considered unlabeled - we don't know in advance which Gaussian component generated each point
  - The GMM could be modeling data from a single class (if used within a classification framework) or an entire dataset without any class distinctions

Note: What does  $x_i \sim \mathcal{N}(600, 3^2)$  mean?

$$P(\text{Coh}) = \omega_W$$

we want to find  
the parameter of  
each distribution



Just as with "simpler" gaussian model we want to find the parameters that maximize the likelihood of observing our data. However, GMs model introduce a significant challenge that single gaussian don't have: the likelihood function is unbounded.

This creates what mathematicians call an "ill-posed problem".

Imagine we place one Gaussian component directly on a single data point and make its variance extremely small (approaching zero). This creates an extremely narrow tall peak. As the variance gets smaller and smaller:

- The probability density at that exact point approaches infinity
  - The likelihood becomes unbounded.
  - The solution becomes degenerate (mathematically valid but practically useless)

In practice we address this issue with regularization techniques and constraint on the covariance matrices to prevent these degenerate solutions. With these safeguards, maximum likelihood estimation works well for GMMs.

## Going deeper

A degenerate solution occurs when one of the Gaussian components becomes pathologically "concentrated" on a very small number of data points (often just one point).

Here's what happens step by step:

### The Pathological Scenario

Imagine we have a GMM with K=2 components, and we place one component directly on a single data point, say  $x_1$ :

1. Set the mean:  $\mu_1 = x_1$  (place the Gaussian exactly on the data point)
2. Shrink the variance: Make  $\Sigma_1 = \sigma^2 I$  where  $\sigma \rightarrow 0$  (make the Gaussian extremely narrow)
3. Keep reasonable weight:  $w_1$  could be any reasonable value, say 0.3

### What Happens to the Likelihood?

The probability density of a Gaussian at point  $x_1$  is:

$$N(x_1 | \mu_1, \Sigma_1) = \frac{1}{(2\pi)^{d/2} |\Sigma_1|^{1/2}} \exp\left(-\frac{1}{2}(x_1 - \mu_1)^T \Sigma_1^{-1} (x_1 - \mu_1)\right)$$

When we place  $\mu_1 = x_1$  and make  $\Sigma_1 = \sigma^2 I$  with  $\sigma \rightarrow 0$ :

• The exponent becomes:  $\exp\left(-\frac{1}{2} \frac{\|x_1 - \mu_1\|^2}{\sigma^2}\right) = \exp(0) = 1$

• The normalization factor becomes:  $\frac{1}{(2\pi\sigma^2)^{d/2}} = \frac{1}{(2\pi)^{d/2}\sigma^d}$

So:  $N(x_1 | \mu_1, \Sigma_1) = \frac{1}{(2\pi)^{d/2}\sigma^d}$

As  $\sigma \rightarrow 0$ , this probability density  $\rightarrow \infty$ .

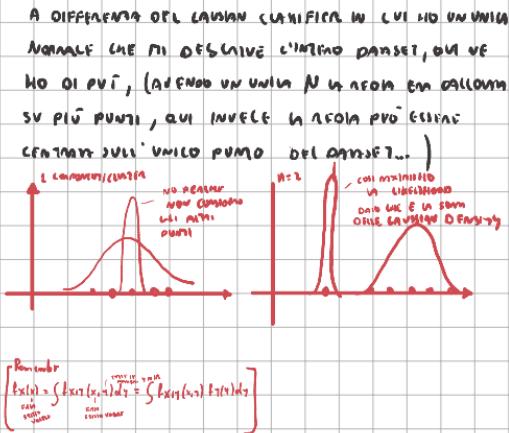
### Impact on the Total Likelihood

The likelihood contribution from point  $x_1$  becomes:

$$L_1 = w_1 N(x_1 | \mu_1, \Sigma_1) + w_2 N(x_1 | \mu_2, \Sigma_2) \approx w_1 \cdot \frac{1}{\sigma^d} + \text{bounded term}$$

As  $\sigma \rightarrow 0$ , this grows without bound, making the total likelihood unbounded.

Indeed, as long as we have more than 1 component, we can devise degenerate solutions for which the likelihood is not bounded above



Let's define our model parameters as  $\Theta = [\pi, \Sigma, w]$  which includes all component means, covariance and weights. The likelihood function is:

$$L(\Theta) = \prod_{i=1}^n f_{X_i}(x_i) = \prod_{i=1}^n \text{cm}(x_i | \pi_i, \Sigma_i, w) = \prod_{i=1}^n \left( \sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c) \right)$$

Taking the logarithm gives us the log-likelihood:

$$\ell(\Theta) = \log L(\Theta) = \sum_{i=1}^n \log \left( \sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c) \right)$$

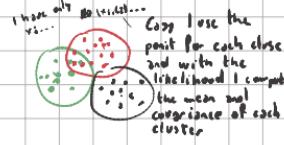
Using this function, we can now fit the parameters  $\pi, \Sigma, w$  to the data. This involves maximizing the log-likelihood  $\ell(\Theta)$  with respect to  $\Theta$ .

A key insight for working with GMMs is interpreting them as a marginal distribution of a joint distribution between observed data points and unobserved (latent) cluster assignments.

For each data point  $x_i$ :

- We have an associate (but unlabeled) cluster label  $c_i$
- The probability that  $x_i$  belongs to cluster  $c$  is  $P(c_i = c) = w_c$
- The conditional distribution of  $x_i$  given  $c_i = c$  is  $N(x_i | \mu_c, \Sigma_c)$

Visually, imagine our complex data distribution as composed of several simpler Gaussian "sub-populations" or clusters:



- Each cluster follows its own Gaussian distribution
- Points are probabilistically assigned to one of these clusters
- The overall distribution is the weighted combination of these cluster-specific distributions.

The problem is that we don't know which data point belongs to which cluster. If we knew the cluster assignments (so don't know nothing, even the number K of clusters)

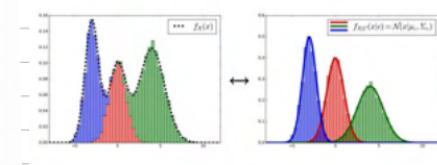
- ① We could divide the data into K separate groups
- ② For each group we could easily estimate the mean and the covariance using standard ML techniques for Gaussians.
- ③ The mixture weights would simply be the proportion of points in each cluster. (if we know what cluster belongs each point w will be 0 or 1)

If we knew the component responsible for each sample (i.e. its cluster label), we could estimate the parameters of each Gaussian by ML from the points of each cluster

Unfortunately, in general the clusters are **unknown**

We treat cluster membership as an **unobserved (latent)** random variables<sup>1</sup>

Intuitively, we want to estimate both cluster assignments and model parameters as to maximize the *marginal* distribution of the data



<sup>1</sup>Note that the model is **not identifiable**: for example, exchanging any two components results in the same marginal likelihood

$$\left. \begin{array}{l} w_1 \in \Sigma_C \\ w_2 \in \Sigma_C \\ w_3 \in \Sigma_C \end{array} \right\} \quad \stackrel{\text{!}}{=} \quad \left. \begin{array}{l} w_1 \in \Sigma_C \\ w_2 \in \Sigma_C \\ w_3 \in \Sigma_C \end{array} \right\}$$

Let's consider a set of GM parameters  $\Theta = (\pi, \mu, \Sigma)$ , we said that we don't know the cluster assignments, so we suppose that we know that.  
For any sample  $x_i$  and component  $c$ , we can define a joint density

$$f_{X_i, C_i}(x_i, c) = w_c N(x_i | \mu_c, \Sigma_c)$$

This represents the probability of both observing data point  $x_i$  AND having it come from component  $c$ . Think of this as answering the question: "What's the probability that we see this specific data point and it was generated by this specific component?"

From the joint probability using the Bayes theorem we can compute the posterior probabilities or responsibilities:

$$\gamma_{c,i} = P(C_i | c / X_i = x_i) = \frac{f_{X_i, C_i}(x_i, c)}{f_{X_i}(x_i)}$$

$$= \frac{w_c N(x_i | \mu_c, \Sigma_c)}{\sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c)} \text{ [our GM density]}$$

This tells us, given that we observed data point  $x_i$ , what's the probability it came from component  $c$ ?

Now, as a first approximation, we can make hard assignments by giving each point to the component with the highest posterior probability:

$$c_i^* = \arg \max_c P(C_i = c | x_i = x_i)$$

So even though we're uncertain about cluster membership, let's pretend each point definitely belongs to its most likely cluster" (this is the absolute truth)

Given that we can then estimate by ML the new GM parameters  $\Theta^{new} = (\boldsymbol{\mu}^{new}, \boldsymbol{\Sigma}^{new}, \boldsymbol{w}^{new})$

The log-likelihood becomes:

$$l(\Theta) = \prod f_{x_i | c_i}(x_i | c_i)$$

$$l(\Theta) = \sum_{i=1}^n \log f_{x_i | c_i}(x_i | c_i) + \sum_{i=1}^n \log P(C_i = c_i)$$

$$= \sum_{i=1}^n [\log N(x_i | \mu_{c_i}, \Sigma_{c_i})] + \sum_{i=1}^n [\log w_{c_i}]$$

DIFFERENTIA. DISPERSION  
MVE  
W PROB. PROPORTION  
CAN ADD UP  
IN OBTAINING SUM  
NEAR 0 VARIANCE  
A POINTS (OUR POSTIO.  
TOMMING) AS IN THE MVE

The first term depends only on the Gaussian parameters ( $\boldsymbol{\mu}$  and  $\boldsymbol{\Sigma}$ ), while the second term depends only on the mixture weights ( $w$ ). They don't interact with each other mathematically, which means we can optimize them completely independently.

$$l(\Theta) = l_N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) + l(w)$$

Let's focus on the first term:

$$l_N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{i=1}^n [\log N(x_i | \mu_{c_i}, \Sigma_{c_i})]$$

corresponds to the log-likelihood of a (multivariate) Gaussian classification model, where the class labels are assumed to be the estimated  $c_i^*$ .

$$l_N(\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \sum_{c=1}^C \sum_{i: c_i^* = c} \log N(x_i | \mu_{c_i}, \Sigma_{c_i})$$

This is exactly the same as if I had  $K$  separate datasets, one for each cluster, and I was fitting a single Gaussian to each dataset independently.

Let  $N_c$  be the number of samples for which  $c_i^* = c$ . The solution for  $\mu_c$  and  $\Sigma_c$  is thus [the maximum likelihood estimate]

$$\mu_c^* = \frac{1}{N_c} \sum_{i|c_i^*=c} x_i, \quad \Sigma_c^* = \frac{1}{N_c} \sum_{i|c_i^*=c} (x_i - \mu_c^*)(x_i - \mu_c^*)^T$$

$$\begin{aligned} c_1^* &= 2 \\ c_2^* &= 2 \\ c_3^* &= 2 \\ c_4^* &= 3 \end{aligned}$$

$\sum_{c=1}^4 N_c \log w_c$   
"log w + log w + log w + ..."

For the second term, we have  $L(w) = \sum_{i=1}^n [\log w_{c_i^*}] = \sum_{c=1}^4 \sum_{i|c_i^*=c} \log w_c$

This is the log-likelihood of a categorical distribution with  $n$  categories, where we observed  $N_c$  occurrences of category  $c$ . The maximum likelihood estimate for a categorical distribution is simply the relative frequency

$$w_c^* = \frac{N_c}{\sum_{c=1}^4 N_c}$$

Total "number of samples"

### What is a Categorical Distribution? Recap

The categorical distribution is one of the simplest probability distributions you can imagine. It describes the outcome of a single trial where exactly one of  $K$  possible outcomes can occur, and each outcome has its own probability.

Think of it as the mathematical way to describe any situation where you're choosing one option from several mutually exclusive possibilities. The classic example is rolling a die you have 6 possible outcomes (1, 2, 3, 4, 5, 6), and exactly one will happen on each roll.

### The Mathematical Setup

For a categorical distribution with  $K$  categories, we need  $K$  parameters:  $p_1, p_2, \dots, p_K$ , where  $p_i$  represents the probability of outcome  $i$  occurring. These probabilities must satisfy two important constraints:

First, each probability must be non-negative:  $p_i \geq 0$  for all  $i$ .

Second, since exactly one outcome must occur, the probabilities must sum to one:  $p_1 + p_2 + \dots + p_K = 1$ .

When we observe a single outcome, say outcome  $j$ , the probability of this observation is simply  $p_j$ . This seems almost trivially simple, but it becomes more interesting when we have multiple observations.

### Multiple Observations and Maximum Likelihood

Now suppose we conduct  $n$  independent trials and observe the following counts: we see outcome 1 exactly  $N_1$  times, outcome 2 exactly  $N_2$  times, and so on, where  $N_1 + N_2 + \dots + N_K = n$ .

The likelihood of observing this particular sequence of outcomes is the product of all individual probabilities. Since we observed outcome 1 exactly  $N_1$  times, outcome 2 exactly  $N_2$  times, and so forth, the likelihood becomes:

$$L(p_1, p_2, \dots, p_K) = p_1^{N_1} p_2^{N_2} \dots p_K^{N_K}$$

Taking the logarithm gives us the log-likelihood:

$$\ell(p_1, p_2, \dots, p_K) = N_1 \log p_1 + N_2 \log p_2 + \dots + N_K \log p_K$$

### The Maximum Likelihood Solution

To find the parameters that maximize this log-likelihood, we need to use calculus, but we must respect the constraint that the probabilities sum to one. Using the method of Lagrange multipliers (which I can explain separately if you'd like), the solution turns out to be beautifully simple:

$$\hat{p}_i = \frac{N_i}{n}$$

In other words, the maximum likelihood estimate for the probability of outcome  $i$  is simply the relative frequency with which we observed outcome  $i$ . This matches our intuitive understanding perfectly.

### Connecting Back to GMM Weights

Now let's see how this connects to our Gaussian Mixture Model context. In our GMM, each data point is assigned to exactly one cluster. We can think of this assignment process as a categorical distribution where:

The  $K$  categories correspond to the  $K$  mixture components. The probability parameters correspond to the mixture weights  $w_1, w_2, \dots, w_K$ . When we observe that  $N$  data points are assigned to component  $c$ , we're essentially observing  $N_c$  occurrences of "category  $c$ " in our categorical distribution.

Following the maximum likelihood principle for categorical distributions, our estimate for the mixture weight  $w_c$  becomes:

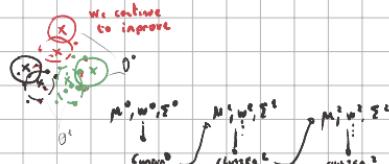
$$w_c^* = \frac{N_c}{n}$$

This is exactly the formula we saw in the GMM parameter estimation! The mixture weights are simply the maximum likelihood estimates for a categorical distribution where the "outcomes" are cluster assignments.

### Example

Hard assignment

$$\begin{aligned} \text{We have } 0 \text{ (initials)} & \quad \text{No hard fit} \quad \text{We said} \\ w_1^0, w_2^0, \dots, w_K^0 & \quad \text{before with} \\ \{x_1, x_2, \dots, x_n\} & \quad \text{hard assignment} \\ w_1^0, w_2^0, \dots, w_K^0 & \quad \text{so } c_i^* = \arg \max P(c|x_i) \\ & \quad \text{if } P(c|x_i) > P(c'|x_i) \quad \text{then } c_i^* = c \\ & \quad \text{if } P(c|x_i) \leq P(c'|x_i) \quad \text{then } c_i^* = c' \\ (x_i, c_i^*) & \quad \text{if } P(c|x_i) > P(c'|x_i) \quad \text{then } c_i^* = c \\ \text{In total no move} & \quad \text{if } P(c|x_i) \leq P(c'|x_i) \quad \text{then } c_i^* = c' \\ \text{Dataset } X, C \text{ known} & \quad \text{if } P(c|x_i) > P(c'|x_i) \quad \text{then } c_i^* = c \\ \text{1 parameter} & \quad \text{if } P(c|x_i) \leq P(c'|x_i) \quad \text{then } c_i^* = c' \\ \left[ \begin{array}{c} w_1^0, \mu_1^0, \Sigma_1^0 \\ w_2^0, \mu_2^0, \Sigma_2^0 \\ \vdots \\ w_K^0, \mu_K^0, \Sigma_K^0 \end{array} \right] & \quad \text{known} \\ \text{known cluster} & \quad \text{cluster} \end{aligned}$$



## The Chicken-and-Egg Problem Visualized

Imagine you're looking at a dataset of people's heights and weights, and you suspect there are three groups hidden in your data - perhaps children, women, and men. But you don't know who belongs to which group! This creates our fundamental problem:

To figure out the parameters of each group (like average height and weight), you need to know who belongs to each group. But to assign people to groups, you need to know what each group looks like. It's like trying to sort mail without knowing the addresses, but needing to know the addresses to sort the mail.

### Breaking the Deadlock with Iteration

The brilliant insight is that we can break this deadlock by making educated guesses and gradually improving them. Let me walk you through exactly how this works with a concrete example.

#### Step 1: Make an Initial Guess

Let's say we start by randomly guessing that our three groups have these characteristics:

Group 1 (children): average height 120cm, average weight 30kg

Group 2 (women): average height 160cm, average weight 60kg

Group 3 (men): average height 180cm, average weight 80kg

These initial guesses are probably wrong, but that's okay - we'll improve them.

#### Step 2: Compute Responsibilities

Now, for each person in our dataset, we ask: "Given their actual height and weight, which group are they most likely to belong to?"

For example, if we see someone who is 175cm tall and weighs 75kg:

Probability they're from Group 1 (children): very low (they're too tall and heavy)

Probability they're from Group 2 (women): moderate (possible, but on the tall/heavy side)

Probability they're from Group 3 (men): high (fits well with our current guess)

We compute these probabilities mathematically using our current parameter estimates.

#### Step 3: Make Hard Assignments

We then assign each person to their most likely group. Our 175cm, 75kg person gets assigned to Group 3 (men) because that had the highest probability.

#### Step 4: Update Parameters ~~our assignment was based on the probabilities so it wasn't correct for force.~~

Now here's the key insight: we use these assignments to compute better parameter estimates. We look at all the people we assigned to Group 3 and compute their actual average height and weight.

Maybe it turns out to be 177cm and 78kg instead of our initial guess of 180cm and 80kg.

We do this for all three groups, getting updated estimates based on who we think belongs to each group.

#### Step 5: Repeat and Improve

We then go back to Step 2, but now with our improved parameter estimates. This time, when we compute responsibilities for our 175cm, 75kg person, the probabilities might be different because our group characteristics have changed. Maybe now they're more likely to belong to Group 2 than we initially thought.

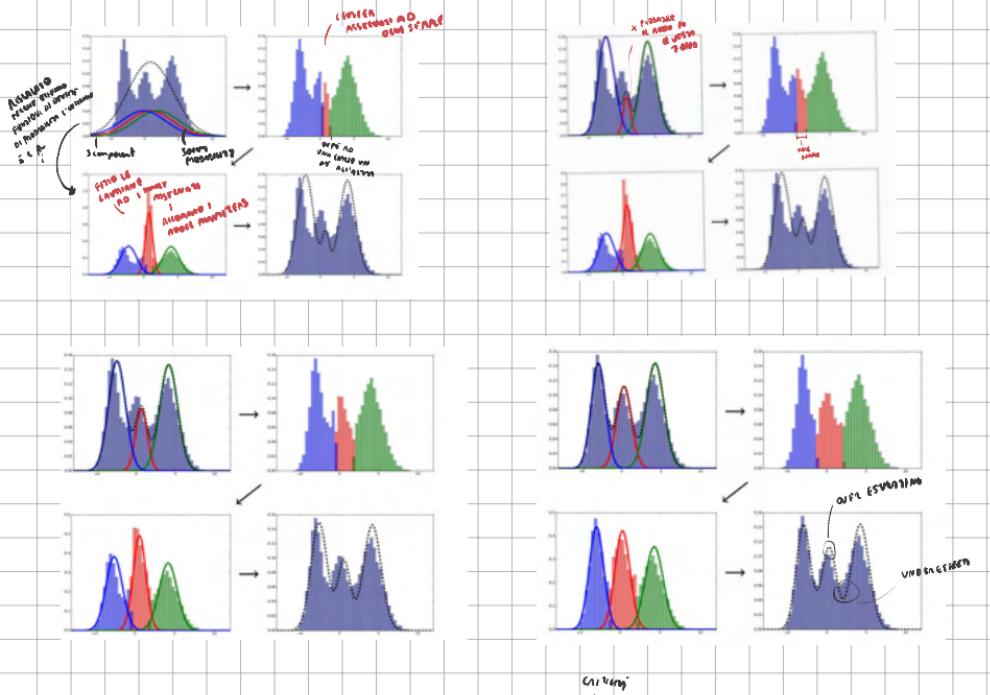
### Why This Creates Progress

Each iteration improves our estimates in a specific way. The responsibilities in Step 2 are computed using the best information we currently have about each group. The parameter updates in Step 4 are the mathematically optimal estimates given our current assignments. So each step is making the best possible decision given what we know at that moment.

The beautiful mathematical property is that each iteration increases the likelihood of our data under the model. Think of likelihood as a measure of "how well does our current model explain what we actually observed?" Each iteration makes our model explain the data better.

### Why It Eventually Stops

The process converges because there are only so many ways to improve. Eventually, the assignments become stable - people stop switching between groups, and the parameter estimates stop changing significantly. At this point, we've found a local optimum where our model does a good job of explaining the data structure.



The hard assignment approach has a critical flaw that became apparent when we examined what happens at cluster boundaries. Consider a data point that lies exactly between two cluster centers. When we compute the posterior probability, we might get something like

- $P(C_1=1|x_i=x_i) = 0.52$  (52% of belonging to  $C_1$ )
- $P(C_1=2|x_i=x_i) = 0.48$  (48% of belonging to  $C_2$ )

With hard assignments we would definitively assign this point to cluster 1 since  $0.52 > 0.48$ . But think about how cruel this approximation is! We're treating a point that's almost equally likely to belong to either cluster as if it definitely belongs to just one cluster. This throws away valuable information about the uncertainty in our assignments.

This becomes even more problematic when posterior probabilities are more evenly distributed, such as  $P(C_1=1|x_i=x_i) = 0.4$ ,  $P(C_1=2|x_i=x_i) = 0.35$ ,  $P(C_1=3|x_i=x_i) = 0.35$ . We will assign the point to cluster 1 but is actually quite uncertain which cluster generated it.

(X) ... (green circle)

So the hard assignment algorithm is not actually maximizing the likelihood of our observed data. We are making approximation that lead us away from the optimal solution.

So before we develop a full solution to estimate a local maximum to the likelihood, let's see what happens when we make some simplifying assumptions.

We still consider hard assignments,

We constraint our GM so that:

- All covariance matrices are identity matrices.  $\Sigma_c = I$
- All mixture weights are equal:  $w_c = \frac{1}{K}$

In this case cluster assignment corresponds to the rule

$$c_i^* = \arg \max_c P(c_i = c | x_i = x_i) = \arg \min_c \|x_i - \mu_c\|^2$$

So the posterior probabilities  $P(c_i = c | x_i = x_i)$  becomes proportional to the Gaussian density.

Our algorithm becomes

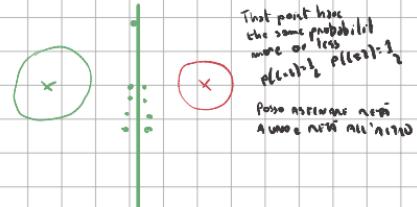
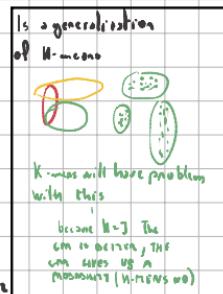
- Compute the component or cluster  $c_i^*$  whose centroid  $\mu_{c_i^*}$  is closest to our point and assign  $x_i$  to that cluster
- Re-estimate the cluster centroids from the given points, and iterate until convergence

This is the K-Means clustering algorithm

GMMs can also be applied to clustering tasks as a generalization of K-Means

The algorithm we considered can be extended to handle soft assignments

We will see that a point is not completely associated to a single Gaussian component, but contributes to the estimation of different components according to its cluster (component) posterior probability



SE NO NO SOLO UNO? [OPORTUNIDAD]

Lo splitton  
2 nosotro nosotro  
uno o otro nulo  
una classe  
1  
LOST USO DE UNIC  
MUCHAS PARAMETR  
DE UN CLUSTRO DE REP  
DE UN CLUSTRO

The key insight for improving our algorithm is to embrace uncertainty rather than eliminate it. Instead of forcing each point to belong to exactly one cluster, we can allow points to contribute to multiple clusters, according to its posterior probability.

We consider the full log-likelihood of our data

$$\sum_{i=1}^n \log f_{x_i}(x_i | \theta) = \sum_{i=1}^n \log \left( \sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c) \right)$$

CONSIDERA TANTO 1 CLUSTRO  
MAS ALMACENAR CLUSTROS  
APROXIMADAMENTE DIFERENTES

This expression capture the probability of observing our entire dataset under the current parameter estimates. Each term in the outer sum represents the log probability of one data point, where that probability comes from the weighted mixture of all  $N$  Gaussian components.

Now to find the optimal parameter we need to take the derivative with respect to each parameter and set it equal to zero.

Let's take the gradient with respect to  $\mu_c$

$$\frac{\partial L}{\partial \mu_c} = \sum_{i=1}^n \frac{\partial}{\partial \mu_c} \log \left( \sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c) \right) =$$

$$\sum_{i=1}^n \frac{1}{\sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c)} \frac{\partial}{\partial \mu_c} \left( \sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c) \right) =$$

only the  $c$ th term in the inner sum depends on  $\mu_c$ ,  $w_c \neq 0$

$$= \frac{\partial}{\partial \mu_c} \left( \sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c) \right) = w_c \frac{\partial}{\partial \mu_c} N(x_i | \mu_c, \Sigma_c) = w_c \frac{\partial}{\partial \mu_c} N(x_i | \mu_c, \Sigma_c)$$

$$N(x_i | \mu_c, \Sigma_c) = \frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma_c|^{\frac{1}{2}}} e^{-\frac{1}{2} (x_i - \mu_c)^T \Sigma_c^{-1} (x_i - \mu_c)}$$

$$\frac{\partial}{\partial \mu_c} N(x_i | \mu_c, \Sigma_c) = -N(x_i | \mu_c, \Sigma_c) \Sigma_c^{-1} (x_i - \mu_c)$$

The derivative of the quadratic form  $(x_i - \mu_c)^T \Sigma_c^{-1} (x_i - \mu_c)$  with respect to  $\mu_c$  is  $-2 \Sigma_c^{-1} (x_i - \mu_c)$

Putting all together

responsibility  
formula

$$\frac{\partial L}{\partial \mu_c} = - \sum_{i=1}^n \frac{w_c N(x_i | \mu_c, \Sigma_c)}{\sum_{c=1}^C w_c N(x_i | \mu_c, \Sigma_c)} \Sigma_c^{-1} (x_i - \mu_c) = - \sum_{i=1}^n \gamma_{ci} \Sigma_c^{-1} (x_i - \mu_c)$$

$$\frac{\partial L}{\partial \mu_c} = 0 \quad \sum_{i=1}^n \gamma_{ci} \Sigma_c^{-1} (x_i - \mu_c) = 0 \Rightarrow \Sigma_c^{-1} \sum_{i=1}^n \gamma_{ci} x_i = \Sigma_c^{-1} \mu_c \sum_{i=1}^n \gamma_{ci} =$$

$$\left( \mu_c - f(\mu_c) \right) = 0 \Rightarrow \mu_c = \frac{\sum_{i=1}^n \sum_{j=1}^n \gamma_{ci} x_i}{\sum_{i=1}^n \sum_{j=1}^n \gamma_{ci}} = \frac{\sum_{i=1}^n \gamma_{ci} x_i}{\sum_{i=1}^n \gamma_{ci}} \quad (3)$$

The optional mean for component  $c$  is simply the weighted average of all data points, where each point is weighted by its responsibility to that component. Points that are very likely to belong to component  $c$  ( $high \gamma_{ci}$ ) contribute strongly to determining that component's center, while points that are unlikely to belong to component  $c$  ( $low \gamma_{ci}$ ) have minimal influence.

This makes perfect intuitive sense. If we were trying to find the center of a cluster and we knew exactly which points belonged to it, we would simply average those points. Here, we're doing the same thing, but accounting for our uncertainty about cluster membership by using weighted averages instead of simple averages.

Notice that the responsibilities  $\gamma_{c,i}$  depend on  $\mu_c$ . If we knew the responsibilities we could compute  $\mu_c$  as in (3)

We can interpret (3) as a **weighted** empirical mean. The weight of each sample is the corresponding **responsibility**.

The terms

$$N_c = \sum_{i=1}^N \gamma_{c,i}$$

and

$$\mathbf{F}_c = \sum_{i=1}^N \gamma_{c,i} \mathbf{x}_i$$

are also called **zero and first order statistics**

Note that we are **summing over all samples** in  $\mathcal{D}$

We can adopt a similar strategy for the covariance matrix, obtaining

$$\Sigma_c = \frac{1}{N_c} \sum_i \gamma_{c,i} (\mathbf{x}_i - \mu_c) (\mathbf{x}_i - \mu_c)^T = \frac{1}{N_c} \sum_i \gamma_{c,i} \mathbf{x}_i \mathbf{x}_i^T - \mu_c \mu_c^T$$

The terms

$$S_c = \sum_i \gamma_{c,i} \mathbf{x}_i \mathbf{x}_i^T \quad Z_c = \frac{S_c}{N_c} = \mu_c \mu_c^T$$

are also called **second order statistics**

The weights can be re-estimated as

$$w_c = \frac{N_c}{N}$$

where  $N$  is the number of samples  $N = \sum_{c=1}^K N_c$

## Covariance Formula Intuitively

Think about what this formula is telling us. If we knew exactly which points belonged to component  $c$ , we would compute the sample covariance as the average of the outer products  $(\mathbf{x}_i - \mu_c)(\mathbf{x}_i - \mu_c)^T$  for all points in that component. Here, we're doing exactly the same thing, but we're weighting each point's contribution by how likely it is to belong to component  $c$ .

This makes perfect sense. A point that's very likely to belong to component  $c$  (high  $\gamma_{c,i}$ ) should strongly influence our estimate of that component's covariance structure. A point that's unlikely to belong to component  $c$  (low  $\gamma_{c,i}$ ) should have minimal impact on that component's covariance estimate.

## Weight Formula Intuitively

This result is remarkably intuitive. The weight for component  $c$  is simply the proportion of the total "responsibility mass" that component  $c$  accounts for. Since  $N_c$  represents the effective number of points assigned to component  $c$ , and  $n$  is the total number of points,  $w_c = N_c/n$  tells us what fraction of the data is effectively explained by component  $c$ .

This makes perfect sense from a probabilistic standpoint. If component  $c$  is responsible for explaining a large portion of the data points (high  $N_c$ ), it should have a correspondingly large weight in our mixture. If component  $c$  explains very few data points (low  $N_c$ ), its weight should be small.

$N_c$  represent the effective number of point assigned to component  $c$ .

$\mathbf{F}_c$  represent the weighted sum of points assigned to component  $c$ .



I take the mean  
of the fraction  
that have been  
assigned

$\bar{\mathbf{x}}$  is some  
of the LF  
fraction and  
can be used  
as sample  
assignment

[On the slide there is an example]

**Note:** We didn't find a solution, we still need a way to compute  $y_{c,i} = \frac{w_c N(x_i | m_c, \Sigma_c)}{\sum_c w_c N(x_i | m_c, \Sigma_c)}$ . so is a  $f(m_c, \Sigma_c)$ , so we need initial parameter  $\Theta$   
(This not stop us we will have a starting point  $w_0 \neq 0$  to compute  $m_0$ )

Think back to our own log-likelihood function (DISTRIBUTED PENALTY FUNCTIONS + PENALTY AGAINST INCONSISTENCIES OR BIAS IN LOG-LIKELIHOOD)  
 $l(\Theta) = \sum_{i=1}^n \log \left( \sum_{c=1}^C w_c N(x_i | m_c, \Sigma_c) \right)$

The fundamental problem lies in that inner summation inside the logarithm. When we try to derive it to find maximum, we end up with expressions that involve ratios of complicated sums, leading to equations that have no closed-form solutions. (it's not solvable using a finite number of standard mathematical operations (addition, subtraction, multiplication...))

Contrast this with what happens when we know the cluster assignments. if we knew that point  $x_i$  came from cluster  $c_i$ , then our log-likelihood would become beautifully simple:

$$l(\Theta) = \sum_{i=1}^n [\log N(x_i | m_{c_i}, \Sigma_{c_i}) + \log w_{c_i}]$$

As we said we don't have this cluster assignments - they are hidden from us

**Note:** What does it mean Closed-Form Solution?

What Does "Closed-Form Solution" Mean?

is a formula or algorithm that gives you the exact answer using a finite number of standard arithmetic operations, exponentiation, logarithms, trigonometric functions, etc.

Examples of closed-form solutions:

- The solution to  $a + b = c$  is  $b = c - a$  (simple algebra)
- The solution to  $x^2 = c$  is  $x = \sqrt{c}$  (square root)
- The solution to  $x^3 = c$  is  $x = \sqrt[3]{c}$  (cube root)
- The solution to  $x^4 = c$  is  $x = \sqrt[4]{c}$  (fourth root)

Examples where NO closed-form solution exists:

- The equation  $x^2 + x + 1 = 0$  is a general formula for polynomials degree 2 or higher
- Finding where  $x^2 = x$  intersects (must use numerical methods)
- Most optimization problems are machine learning

Why is "Log Inside a Log" So Problematic?

Let me show you exactly why this mathematical structure causes such difficulties

Simple case to try to solve:

$$\log(\log(x)) = 0$$

Taking the exponent of both sides:  $\log(\cdot) = 0 \rightarrow \cdot = 1$

Setting  $x = e^{1/2}$

This has a nice closed-form solution

Complex case (the CIR problem):

$$\log(\log(\log(x))) = \log(\log(\log(e))) = 0$$

Where even  $\log(\cdot)$  is linear in  $\log(\cdot)$ , taking derivatives can be even more difficult. When we try to take the derivative to the right, we get expressions like

$$\frac{d}{dx} \log(\log(\log(x))) = \log(\log(x))^{-1} \cdot \frac{d}{dx} \log(\log(x)) = \frac{\log(\log(x))^{-1}}{\log(x)^{-1}} \cdot \frac{d}{dx} \log(x) = \frac{1}{\log(x)^2}$$

This is really bad because the denominator is zero when  $x = 1$ . This is a singularity that gives us an error.

So we can't take the derivative with standard algorithms.

In fact just for example

so compute that graph

The EM is an iterative procedure suited for the ML estimation of the parameters of complex likelihoods<sup>2</sup>  $f_X(x|\theta)$  that can be expressed through marginalization of joint likelihoods  $f_{X,H}(x, h|\theta)$ :

$$f_X(x) = \int f_{X,H}(x, h) dh = \int f_{X|H}(x|h) f_H(h) dh$$

**NOTE:** Here and in the following we do not make any assumption on what  $X$  represents. We simply assume that it's a random variable or a random vector, for which we **observe a value**  $x$ . We shall see later that, for our GMM estimation task,  $X$  represents the set of random vectors that describe the feature vectors in our dataset  $\mathcal{D}$

The EM algorithm applies to any situation where we have observed data  $x$  and hidden (latent) variables  $H$ , and where the joint likelihood  $f_{X,H}(x, h|\Theta)$  is much easier to work with than the marginal likelihood  $f_X(x|\Theta)$

The derivation holds for both continuous and discrete latent variable replacing integrals

with sums.

$H$  represents a **latent (or hidden) random variable (or vector)** i.e.  
a R.V. whose **value has not been observed**, i.e., is unknown

As we shall see, the EM transforms the maximization of a log-likelihood  $\log f_X(x|\theta)$  into a sequence of optimizations of expectations of the joint log-likelihood  $\log f_{X,H}(x, h|\theta)$

Let's consider again the marginal log-likelihood

$$\ell(\theta) = \log f_X(x|\theta) = \log \frac{f_{X,H}(x,h|\theta)}{f_{H|x}(h|x,\theta)}$$

once theorem

Given a density  $Q(h)$  with the same support of  $f_{H|x}(h)$  ( $Q(h)$  and  $f_{H|x}(h)$  are both no zero on exactly the same domain), this prevent us to have  $Q(h)>0$  in regions where  $f_{H|x}(h)=0$ )  
we can rewrite the log-pdf as

$$\begin{aligned} \log f_X(x|\theta) &= \log f_X(x|\theta) \int Q(h) dh = \\ &= \int Q(h) \log f_{X,H}(x,h|\theta) dh = \int Q(h) \cdot \log \frac{f_{X,H}(x,h|\theta)}{f_{H|x}(h|x,\theta)} dh = \\ &= \int Q(h) \log f_{X,H}(x,h|\theta) dh - \int Q(h) \log f_{H|x}(h|x,\theta) dh = \\ &= \int Q(h) \log f_{X,H}(x,h|\theta) dh - \left( \int Q(h) \log f_{H|x}(h|x,\theta) dh \right) = \\ &\quad + \int Q(h) \log Q(h) dh - \int Q(h) \log Q(h) dh = \\ &= \int Q(h) \log \frac{f_{X,H}(x,h|\theta)}{Q(h)} dh - \int Q(h) \log \frac{f_{H|x}(h|x,\theta)}{Q(h)} dh \end{aligned}$$

[  $\int Q(h) dh = 1$ , is a probability distribution, so we can use it as a weighting function ]

In our case both  $Q(h)$  and  $f_{H|x}(h)$  represent the distribution over the cluster assignment  $(1, 2, \dots, n)$  so they naturally have the same support

The term:

$$D_h(Q(h) || f_{H|x}(h|x,\theta)) = - \int Q(h) \log \frac{f_{H|x}(h|x,\theta)}{Q(h)} dh = - E_Q \left[ \log \frac{f_{H|x}(h|x,\theta)}{Q(h)} \right]$$

This term is called Kullback-Leibler (KL) divergence (usually denoted simply as  $D(Q||f_{H|x})$ )

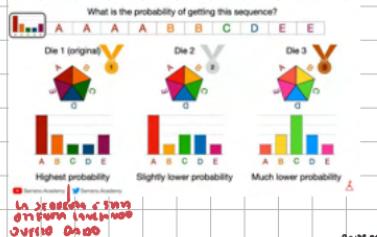
$D$  stand for divergencies standing that this term measure the difference between the two distributions "the distance".

The  $||$  bars mean "measured against", or "compared to", so when we write  $D_h(Q(h) || f_{H|x}(h|x,\theta))$  we are asking the question how far away is our approximation  $Q$  from the true target distribution  $f_{H|x}(h|x,\theta)$

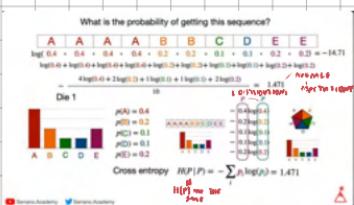
So how goal will be to minimize  $D_{KL}(Q||P_{\text{true}})$ , which means making our auxiliary distribution as close as possible to the posterior distribution  $P_{\text{true}}(b|x, \theta)$ . [now considerando il rango quantitativo questo è facile coniugare ci sarà una minimizzazione logaritmica]

BRO È UNA CURE E' UNA DIFFERENZA, STAI CONSIDERANDO FUNZIONI DI PROBABILITÀ O PROBABILITÀ IL CUI VALORE DI PROBABILITÀ EFFETTIVO È IL PUNTO INTELLIGIBILE, HAI MESSO IL COL DELL'ALLENAMENTO E SEMPRE MA LA PROPRIETÀ DI COL VERSO TI DICE?  $\log_2 - \log_3 - \log_4 =$  DIFFERENZA DA DISTINZIONI (IN PROBABILITÀ SONO L'INTELLIGIBILE). QUESTO (HE) È IL PESO (WEIGHT) AL PENSIERE DELLA DECISIONE)

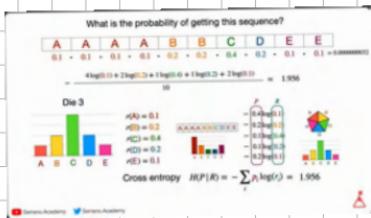
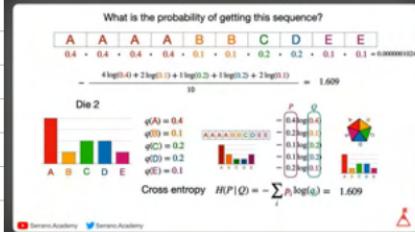
## Example



## LA JEFERIA DE LOS DIFUSORES DADO



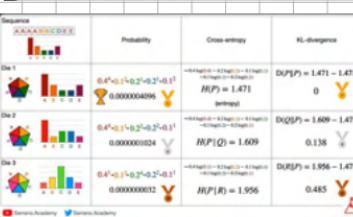
PROBABILITÀ DI OTTENERE  
LA STEMA SEQUENZA  
CON I DIVERSI  
DATI.



So KL divergencies is calculated as a difference of the cross-entropy is not symmetric (we can change the order of the terms)

$\delta_{(n)}$  è la distribuzione che esiste che venga tirata e che si sposa altra con i dati che abbiamo

In 1917 the 1st S.S. was formed.



What is cross-entropy?? How it connect with oligvergencis?

## **Il Concetto Fondamentale: La Sorpresa**

Prima di immergerti nelle formule matematiche, dobbiamo comprendere il concetto di sorpresa (o "surprise"), che è il cuore di tutto. La sorpresa misura quanto è inaspettato un evento. Pensa a quando lanci una moneta: se esce testa, non sei molto sorpreso perché sapevi che aveva il 50% di probabilità. Ma se lanci un dado e esce un sei specifico, sei più sorpreso perché aveva il 1/6 di probabilità.

La sorpresa deve avere proprietà molto specifiche che riflettano la nostra intuizione. Primo, eventi più rari dovrebbero causare più sorpresa. Secondo, eventi certi (con probabilità 1) non dovrebbero causare alcuna sorpresa. Terzo, e questo è cruciale, la sorpresa deve essere additiva per eventi indipendenti.

Cosa significa "additiva"? Se lanci tre dadi simultaneamente e prevedi correttamente tutti e tre i risultati, la sorpresa totale dovrebbe essere circa tre volte la sorpresa di prevedere un singolo dado. Ma qui nasce un problema matematico interessante: mentre le sorprese si sommano, le probabilità di eventi indipendenti si moltiplicano. La probabilità di ottenere tre risultati specifici è  $(1/6)^3 = (1/6) \times (1/6)$ .

### Perché il Logaritmo È Essenziale

Per risolvere questo problema matematico, abbiamo bisogno di una funzione che trasformi la moltiplicazione in addizione. Questa funzione magica è il logaritmo! La sorpresa di un evento con probabilità  $p$  è definita come  $-\log(p)$  o  $\log(1/p)$ .

Questa definizione è brillante perché soddisfa tutte le nostre esigenze. Primo,  $\log(p)$  è sempre positivo quando  $0 < p < 1$ , perché il logaritmo di un numero minore di 1 è negativo. Secondo, diminuisce quando  $p$  aumenta, raggiungendo zero quando  $p = 1$  (poiché  $\log(1) = 0$ ). Terzo, e più importante, trasforma la moltiplicazione in addizione:  $\log(1/p_1 \cdot p_2 \cdots p_n) = \log(1/p_1) + \log(1/p_2) + \cdots + \log(1/p_n)$ .



## L'Entropia: La Sorpresa Media

Ora che abbiamo definito la sorpresa per eventi singoli, possiamo estendere il concetto a intere distribuzioni di probabilità. L'**entropia** di una distribuzione  $P$  misura la sorpresa media che ci aspettiamo da quella distribuzione.

Matematicamente, l'entropia è:

$$H(P) = \sum_i P(x_i) \times (-\log P(x_i)) = -\sum_i P(x_i) \log P(x_i)$$

Questa formula ci dice di prendere ogni possibile risultato, moltiplicare la sua probabilità per la sua sorpresa, e sommare tutto. Il risultato ci dice quanto è "incerta" o "imprevedibile" la distribuzione. Una distribuzione uniforme (dove tutti gli esiti sono ugualmente probabili) ha alta entropia perché è molto imprevedibile. Una distribuzione dove un esito ha probabilità molto alta e gli altri molto bassa ha bassa entropia perché è molto prevedibile.

## La Divergenza KL: Il Ponte Tra i Concetti

Ora arriviamo alla connessione elegante che lega tutto insieme. La **divergenza di Kullback-Leibler** è definita come:

$$D_{KL}(P||Q) = H(P, Q) - H(P)$$

Questa relazione rivelà che la divergenza KL misura il costo informativo extra quando usiamo la distribuzione  $Q$  invece della distribuzione ottimale  $P$ . Matematicamente:

$$D_{KL}(P||Q) = \sum_i P(x_i) \log \frac{P(x_i)}{Q(x_i)} = -H(Q) + H(P, Q)$$

Il primo termine,  $-H(Q)$ , rappresenta il costo informativo minimo possibile quando usiamo lo schema di codifica ottimale basato sulla vera distribuzione  $P$ . Il secondo termine,  $H(P, Q)$ , rappresenta il costo informativo effettivo quando usiamo lo schema subottimale basato sulla distribuzione  $Q$ . La loro differenza ci dà la divergenza KL, che misura esattamente quanto costo extra sosteniamo usando la distribuzione sbagliata.

## Perché Spesso Usiamo Cross-Entropia Invece di Divergenza KL

Nella pratica, il codice per l'addestramento dei modelli spesso non usa la divergenza KL come funzione di errore da minimizzare, ma piuttosto la cross-entropia. Questo accade per una ragione molto pratica e elegante.

Guardando la definizione di divergenza KL, il termine entropia di  $P(H|P)$  non dipende dal modello  $Q$ . È una costante determinata dalla vera distribuzione  $P$ , deriva dall'incertezza intrinseca nei dati di addestramento e non può essere modificata modificando i parametri del modello.

Pertanto, qualsiasi modello  $Q$  che minimizza la divergenza KL minimizza anche, per definizione, la cross-entropia. L'entropia dei dati  $H|P$  non influenza su quale  $Q$  è ottimale; sposta solo il valore della divergenza KL di una quantità costante.

Poiché il calcolo esatto dell'entropia dei dati  $H|P$  da un numero finito di campioni richiede una stima che di solito non è necessaria per l'ottimizzazione, si evita di stimarla per risparmiare computazioni. In sintesi, minimizzare la cross-entropia è equivalente a minimizzare la divergenza KL rispetto al modello  $Q$ .

## La Cross-Entropia: Quando Usiamo il Modello Sbagliato

La **cross-entropia** estende il concetto di entropia a situazioni dove abbiamo due distribuzioni diverse. Immagina di avere la vera distribuzione  $P$ , ma stai usando una distribuzione modello  $Q$  per fare previsioni. La cross-entropia misura la sorpresa media che ti aspetti quando la realtà segue  $P$  ma tu credi che segua  $Q$ .

La formula è:

$$H(P, Q) = \sum_i P(x_i) \times (-\log Q(x_i)) = -\sum_i P(x_i) \log Q(x_i)$$

Nota la differenza sottile ma cruciale: stiamo ancora pesare con le probabilità vere  $P(x_i)$ , ma ora prendiamo il logaritmo delle probabilità del modello  $Q(x_i)$ . Questo cattura il costo di usare la distribuzione sbagliata per le nostre previsioni.

Pensa a un esempio pratico di previsioni meteorologiche. Supponi che la vera probabilità di pioggia domani sia 70%, ma il tuo modello prevede solo 30%. La cross-entropia misura quanto costo informativo extra paghi perché le tue previsioni non corrispondono alla realtà.

## La Connessione con l'Apprendimento Automatico

Questi concetti diventano particolarmente potenti quando comprendiamo il loro ruolo nell'ottimizzazione dell'apprendimento automatico. In molti problemi di machine learning, stiamo cercando di trovare una distribuzione modello  $Q$  che approssimi il più possibile una distribuzione target  $P$ . La divergenza KL fornisce una misura naturale di quanto bene la nostra approssimazione si comporta.

Nei problemi di classificazione, la loss della cross-entropia è una delle funzioni obiettivo più comuni. Quando minimizziamo la cross-entropia tra le probabilità predette delle classi e le etichette vere delle classi, stiamo essenzialmente minimizzando la divergenza KL tra le previsioni del nostro modello e la distribuzione vera.

Questa connessione spiega perché la cross-entropia funziona così bene come funzione di loss. Non stiamo solo minimizzando una misura di distanza arbitraria; stiamo minimizzando il costo teorico-informativo di usare il nostro modello invece del modello ottimale.

The term

$$L_h(Q(h), \theta) = \int Q(h) \log \frac{f_{x,h}(x|h|\theta)}{Q(h)} dh = \text{Each} [f_{x,h}(x|h|\theta)] + H(Q(h))$$

where  $H(Q(h))$  is the entropy of distribution  $Q(h)$ , provides a lower bound of the log-likelihood (again, the suffix  $h$  is usually omitted, but we keep it to remember we are integrating w.r.t.  $h$ )

$L_h(Q(h), \theta)$  is called evidence Lower Bound or ELBO

When we write

$$L_h(Q(h), \theta) = \int Q(h) \log \frac{f_{x,h}(x|h|\theta)}{Q(h)} dh$$

We can actually split this integral into two meaningful parts, using the logarithm property that

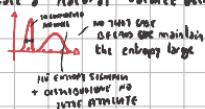
$$\log \left( \frac{A}{C} \right) = \log(A) - \log(C)$$

$$L_h(Q(h), \theta) = \int Q(h) \log f_{\text{mix}}(x_i | h, \theta) dh - \int Q(h) \log Q(h) dh$$

The first integral,  $\int Q(h) \log f_{\text{mix}}(x_i | h, \theta) dh$  represent the expected value of the joint log-likelihood under the distribution  $Q(h)$ . This is exactly what we would want to maximize, if we knew the hidden variable - it's expectation of our complete log-likelihood function.

The second integral  $-\int Q(h) \log Q(h) dh$  is the entropy  $H(Q(h))$  of our auxiliary distribution  $Q$ . Remember that entropy measures the uncertainty or randomness in a distribution. Higher entropy means the distribution  $Q$  is more spread out and uncertain about hidden variables.

So ELBO as we well see, we want to find parameters  $\theta$  that make the joint likelihood large (first term) but we also want to maintain high entropy in our auxiliary distribution (second term). This creates a natural balance between fitting the data well and not being overly confident about our hidden variable assignments.



Why ELBO provides a valid lower bound?

The key insight relies on the fundamental inequality  $\log z \leq z-1$  for any positive number  $z$ , with equality if and only if  $z$  equals one.  $[\log z + z-1 = z(z-1)]$



When we apply this inequality to our KL divergence, we set  $z = \frac{f_{\text{mix}}(h|x, \theta)}{Q(h)}$  and get:

$$\left[ -\log f_{\text{mix}}(h|x, \theta) \leq -\frac{\ln f_{\text{mix}}(h|x, \theta) + 1}{Q(h)} \right] \text{ with the equality holding if and only if } z=1 \quad Q(h) = f_{\text{mix}}(h|x, \theta)$$

This inequality holds pointwise for every value of  $h$  in the support of  $Q(h)$ . When we integrate both sides with respect to  $Q(h)$ ,

$$\begin{aligned} - \int Q(h) \log \frac{f_{\text{mix}}(h|x, \theta)}{Q(h)} dh &\geq - \int \frac{Q(h) \ln f_{\text{mix}}(h|x, \theta)}{Q(h)} dh + \int Q(h) dh \\ &\quad \left. \begin{aligned} &= \int f_{\text{mix}}(h|x, \theta) dh - \\ &= D_{\text{KL}}(Q(h) || f_{\text{mix}}(h|x, \theta)) \geq 1 - 0 \end{aligned} \right\} \end{aligned}$$

So, the KL divergencies is always non-negative. The equality conditions tell us also that  $D_h=0$  if and only if  $Q(h) = f_{H|X}(h|x, \theta)$  (AVER INTEGRAZIONE NELL'INTESA HEMPF, AVVOLGENDO I QUADRATI DELL'ESPRESSIONE DELLA DIVERGENZA CON UNA EQUAZIONE LOGARITMICA).

$$D_h(Q(h)||f_{H|X}(h|x, \theta)) = 0 \iff Q = f_{H|X} \text{ a. e.}$$

We have decomposed the log-likelihood as

$$\log f_X(x|\theta) = \mathcal{L}_h(Q(h), \theta) + D_h(Q(h)||f_{H|X}(h|x, \theta))$$

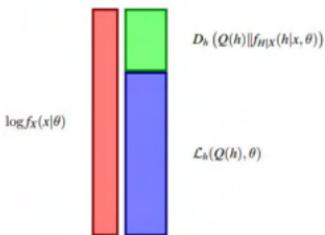
Notice that the left side of the equation does not depend on the choice of  $Q$  at all. Furthermore since we've proven that the KL divergence is always non-negative, we immediately know that

$$D_h(Q(h)||f_{H|X}(h|x, \theta)) \geq 0 \implies \mathcal{L}_h(Q(h), \theta) \leq \log f_X(x|\theta)$$

*È dunque se questo è minima somma di 2 quantità di cui è sempre positiva.  
 $\beta_1 = L + D$ , ponendo  $\geq 0$   
 perché si tratta  
 delle quantità*

So we have proved that ELBO provides a lower bound that we can optimize instead of the original intractable log-likelihood.

Decomposition of  $\log f_X(x|\theta) = \mathcal{L}_h(Q(h), \theta) + D_h(Q(h)||f_{H|X}(h|x, \theta))$



Our objective was to find the maximum log-likelihood we can do this directly maximizing the lower bound  $\mathcal{L}_h(Q(h), \theta)$ . We can make this bound arbitrarily tight by choosing  $Q$  appropriately.

The EM algorithm optimizes the log-likelihood by iteratively

- Maximizing the lower bound  $\mathcal{L}_h(Q(h), \theta)$  with respect to  $Q$ .
- Maximizing the lower bound  $\mathcal{L}_h(Q(h), \theta)$  with respect to  $\theta$ .

From an initial sets of parameters  $\theta_0$ :

- $Q_0 = \arg \max_Q \mathcal{L}_h(Q(h), \theta_0)$
- $\theta_1 = \arg \max_\theta \mathcal{L}_h(Q_0(h), \theta)$
- $Q_1 = \arg \max_Q \mathcal{L}_h(Q(h), \theta_1)$
- $\theta_2 = \arg \max_\theta \mathcal{L}_h(Q_1(h), \theta)$
- ...

Think of this as a coordinated dance where each partner takes turns leading while the other follows (it's exactly what we have done until now, when we have tried to compute directly the sum of our dataset  $x$  we are just giving the mathematical interpretations).

There are there are two step:

### The Expectation Step (E-step)

Given our current parameter estimates  $\theta^{(t)}$  (for example  $(\mu_1, \Sigma_1, \pi_1)$  in our case) we choose  $Q$  to make the lower bound as tight as possible. (we maximize our ELBO, setting it equal to its maximum possible value).

We have shown that

$$L_h(Q(h), \theta_t) \leq f_{\text{mix}}(h|x, \theta_t)$$

We can maximize  $L_h(Q(h), \theta_t)$  with respect to  $Q$  by simply selecting

$$Q(h) = f_{\text{mix}}(h|x, \theta_t)$$

We set our auxiliary distribution  $Q$  equal to the posterior distribution of the hidden variables given our current parameters.

But what this posterior actually represent? It answer the question: "Given my current model parameters and the data I've observed, what's the probabilities that each hidden variable take each possible value?" In our context, this becomes: "Given my current estimates of the cluster means, covariances, and weights, what's the probability that each datapoint belong to a cluster?"

This step is called expectation because we're computing expected values of the sufficient statistic under this posterior distribution.

We're calculating things like "the expected number of points assigned to cluster 1" and "the expected number of points assigned to cluster 2," where the expectations account for our uncertainty about cluster assignments.

it's not the effective number but the portion instead.

The KL divergence reaches its minimum value of zero when the two distributions are identical. Therefore, the optimal choice is always

$$Q_t(h) = f_{\text{mix}}(h|x, \theta_t)$$

This makes intuitive sense - we're setting our auxiliary distribution equal to what we currently believe the true posterior distribution of the hidden variables should be. This decrease the divergence between the two distribution

This step doesn't change the log-likelihood value at all, since we're only adjusting our auxiliary distribution  $Q$ , not the actual parameters  $\theta$ . But it sets us up perfectly for the M-step by ensuring that our lower bound is as tight as possible.

### M-Step (Maximization with respect to Parameters)

Now given the distribution  $Q(h)$  we can maximize the log-likelihood computing the new model parameter  $\theta$  (= actual maximize the log-likelihood following the new  $Q(h)$  distribution).

We fix our auxiliary distribution  $Q(h)$  and find the parameters  $\theta_{\text{new}}$  that maximize the lower bound. The optimization problem becomes:

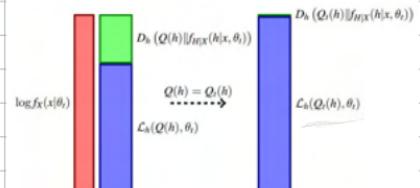
$$\theta_{\text{new}} = \arg \max \mathbb{E}_{Q(h)} [\log f_{\text{mix}}(x, h|\theta) dh] = \int f_{\text{mix}}(h|x, \theta_t) \log f_{\text{mix}}(h|x, \theta) dh$$

so  $\theta_{\text{new}}$  does not depend on  $Q$  but on  $Q_t$  that is fixed.

This expectation has a beautiful interpretation. We're finding the parameters

that maximize the expected complete data log-likelihood, where the expectation is take with respect to our current

Setting  $Q_t(h) = f_{\text{mix}}(h|x, \theta_t)$  does not change the log-likelihood  $\log f_X(x|\theta_t)$ , but reduces to zero the KL divergence:



best guess about the hidden variables.

The expectation  $E_{\{Q^t(h)\}}$  appears because we don't know for certain which cluster each point belongs to. Instead of using hard assignments ("this point definitely belongs to cluster 1"), we use soft assignments ("this point is 70% likely to belong to cluster 1, 30% to cluster 2").

The expectation weights each possibility by its probability:

If a point is 70% likely to be chocolate, it contributes 70% of its weight to estimating chocolate properties

If a point is 30% likely to be vanilla, it contributes 30% of its weight to estimating vanilla properties

The  $E_{Q^t(h)}$  operator weights each possible configuration of hidden variables by its probability under  $Q^t(h)$ . This creates a weighted average where more probable configurations have more influence on the parameter update.

Now we can also rewrite the problem as

$$E_{Q^t(h)} [\log f_{H|X}(h|x, \theta)] + E_{Q^t(h)} [\log f_H(h|\theta)]$$

Which corresponds exactly to what we used in EM derivation ( $S = \sum_{i=1}^n \gamma_{i,t} = Q^t(h)$ )

The first term represents the expected log-likelihood of the observed data given the hidden variables and parameters. The second term represents the expected log-likelihood of the hidden variables themselves given the parameters.

Since we are maximizing  $L_h(Q^t(h), \theta)$  we have:

$$L(\theta_t, \theta_{t+1}) \geq L(\theta_t, \theta_t)$$

But as we said earlier  $L_h(Q^t, \theta_{t+1})$  is the lower bound of  $L_h(Q^t, \theta_t)$ :  $L_h(Q^t, \theta_{t+1}) \leq L_h(Q^t, \theta_t) \leq L_h(Q^t, \theta_{t+1}) + D(Q^t, f_{H|X}(h|x, \theta_{t+1}))$

↳ therefore

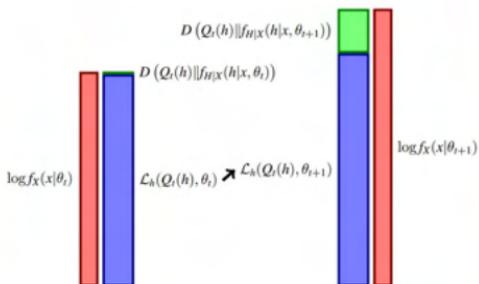
$$\log f_X(x|\theta_{t+1}) = L(Q^t(h), \theta_{t+1}) + D(Q^t, f_{H|X}(h|x, \theta_{t+1}))$$

Maximization with respect to  $\theta$  has increased both  $L_h$  and the KL-divergence, since, in general,  $Q_t(h) \neq f_{H|X}(h|x, \theta_{t+1})$  unless we reached convergence

We thus have that

$$\log f_X(x|\theta_{t+1}) \geq \log f_X(x|\theta_t)$$

Maximization w.r.t.  $\theta$  of  $\mathcal{L}_h(Q_t, \theta)$  increases the log-pdf  $\log f_X(x|\theta)$



The algorithm iterates between two steps

- **Expectation (E) step:** Compute the posterior distribution  $f_{H|X}(h|x, \theta_t)$  and compute the auxiliary function:

$$Q(\theta, \theta_t) = \mathbb{E}_{f_{H|X}(h|x, \theta_t)} [\log f_{X,H}(x, h|\theta)] \quad \text{start c19}$$

- **Maximization (M) step:** Maximize  $Q(\theta, \theta_t)$  w.r.t.  $\theta$  to obtain

$$\theta_{t+1} = \arg \max_{\theta} Q(\theta, \theta_t)$$

This creates the sequence  $\log f_X(x|\theta^0) \leq \log f_X(x|\theta^1) \leq \log f_X(x|\theta^2) \leq \dots$  guaranteeing that each iteration improves our objective function or at least leaves it unchanged.

The EM algorithm converges to a stationary point of the log-likelihood function under very general conditions. However, this stationary point might be a local maximum, a global ( $\theta^0$ ) maximum, or in rare cases, a saddle point. The algorithm cannot distinguish between these possibilities based solely on the convergence behavior.

This dependence on local optima makes initialization crucial for EM's success. Different starting points  $\theta^0$  can lead to different final solutions, and some of these solutions will be better than others. In practice, it's common to run EM multiple times with different random initializations and select the solution that achieves the highest final log-likelihood.

Think of the log-likelihood surface as a mountainous landscape with multiple peaks. EM is like a hiker who always walks uphill but can only see the immediate terrain around their current position. Depending on where the hiker starts, they might end up on a small local peak or discover the highest mountain in the region. Running multiple hiking expeditions from different starting points increases the chances of finding the tallest peak.

We have  $n$  hidden variables  $h = (c_1, c_2, \dots, c_n)$  where each  $c_i$  represents the cluster assignment for the  $i$ -th data point. Think of  $c_i$  as answering the question: "Which Gaussian component generated data point  $x_i$ ?"

Each cluster assignment  $c_i$  can take one of  $K$  possible values  $[1, 2, \dots, K]$  corresponding to the  $K$  components in our mixture model; if  $c_i = j$ , it means that data point  $x_i$  was generated by the  $j$ -th Gaussian component in our mixture.

The joint-likelihood for a single sample and its cluster assignment is:

$$f_{x_i, c_i}(x_i, c_i) = w_c N(x_i | \mu_c, \Sigma_c)$$

The term  $w_c$  represents the prior probability that any randomly selected point belongs to cluster  $c$ . ( $P(c_i = c)$ ) The term  $N(x_i | \mu_c, \Sigma_c) = f_{x_i, c_i}(x_i)$  represents the probability of observing data point  $x_i$ , given that it came from cluster  $c$ . Together they give the joint probability of observing both, the data point and its cluster assignment.

Think about this intuitively. If you randomly select a person from a mixed population of children and adults,  $w_c$  tells you the probability they're from group  $c$  (children or adults), and  $N(x_i | \mu_c, \Sigma_c)$  tells you the probability of observing their specific height and weight given their group membership.

We assume that samples are independent given the model parameters, so that we can express the joint log-likelihood for all the training set samples as:

$$\log f_{x_1, \dots, x_N, c_1, \dots, c_N}(x_1, \dots, x_N | c_1, \dots, c_N | \theta) = \sum_{i=1}^N \log f_{x_i, c_i}(x_i, c_i | \theta)$$

The EM algorithm requires computing the posterior for the hidden variables  $C_1, \dots, C_N | X_1, \dots, X_N, \theta$ .

Due to the independence assumptions, also the posterior distribution factorizes as

$$f_{C_1, \dots, C_N | X_1, \dots, X_N}(c_1, \dots, c_N | x_1, \dots, x_N, \theta) = \prod_{i=1}^N P(C_i = c_i | X_i = x_i, \theta)$$

The cluster assignment for one data point doesn't depend on the cluster assignments of other data points, given the model parameters.

We begin with the **E-Step**, that requires to compute the auxiliary function:

$$\begin{aligned} Q(\theta, \theta_t) &= E_{C_1, \dots, C_N | X_1, \dots, X_N, \theta_t} [\log f_{x_1, \dots, x_N, c_1, \dots, c_N}(x_1, \dots, x_N, c_1, \dots, c_N | \theta)] \\ &= \sum_{i=1}^N E_{C_i | x_i = x_i, \theta_t} [\log f_{x_i, c_i}(x_i, c_i | \theta)] \end{aligned}$$

The expectation here is taken over all possible cluster assignments for each data point, weighted by their posterior probabilities under the current parameter estimates. This is where the famous "responsibilities"  $y_{c,i} = P(C_i = c | X_i = x_i, \theta_t)$  come into play.

Think of  $y_{c,i}$  as answering the question: 'Given the current model parameters and the observed data point  $x_i$ , what's the probability that this point belongs to cluster  $c$ ?' These responsibilities capture our uncertainty about cluster membership in a probabilistic way.

For example, if  $y_{1,i} = 0.7$  and  $y_{2,i} = 0.3$  for some data point  $i$ , it means we believe there's a 70% chance this point belongs to cluster 1 and a 30% chance it belongs to cluster 2. This soft assignment allows us to work with uncertainty rather than making premature hard decisions about cluster membership.

**E-step:** Compute  $\gamma_{c,i} = P(C_i = c | X_i = \mathbf{x}_i, \boldsymbol{\theta}^t)$ . The auxiliary function is

$$\begin{aligned}\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}_t) &= \sum_{i=1}^N \sum_{c=1}^K P(C_i = c | X_i = \mathbf{x}_i, \boldsymbol{\theta}_t) \log f_{X_i, C_i}(\mathbf{x}_i, c | \boldsymbol{\theta}) \\ &= \sum_{i=1}^N \sum_{c=1}^K \gamma_{c,i} [\log \mathcal{N}(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c) + \gamma_{c,i} \log w_c] \\ &= \sum_{i=1}^N \sum_{c=1}^K \gamma_{c,i} \left( \frac{1}{2} \log |\boldsymbol{\Lambda}_c| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Lambda}_c (\mathbf{x}_i - \boldsymbol{\mu}_c) \right) \\ &\quad + \sum_{i=1}^N \sum_{c=1}^K \gamma_{c,i} \log w_c\end{aligned}$$

For each data point  $i$  and cluster  $c$ , we're weighting the log-probability terms by how likely the point is to belong to that cluster. Points that are really likely to belong to that cluster contribute strongly to that cluster's parameter estimation, while points that are unlikely to belong to a cluster contribute weakly.

**Note:**

$$\log(N(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)) = -\frac{1}{2} \log |\boldsymbol{\Sigma}_c| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Sigma}_c^{-1} (\mathbf{x}_i - \boldsymbol{\mu}_c)$$

The precision matrix  $\boldsymbol{\Lambda}_c = \boldsymbol{\Sigma}_c^{-1}$

$$\log(N(\mathbf{x}_i | \boldsymbol{\mu}_c, \boldsymbol{\Sigma}_c)) = \frac{1}{2} \log |\boldsymbol{\Lambda}_c| - \frac{1}{2} (\mathbf{x}_i - \boldsymbol{\mu}_c)^T \boldsymbol{\Lambda}_c (\mathbf{x}_i - \boldsymbol{\mu}_c) + \text{constant}$$

We find the parameters that maximize the auxiliary function

**M-step:** Maximize  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}_t)$  w.r.t.  $\boldsymbol{\theta} = (\mathbf{M}, \mathcal{S}, \mathbf{w})$ , subject to  
 $\sum_{k=1}^K w_k = 1$ :

$$\begin{aligned}\boldsymbol{\mu}_c^* &= \frac{\sum_i \gamma_{c,i} \mathbf{x}_i}{\sum_i \gamma_{c,i}} \\ \boldsymbol{\Sigma}_c^* &= \frac{\sum_i \gamma_{c,i} (\mathbf{x}_i - \boldsymbol{\mu}_c)(\mathbf{x}_i - \boldsymbol{\mu}_c)^T}{\sum_i \gamma_{c,i}} \\ w_c^* &= \frac{\sum_i \gamma_{c,i}}{\sum_i \sum_c \gamma_{c,i}}\end{aligned}$$

The new estimate of the parameters is  $\boldsymbol{\theta}_{t+1} = (\mathbf{M}_{t+1}, \mathcal{S}_{t+1}, \mathbf{w}_{t+1})$ :

$$\mathbf{M}_{t+1} = [\boldsymbol{\mu}_1^* \dots \boldsymbol{\mu}_K^*], \quad \mathcal{S}_{t+1} = [\boldsymbol{\Sigma}_1^* \dots \boldsymbol{\Sigma}_K^*], \quad \mathbf{w}_{t+1} = [w_1^* \dots w_K^*]$$

**Updating the Means:** The optimal mean for cluster c is:

$$\mu_c^* = \frac{\sum_i \gamma_{c,i} x_i}{\sum_i \gamma_{c,i}}$$

This is a weighted average of all data points, where each point is weighted by its responsibility to cluster c. Think about how intuitive this is: points that are very likely to belong to cluster c (high  $\gamma_{c,i}$ ) have a strong influence on where we place the cluster center, while points that are unlikely to belong to cluster c (low  $\gamma_{c,i}$ ) have minimal influence.

**Updating the Covariances:** The optimal covariance for cluster c is:

$$\Sigma_c^* = \frac{\sum_i \gamma_{c,i} (x_i - \mu_c)(x_i - \mu_c)^T}{\sum_i \gamma_{c,i}}$$

This is a weighted sample covariance matrix. Each data point contributes to the covariance estimate in proportion to how likely it is to belong to that cluster. Points that are probably in the cluster strongly influence the cluster's shape and orientation, while points that probably aren't in the cluster have minimal impact.

**Updating the Mixture Weights:** The optimal weight for cluster c is:

$$w_c^* = \frac{\sum_i \gamma_{c,i}}{\sum_i \sum_c \gamma_{c,i}} = \frac{\sum_i \gamma_{c,i}}{N}$$

This represents the proportion of the data that is effectively assigned to cluster c. The numerator  $\sum_i \gamma_{c,i}$  can be interpreted as the "effective number of points" assigned to cluster c. Since the responsibilities for each point sum to 1 across all clusters, the denominator equals N (the total number of data points).

The EM algorithm for GMMs creates a beautiful iterative refinement process. In each E-step, we update our beliefs about which cluster each point belongs to based on the current parameter estimates. In each M-step, we update our parameter estimates based on the current cluster assignment beliefs.

This creates a feedback loop where better cluster assignments lead to better parameter estimates, which in turn lead to better cluster assignments. The mathematical guarantee is that each iteration improves the log-likelihood, ensuring that we make steady progress toward a solution that at least locally optimizes our objective function.

The algorithm continues until convergence, which typically occurs when the parameter changes between iterations become negligibly small, or when the log-likelihood improvement falls below a specified threshold. At convergence, we have found parameter estimates that represent at least a local maximum of the likelihood function, giving us a GMM that provides a good probabilistic model of our data

## From mle to nn classification

Traditional mle classification makes a strong assumption that might not hold in practice. For each class  $c, i$ , it assumes that all data points follow a single Gaussian distribution:  $X_i | C=c \sim N(\mu_{c,i}, \Sigma_c)$ . This works well when each class has a roughly elliptical unimodal distribution, but real world data often violates this assumption dramatically.

MVG for classification:

- Fit a gaussian density to samples of each class

$$X_t | C_t = c \sim N(\mu_c, \Sigma_c) \quad \text{NOT ONE GAUSSIAN}$$

$$f_{X_t|C_t}(x_t|c) = N(x_t|\mu_c, \Sigma_c)$$

$$P(C_t = c | X_t = x_t) = \frac{N(x_t|\mu_c, \Sigma_c) P(C_t = c)}{\sum_c' N(x_t|\mu_{c'}, \Sigma_{c'}) P(C_t = c')}$$

This is where GMM classification becomes powerful. Instead of forcing each class to follow a single gaussian distribution, we allow each class to be modeled by a mixture of multiple gaussian components. Each component can capture different "sub-style" or "sub-pattern" within the class.

In GMM classification we model each class  $c$  using its own GMM with parameters  $(\pi_c, \mu_c, \Sigma_c)$

$$X_t | C_t = c \sim \text{GMM}(\pi_c, \mu_c, \Sigma_c)$$

The class conditional distribution:

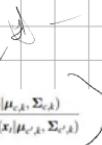
$$f_{X_t|C_t=c} = \sum_{k=1}^{K_c} w_{c,k} N(x_t|\mu_{c,k}, \Sigma_{c,k})$$

Each class can have a different number of components  $K_c$ .

Each class has its own complete set of parameters

The posterior probability:

$$P(C_t = c | X_t = x_t) = \frac{P(C_t = c) \sum_{k=1}^{K_c} w_{c,k} N(x_t|\mu_{c,k}, \Sigma_{c,k})}{\sum_{c'} P(C_t = c') \sum_{k=1}^{K_{c'}} w_{c',k} N(x_t|\mu_{c',k}, \Sigma_{c',k})}$$



This formula elegantly combines the prior probability of each class with the GMM-based likelihood. Each class's likelihood is computed as a weighted combination of multiple Gaussian components, allowing for much more flexible modeling of within-class variation.

## Open-set Classification

We can also exploit GMMs clustering capabilities for **open-set** multiclass classification

The difficulty in open-set classification consists in building a robust model for the **none-of-the-others** class

This class is usually very heterogeneous, and explicit modeling its sub-components requires labeled examples of its possible objects

We can partially alleviate the issue using a GMM.

We can assume that samples of known classes can be modeled by MVG distributions

We can collect a large set of **unlabeled** samples for the **none-of-the-others** class

We can model the samples of the **none-of-the-others** class using a GMM

The GMM will find homogeneous clusters (sub-classes) of the **none-of-the-others** population, and can be used as an estimate for the conditional density of a test sample assuming that it belongs to the **none-of-the-others** class<sup>4</sup>

<sup>4</sup>Given the complexity of the task, and due to the fact that different amounts of parameters are used for modeling the **none-of-the-others** class these kind of models often provide class-conditional likelihoods that are not calibrated, and may require further score processing (e.g. score calibration) to obtain good decisions

59

GMMs offer an elegant solution to this modeling challenge. We can collect a large set of unlabeled samples that we believe represent the **none-of-the-others** class and use EM to fit a GMM with many components. The algorithm will automatically discover homogeneous sub-clusters within this heterogeneous class, effectively learning a representation of the different types of "unknown" samples.

As we did for MVG, we can train GMM with diagonal covariance matrices to reduce the number of parameters to estimate (reducing overfitting and computational costs).

$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \quad \text{no correlation}$$

The solution is again given by the diagonals of  $\Sigma_c$ 's we defined before

We may need more components to model more complex distributions

The diagonal covariance assumption **DOES NOT** correspond to the **Naive Bayes** assumption in this case

The Naive Bayes assumption would correspond to training a **different** GMM (possibly with different number of components) for **each subset of features** that is assumed independent from the other features

We can also assume that all components of a GMM have the same covariance matrix (tied GMM)

Note that in this case we are tying the **components** of a **single GMM**, i.e. the Gaussian components of the GMM of a single class

This is different from the Tied Gaussian model, where parameters were shared **across** classes

Of course, we can extend GMM parameters tying across classes as well — we won't consider this model though, as it would require revisiting the EM estimation procedure, since sharing the parameters across classes would **not allow** us to **independently estimate a GMM over** the samples of each class

## MNIST — GMM (PCA 50)

Components:	1	2	4	8	16
FullCov	3.6%	3.4%	2.8%	2.3%	2.2%
Diagonal	12.3%	10.1%	8.9%	7.6%	6.2%
Components:	32	64	128	256	
FullCov	2.3%				
Diagonal	5.1%	4.3%	4.3%	4.3%	

### Performance Insights from Real Data

The MNIST results shown in your slides reveal important patterns about GMM classification performance. Several key insights emerge from this data:

**The Power of Multiple Components:** Moving from 1 to 16 components dramatically improves performance for both full and diagonal covariance models. With full covariances, error rates drop from 3.6% to 2.2%. This demonstrates that even relatively simple datasets like MNIST benefit significantly from multi-component modeling.

**Covariance Structure Matters:** Full covariance matrices consistently outperform diagonal ones, but the gap narrows as the number of components increases. With 1 component, full covariances achieve 3.6% error versus 12.3% for diagonal. But with 16 components, the gap closes to 2.2% versus 6.2%. This suggests that multiple components can partially compensate for the limitations of diagonal covariances.

**Diminishing Returns:** Performance improvements slow down as we add more components. The jump from 1 to 2 components provides substantial gains, but going from 128 to 256 components shows minimal improvement. This reflects the bias-variance trade-off: more components provide greater flexibility but require more data to estimate reliably.

## Initialization problem

$$\sum_i = I \quad w_i = \frac{1}{k} \quad \text{INITIAL PARAMETERS}$$

Unlike simple optimization problems with convex objective functions, GMM training faces the challenge of a highly non-convex likelihood surface riddled with local maxima. The EM algorithm is guaranteed to find a local maximum, but there's no guarantee it will find the global maximum or even a good local maximum. This makes initialization absolutely critical.

Think of the likelihood surface as a mountainous landscape with many peaks of different heights. The EM algorithm is like a hiker who always walks uphill but can only see their immediate surroundings. Where you start your hike largely determines which peak you'll reach. A poor starting point might lead you to a small hill when there's a towering mountain just a few valleys away.

### The Random Initialization Problem:

Random initialization often leads to poor results because it doesn't provide the algorithm with any guidance about the underlying structure in the data. You might initialize components in empty regions of the feature space, or multiple components might start so close together that they compete for the same data points throughout training.

### K-means as a Natural Initializer:

The connection between GMMs and K-means provides an elegant initialization strategy. When we constrain GMMs to have isotropic (spherical) covariances and make hard cluster assignments, we essentially recover the K-means algorithm. This connection suggests using K-means clustering results to initialize GMM parameters.

The K-means initialization process works as follows:

Run K-means clustering on your data to get initial cluster centers

- Initialize GMM means  $\mu_C$  to the K-means centroids
- Initialize covariances  $\Sigma_C$  based on the sample covariances of the K-means clusters
- Initialize mixture weights  $w_C$  based on the relative sizes of the K-means clusters

This approach ensures that each GMM component starts in a region where there's actual data, and the components are reasonably well-separated from each other.

## An alternative approach is the LBG Algorithm

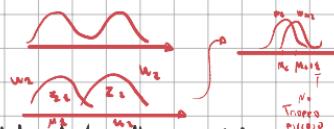
In the LBG instead of jumping directly to your target number of components you start simply and progressive split component to increase model complexity

- ① Start with a single Gaussian component fitted to all your data
- ② Split this component into two by creating a small perturbation:  $\mu_1' = \mu_1 + \epsilon$  and  $\mu_2' = \mu_1 - \epsilon$

- ③ Run EM to converge on this 2-component model

- ④ Split the component that contributes most to the likelihood (or all components) to get 3 components

- ⑤ Continue until you reach your desired number of components.



A good value for  $\epsilon$  can be a displacement along the principal eigenvector of the covariance matrix  $\Sigma_C$ .

We have sample  $x$  with a lot of feature (so  $\mu$  is a vector with the mean for each feature)

The eigenvectors of a covariance matrix tell us the principal direction of variation in the data, while the eigenvalues tell us how much variation exists along each direction.

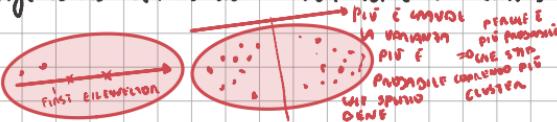
So when we have  $Z = U\Lambda V^T$

eigenvector as column

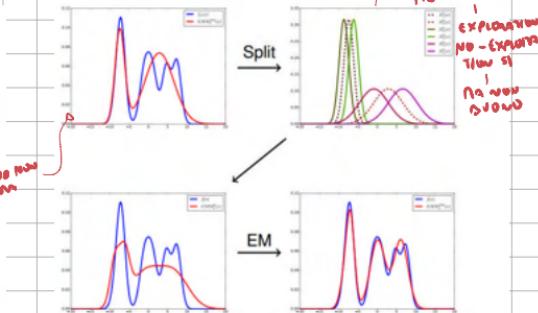
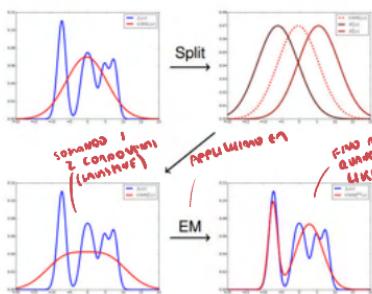
the first eigenvector that corresponds to the highest eigenvalue points in the direction where the data variate the most. This is the direction where it makes most sense to split a gaussian component into two sub-components.

$$\epsilon = \sigma \cdot u_1$$

magnitude of the split



26 → 26 - 0 ... 2336



One of the most difficult practical problems in GMM training is determining the appropriate number of components. This is particularly challenging because the likelihood always increases (or at least doesn't decrease) as you add more components, making it impossible to use likelihood alone for model selection. (overfitting)

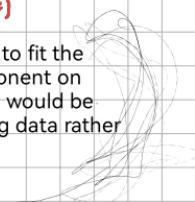
**Why Likelihood Fails as a Selection Criterion:**

Adding more components always provides the model with greater flexibility to fit the training data. In the extreme case, you could place one tiny Gaussian component on each training point and achieve arbitrarily high likelihood. But such a model would be useless for generalization because it would essentially memorize the training data rather than learning underlying patterns.

We cannot choose based on likelihood alone

Several criteria, more or less successful, have been proposed (AIC, BIC)

We can also resort to cross-validation usando questa



We also need to pay attention to degenerate models

As we said at the beginning, the log-likelihood for a GMM, as long as we have at least two components, is unbounded

The EM algorithm will usually find local maxima that are well-behaved, however, especially if we have too many components, we may obtain degenerate models which cause numerical issues

Some heuristics can be used to force models to be well-behaved (e.g. impose minimum values for the eigenvalues of the covariance matrices, tie the covariance of different components)

We can also modify our initialization so that the algorithm may end up in a different local maximum

Note: Degenerate solution typically occur when one or more components become pathologically concentrated on a small number of data points. As discussed earlier, this can lead to covariance matrix with very small eigenvalues, causing numerical instability and unbounded likelihood values

