# Lecture 8 - 2/22

**Soft Clustering**

# Soft Clustering

*every row has a single assignment*

So far, clustering was done using **hard assignments** (1 point -> 1 cluster)

Sometimes this doesn't accurately represent the data: it seems reasonable to have overlapping clusters.

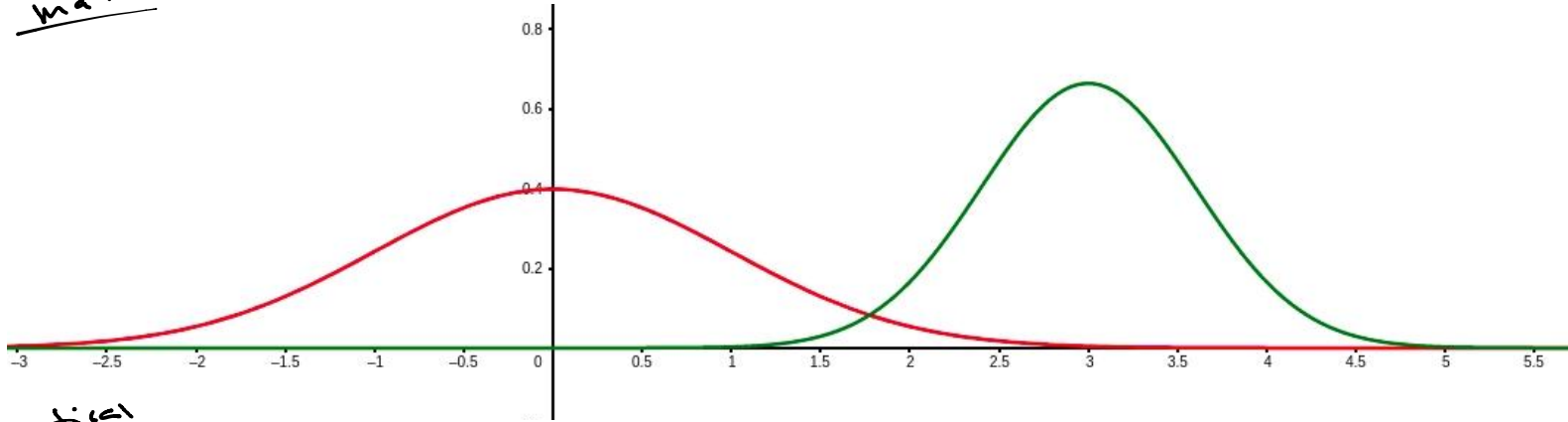In this case, we can use **soft assignment** to assign points to every cluster with a certain probability.

— *the model we will use*

# Soft Clustering - Example

Generate data using $N(\mu_1, \sigma_1)$ and $N(\mu_2, \sigma_2)$
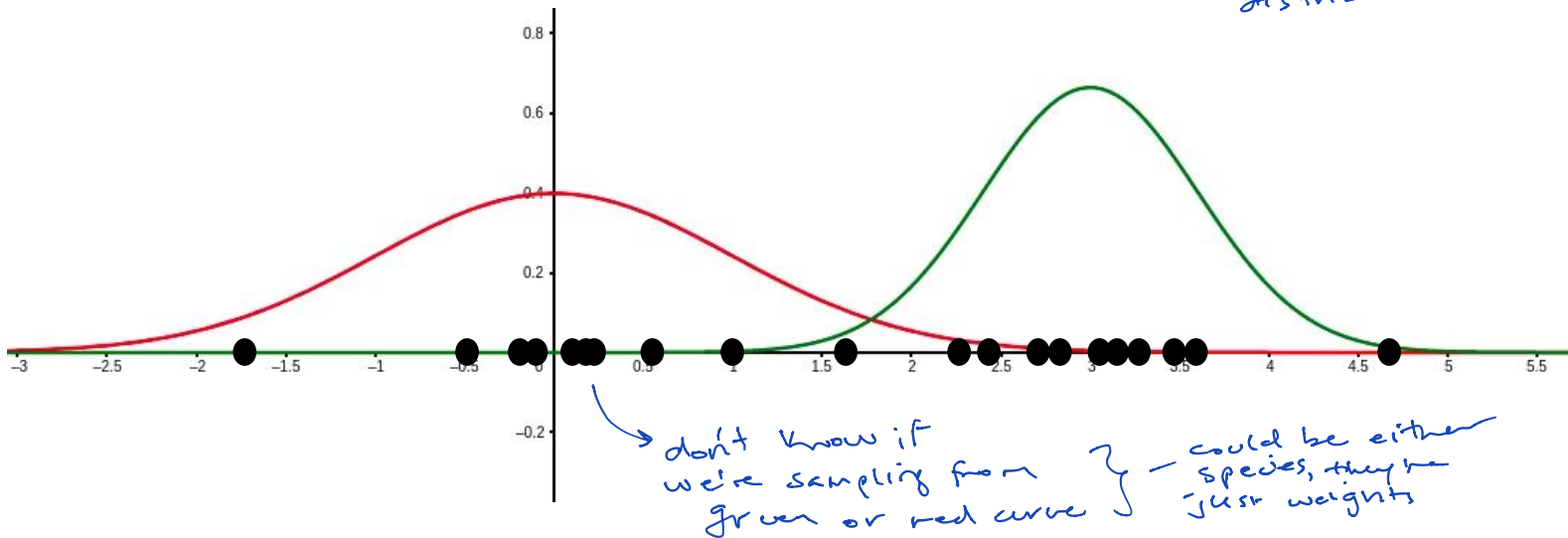
math



practical

Or: we are given the weights of animals. Unknown to us these are weights from two different species.
Can we determine the species (group / assignment) from the height? which animal is more likely?

# Soft Clustering - Example

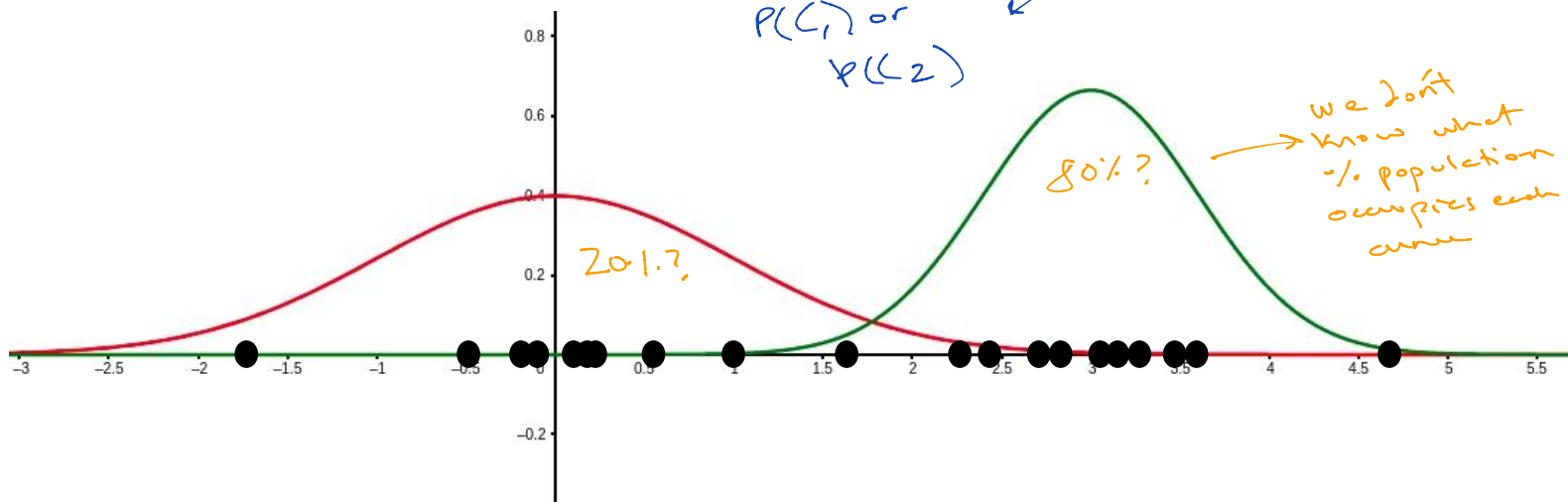Generate data using $N(\mu_1, \sigma_1)$ and $N(\mu_2, \sigma_2)$



— once you know whether you're picking red, or green
↳ Follow that distribution

→ don't know if we're sampling from green or red curve
} — could be either species, they're just weights

# Soft Clustering - Example

Generate data using **N(μ₁, σ₁)** and **N(μ₂, σ₂)**

_— d — this many times to generate dataset (but don't know upfront the prob. of picking from each curve)_

_P(C₁) or P(C₂)_



_80%?_

_20%?_

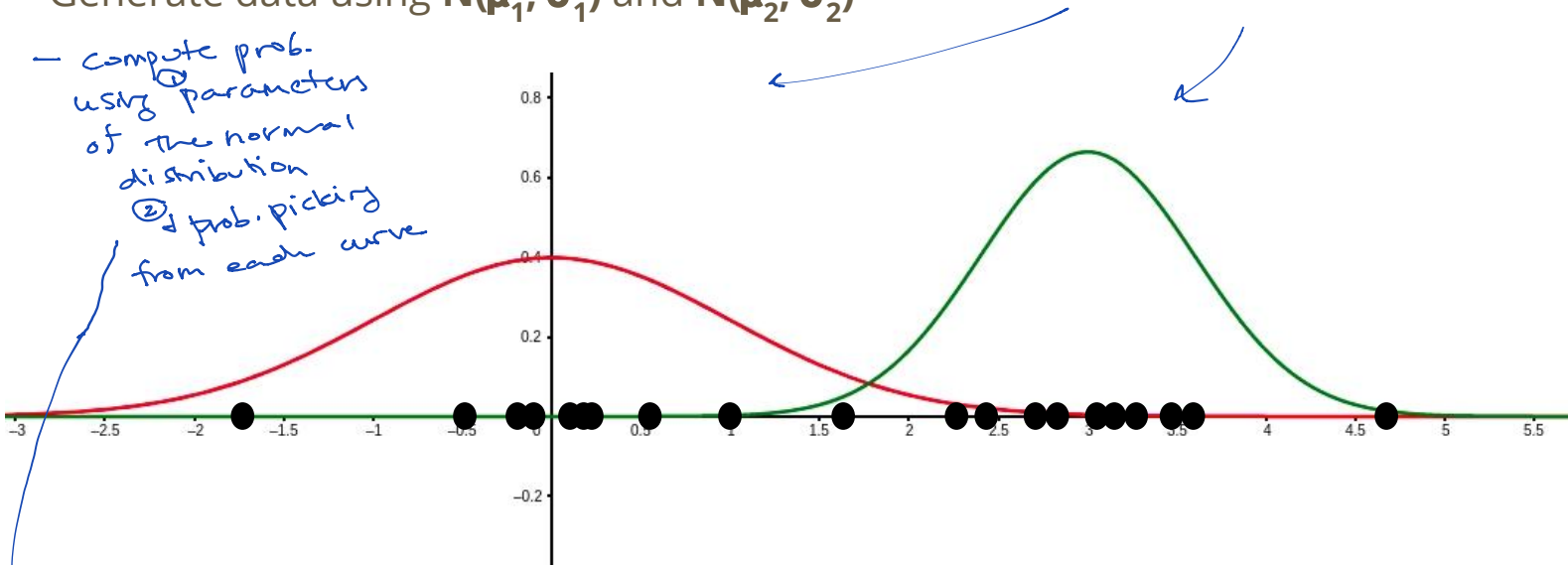_we don't know what % population occupies each curve_

Any of these points could technically have been generated from either curve.

# Soft Clustering - Example

Generate data using $N(\mu_1, \sigma_1)$ and $N(\mu_2, \sigma_2)$

*handwritten annotations:*

— compute prob.
using ① parameters
of the normal
distribution
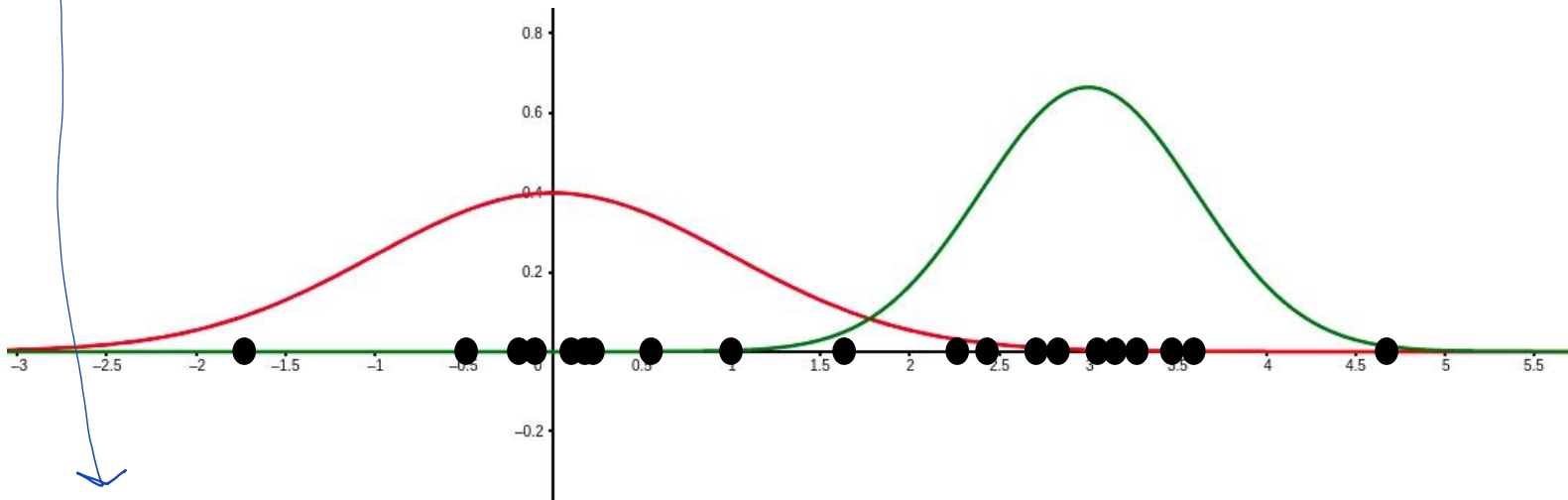② prob. picking
from each curve

*within the probability
of each curve (calculate)
we have this
distribution of weights/points*



For each point we can compute the probability of it being generated from either curve

# Soft Clustering - Example

Generate data using **N(μ₁, σ₁)** and **N(μ₂, σ₂)**



We can create soft assignments based on these probabilities.

# Mixture Model

X comes from a mixture model with k mixture components if the probability distribution of X is:

*# distributions / components / clusters*

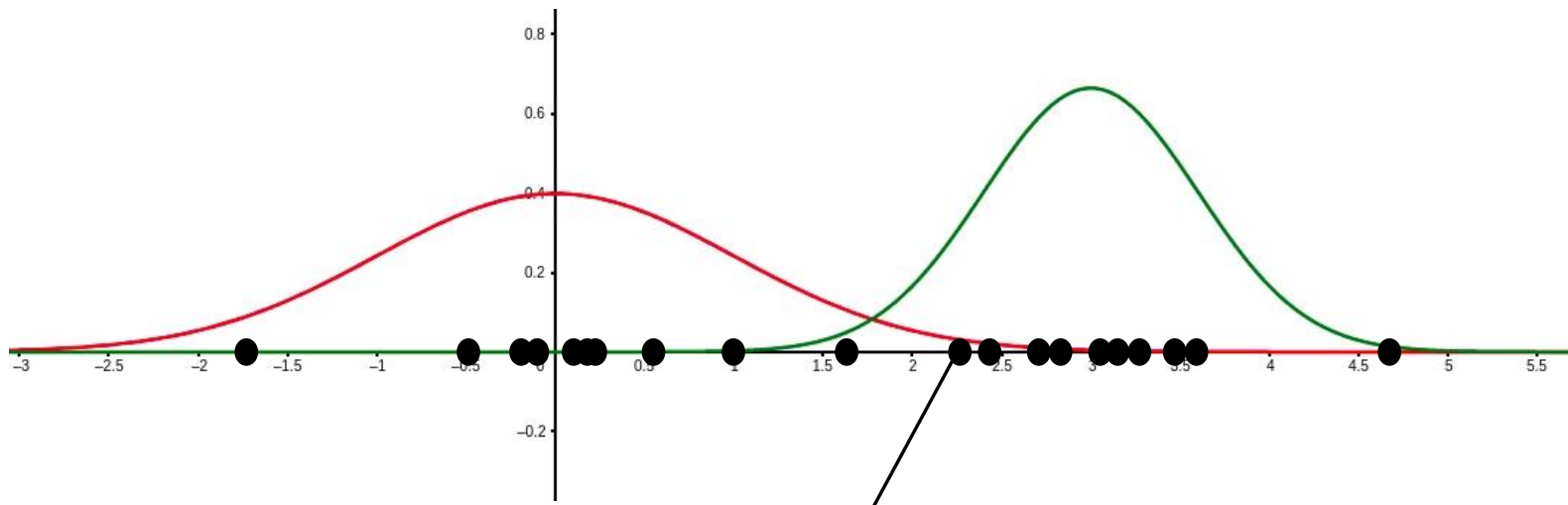$$P(X = x) = \sum_{j=1}^{k} P(C_j) P(X = x | C_j)$$

*once you know you're in that curve*

Mixture proportion
Represents the probability of belonging to $C_j$ *(a curve)*

Probability of seeing x when sampling from $C_j$

*Weighted distribution*

# Example
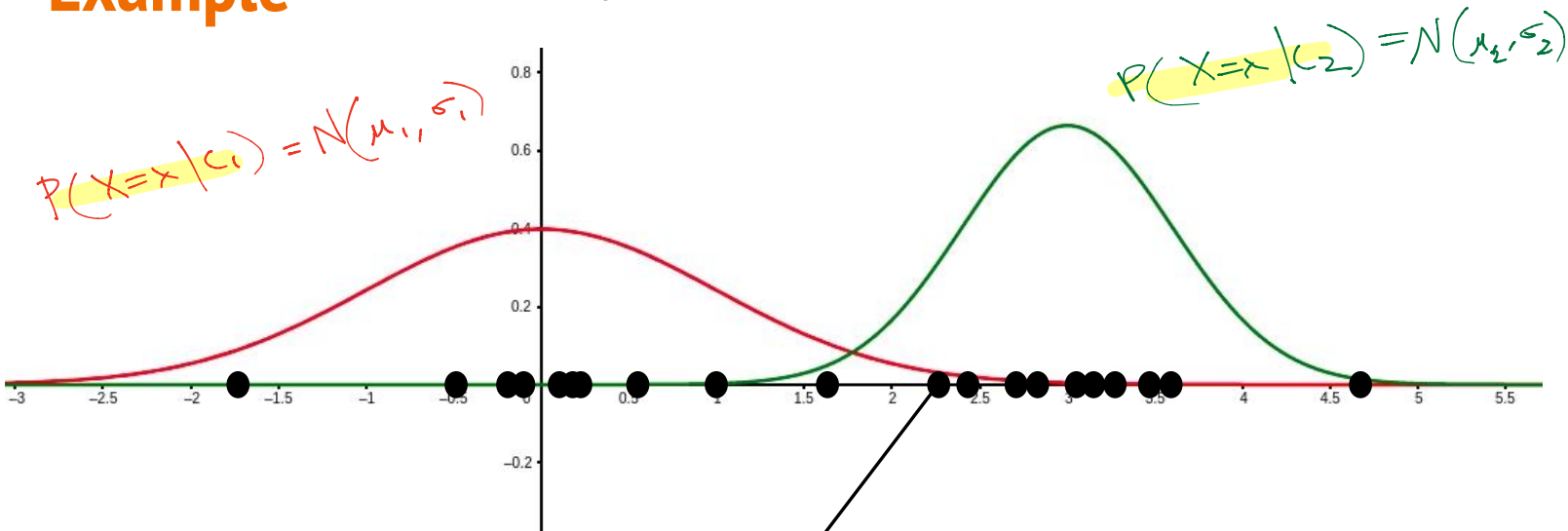


What is the probability distribution here?

# Example

→ clearly a "mixture distribution"
— so use the mixture model



$P(X=x|C_1) = N(\mu_1, \sigma_1)$

$P(X=x|C_2) = N(\mu_2, \sigma_2)$

$$P(X = x) = P(C_1)P(X = x|C_1) + P(C_2)P(X = x|C_2)$$

prob. of $c_1$ given you're in $c_1$
✗ prob. of being in $c_1$ ($P(C_1)$)

prob. of $C_2$ given you're in $C_2$
✗ prob. of being in $C_2$ ($P(C_2)$)

# Example



probability of seeing x

just P(x)

$$P(X = x) = P(C_1)\frac{1}{\sigma_1\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{x-\mu_1}{\sigma_1}\right)^2} + P(C_2)\frac{1}{\sigma_2\sqrt{2\pi}}e^{-\frac{1}{2}\left(\frac{x-\mu_2}{\sigma_2}\right)^2}$$

not given upfront

not given

any possible x that X can take on

# Gaussian Mixture Model

A Gaussian Mixture Model (GMM) is a mixture model where

$$P(X = x | C_i) \sim N(\mu, \sigma)$$

the probability within the cluster...

(once you know $P(C_n)$)

follows a normal distribution

(parameterized by $\mu$ and $\sigma$)

# GMM Clustering

**Goal**: Find the GMM that <mark>maximizes the probability of seeing the data we</mark> <mark>have.</mark>

*goal in all of ML*

Finding the GMM means finding the parameters that uniquely characterize it. What are these parameters?

*ex] what are parameters of normal distr. that maxb. the prob. of seeing the given data?*

$P(C_i)$ & $\mu_i$ & $\sigma_i$ for all **k** components.

*what parameters uniquely characterize the data?*

# GMM Clustering

**Goal**: Maximization function

each has this distribution

n data points from GMM, sampled independently

(size of dataset)

Find the argument $\theta$ that maximizes this

contains all parameters

$$\theta^* = \arg\max_{\theta} \prod_{i=1}^{n} \sum_{j=1}^{k} P(C_j) P(X_i \mid C_j)$$

product over all the datapoint of all their individual probabilities of the GMM

Joint probability distribution of our data

Assuming our data are independent

$\theta^*$ is the arg max of $\theta$

$\theta = \{ \mu_1 \dots \mu_k, \sigma_1 \dots \sigma_k, P(C_1) \dots P(C_k) \}$

, from calculus → take derivative w/ respect to product
set = 0, solve
↳ BUT
very complicated
for a product

# GMM Clustering

How do we find the critical points of this function?
↳ value of $\theta$ that is a local max/min

Notice: taking the log-transform does not change the critical points = trick :)
↳ transform into a space where the problem is easier to solve

Define:

※ transform the space

$$l(\theta) = \log(L(\theta))$$

$$= \sum_{i=1}^{n} \log\left(\sum_{j=1}^{k} P(C_j)P(X_i \mid C_j)\right)$$

↳ the product becomes a sum (much simpler)

# GMM Clustering

For $\boldsymbol{\mu} = [\boldsymbol{\mu}_1, ..., \boldsymbol{\mu}_k]^T$ and $\boldsymbol{\Sigma} = [\boldsymbol{\Sigma}_1, ..., \boldsymbol{\Sigma}_k]^T$

We can solve

$$\frac{d}{d\Sigma} l(\theta) = 0 \qquad \frac{d}{d\mu} l(\theta) = 0$$

$\hookrightarrow$ find best value for $\Sigma$

$\hookrightarrow$ find best value for $\mu$

# GMM Clustering

To get

The Logic

the solutions →

① have an optimization function trying to maximize

② difficult to take max as is, so LOG TRANSFORM

③ take derivative w/ respect to each parameter, set $=0$, solve get CRITICAL POINTS

$$\hat{\mu}_j = \frac{\sum_{i=1}^n P(C_j|X_i)X_i}{\sum_{i=1}^n P(C_j|X_i)}$$

$$\hat{\Sigma}_j = \frac{\sum_{i=1}^n P(C_j|X_i)(X_i - \hat{\mu}_j)^T(X_i - \hat{\mu}_j)}{\sum_{i=1}^n P(C_j|X_i)}$$

→ how do you got this?

$$\hat{P}(C_j) = \frac{1}{n}\sum_{i=1}^n P(C_j|X_i)$$

∧ represents estimated value

# GMM Clustering

Do we have everything we need to solve this?

Still need **P(C$_j$ | X$_i$)** (i.e. the probability that X$_i$ was drawn from C$_j$)

# GMM Clustering

Bayes Rule

$$P(C_j|X_i) = \frac{P(X_i|C_j)}{P(X_i)}P(C_j)$$

$$= \frac{P(X_i|C_j)P(C_j)}{\sum_{j=1}^{k} P(C_j)P(X_i|C_j)}$$

but... you need $P(C_j|X_i)$ to get this? → see eqns on prev. slide

why we need an algorithm to do this

Can't just solve an equation... need an algorithm

Looks like a loop! Seems we need P(C$_j$) to get P(C$_j$ | X$_i$) and P(C$_j$ | X$_i$) to get P(C$_j$)

(stuck)

# Expectation Maximization Algorithm

1. Start with random $\theta$
2. Compute $P(C_j \mid X_I)$ for all $X_i$ by using $\theta$ → *can compute using random values (what we were trying to solve for before)*
3. Compute / Update $\theta$ from $P(C_j \mid X_I)$
4. Repeat 2 & 3 until <mark>convergence</mark>

*θ values help each other by updating throughout the algorithm*

*why would we converge necessarily? (ask in OH) — not as simple as Lloyd's algorithm*

**Demo**

wrote some code
that uses
GMM

( see the repo )

→ Clustering Aggregation