

HW2. Start early. Otherwise you can't finish

FOR LOOP



Parallel (broadcasting) numpy



Gaussian Mixture Model

Mahdi Roozbahani
Georgia Tech

**When you're actually paying attention in class
and you still have no idea what's going on**



Outline

- Overview 
- Gaussian Mixture Model
- The Expectation-Maximization Algorithm

Recap

$$P(A, B | C) = P(A | B, C) P(B | C)$$

Conditional probabilities:

$$p(A, B) = p(A|B)p(B) = p(B|A)p(A)$$

Bayes rule:

$$p(A|B) = \frac{p(A, B)}{p(B)} = \frac{p(B|A)p(A)}{p(B)}$$

$$p(A = 1) = \sum_{i=1}^K p(A = 1, B_i) = \sum_{i=1}^K p(A|B_i) p(B_i)$$

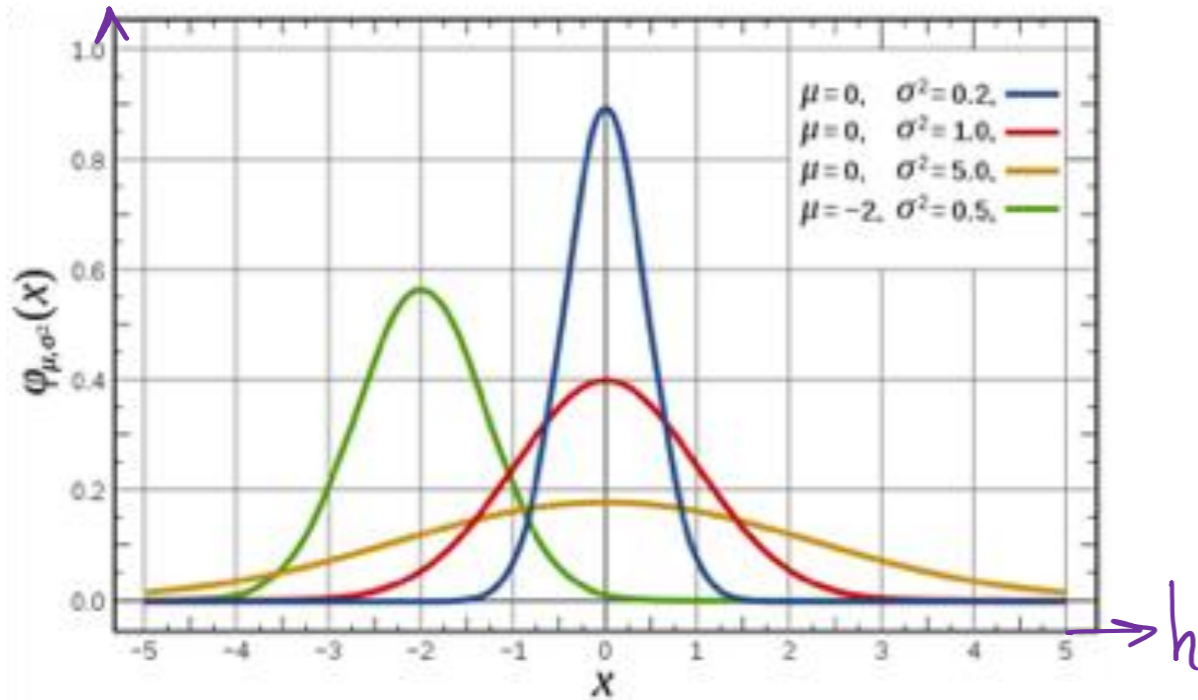
| | | Tom | | |
|-------------|-------------|---|---|---|
| | | Tomorrow=Rainy | Tomorrow=Cold | P(Today) |
| Tod | Today=Rainy | $\frac{4}{9}$ | $\frac{2}{9}$ | $[\frac{4}{9} + \frac{2}{9}] = \frac{2}{3}$ |
| | Today=Cold | $\frac{2}{9}$ | $\frac{1}{9}$ | $[\frac{2}{9} + \frac{1}{9}] = \frac{1}{3}$ |
| P(Tomorrow) | | $[\frac{4}{9} + \frac{2}{9}] = \frac{2}{3}$ | $[\frac{2}{9} + \frac{1}{9}] = \frac{1}{3}$ | |

$$\begin{aligned}
 P(\text{Tomorrow} = \text{Rainy}) &= \sum_{\text{Tod}} P(\text{Tom} = \text{rainy}, \text{Tod}) \\
 &= P(\text{Tom} = \text{rainy}, \text{Tod} = \text{rainy}) + P(\text{Tom} = \text{rainy}, \text{Tod} = \text{cold}) \\
 &= \frac{4}{9} + \frac{2}{9} = \frac{6}{9}
 \end{aligned}$$

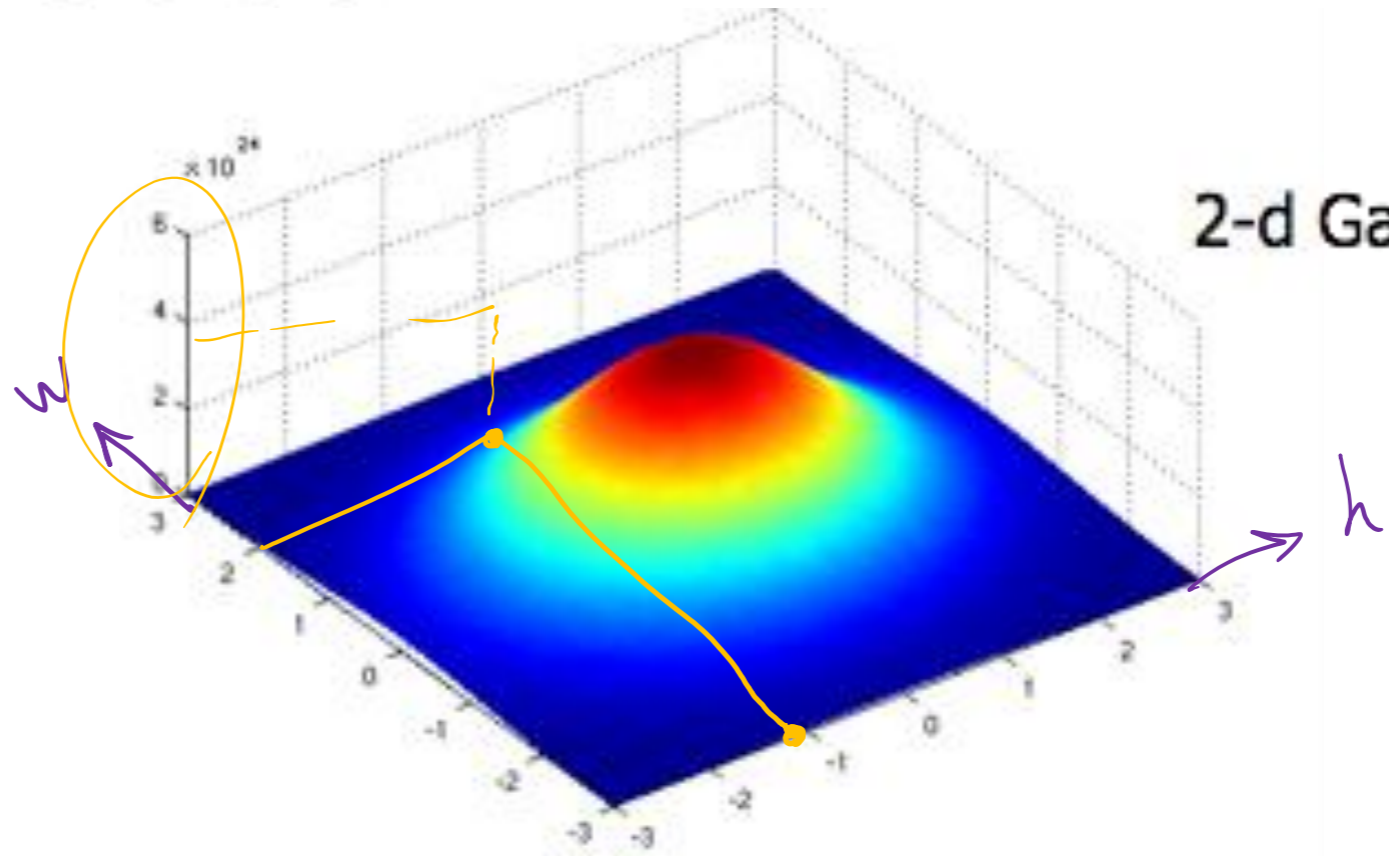
Gaussian Distribution

density

1-d Gaussian



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



2-d Gaussian

$$x = \begin{bmatrix} 0 & 0 & 0 \end{bmatrix}_{n \times d}$$

What is a Gaussian?

For d dimensions, the Gaussian distribution of a vector $x = (x_1, x_2, x_3, \dots, x_d)^T$ is defined by:

$$N(x | \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} \sqrt{|\Sigma|}} \exp\left(-\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu)\right)$$

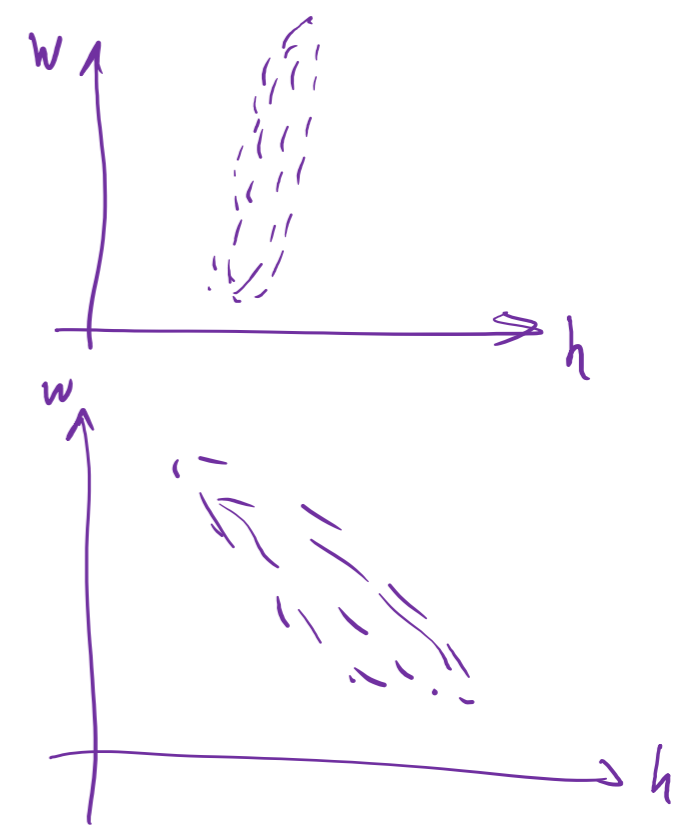
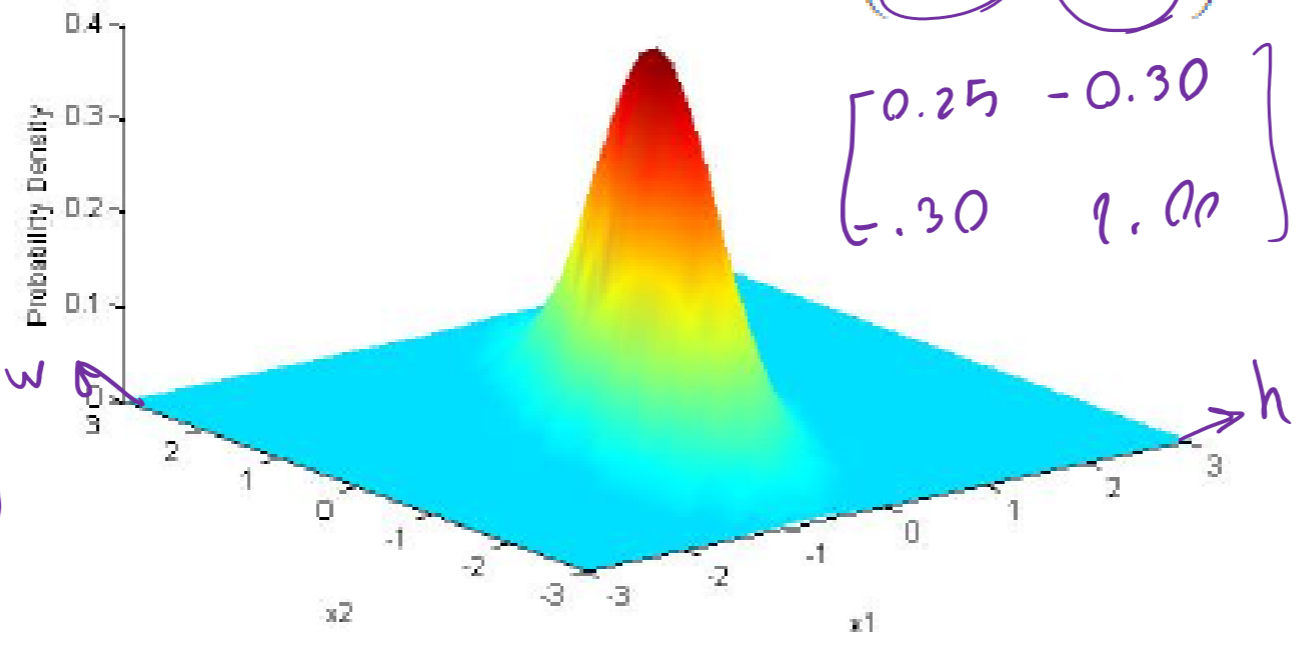
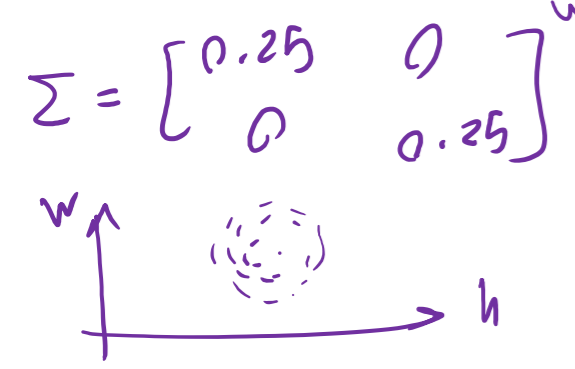
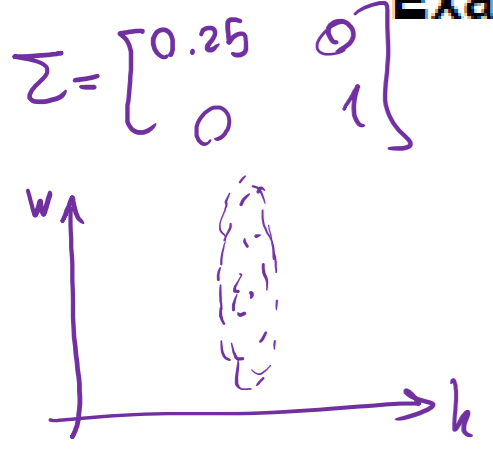
Handwritten annotations: μ is $1 \times d$, Σ is $d \times d$, $(x - \mu)^T$ is $1 \times d$, Σ^{-1} is $d \times d$, $(x - \mu)$ is $d \times 1$.

where μ is the mean and Σ is the covariance matrix of the Gaussian.

Example:

$$\mu = (0, 0)^T$$

$$\Sigma = \begin{bmatrix} 0.25 & 0.30 \\ 0.30 & 1.00 \end{bmatrix}$$

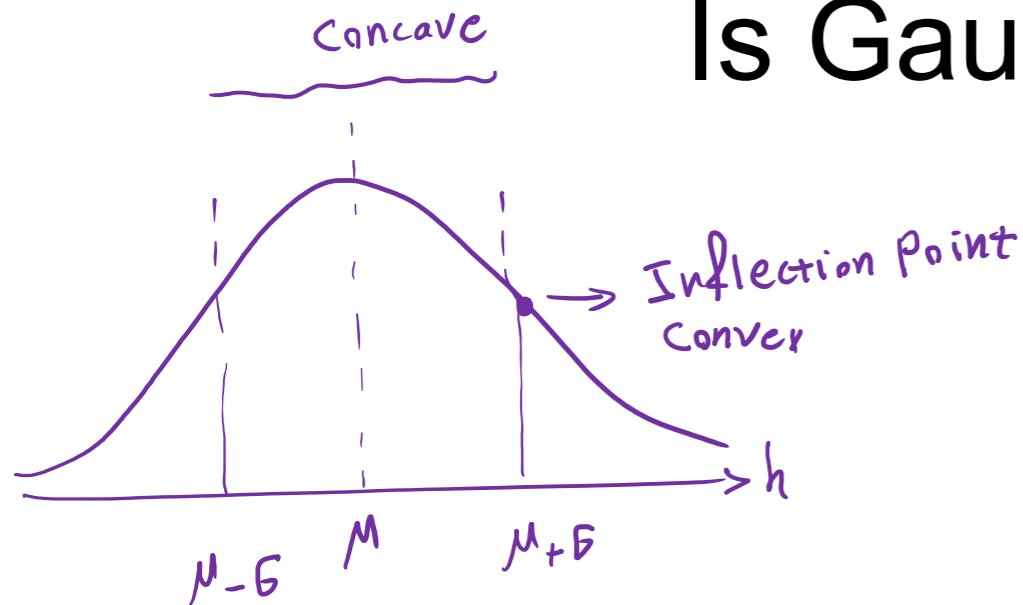


Is Gaussian Concave?

$$\log a * b = \log a + \log b$$

$$\log a/b = \log a - \log b$$

$$\log(\exp(u)) = u$$

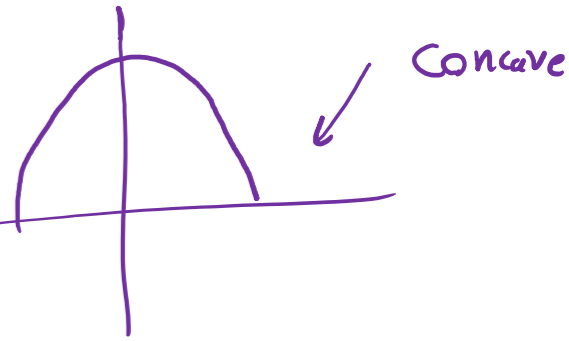


$$f(x | \mu, \sigma) = \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right) \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$$

$$f(x | \mu, \sigma) = c_1 * \exp(-c_2 (x - c_3)^2)$$

$$\log f(x | \mu, \sigma) = \log c_1 - u = \log c_1 - c_2 (x - c_3)^2 \approx -x^2$$

$$-x^2 - x^2 - x^2 = -3x^2$$

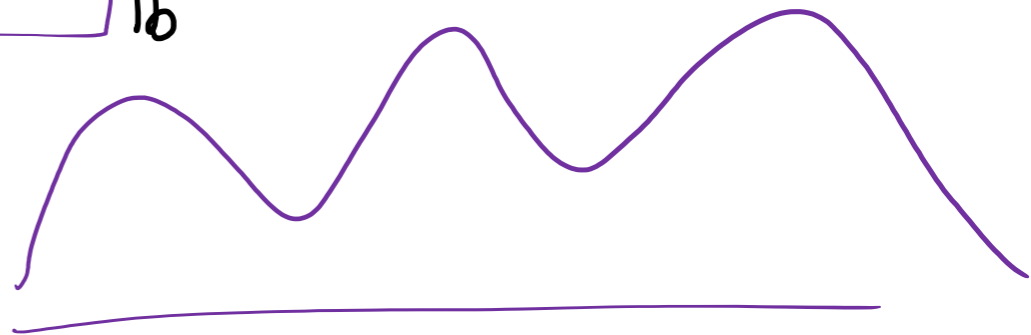


$$\log f(x | \mu_0, \sigma_0) + \log f(x | \mu_1, \sigma_1) + \log f(x | \mu_2, \sigma_2)$$


↑ Jensen Inequality

GMM

$$\log \left(f(x | \mu_0, \sigma_0) + \dots + f(x | \mu_k, \sigma_k) \right)$$

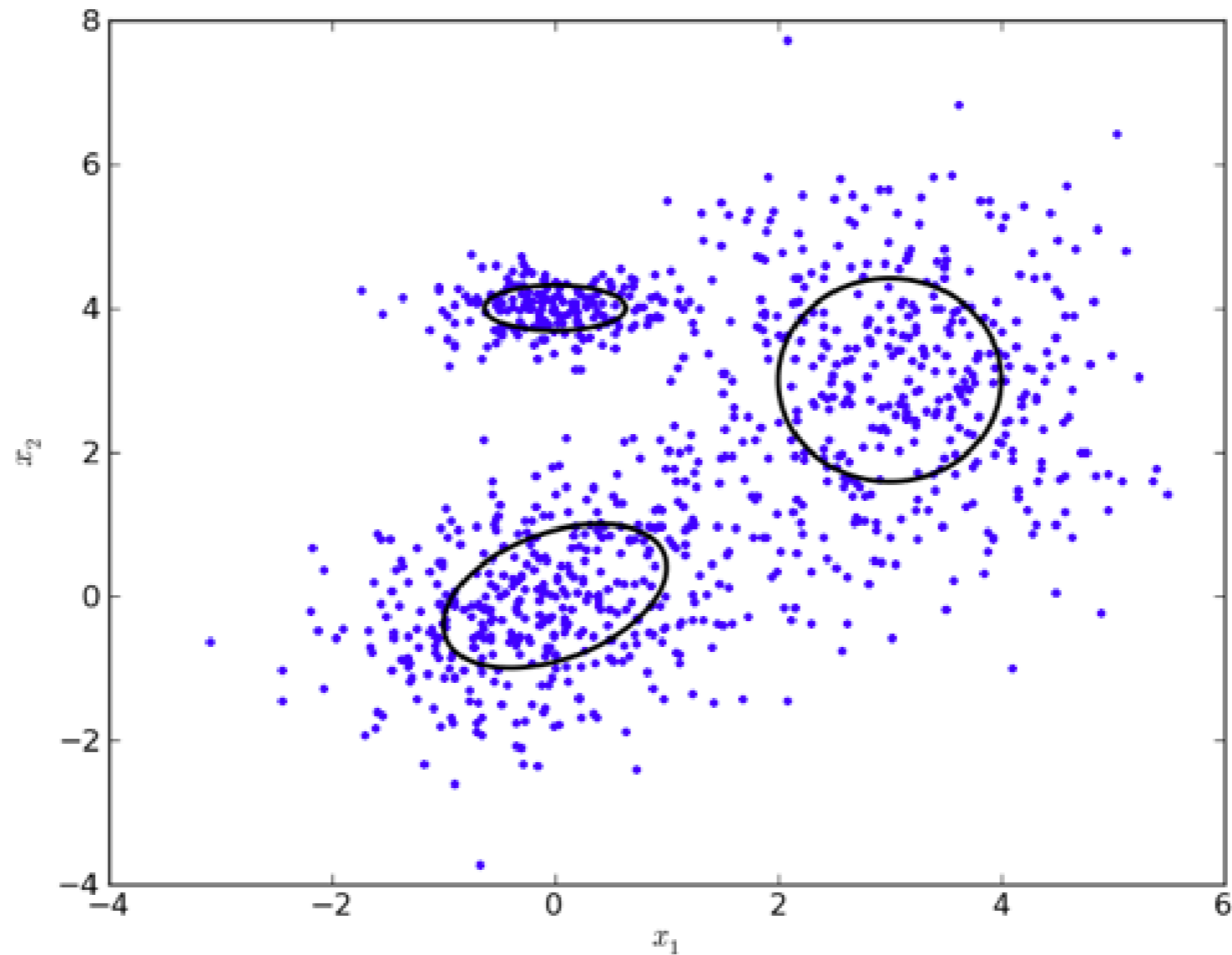


Outline

- Overview
- Gaussian Mixture Model 
- The Expectation-Maximization Algorithm

Hard Clustering Can Be Difficult

- Hard Clustering: K-Means, Hierarchical Clustering, DBSCAN



Towards Soft Clustering

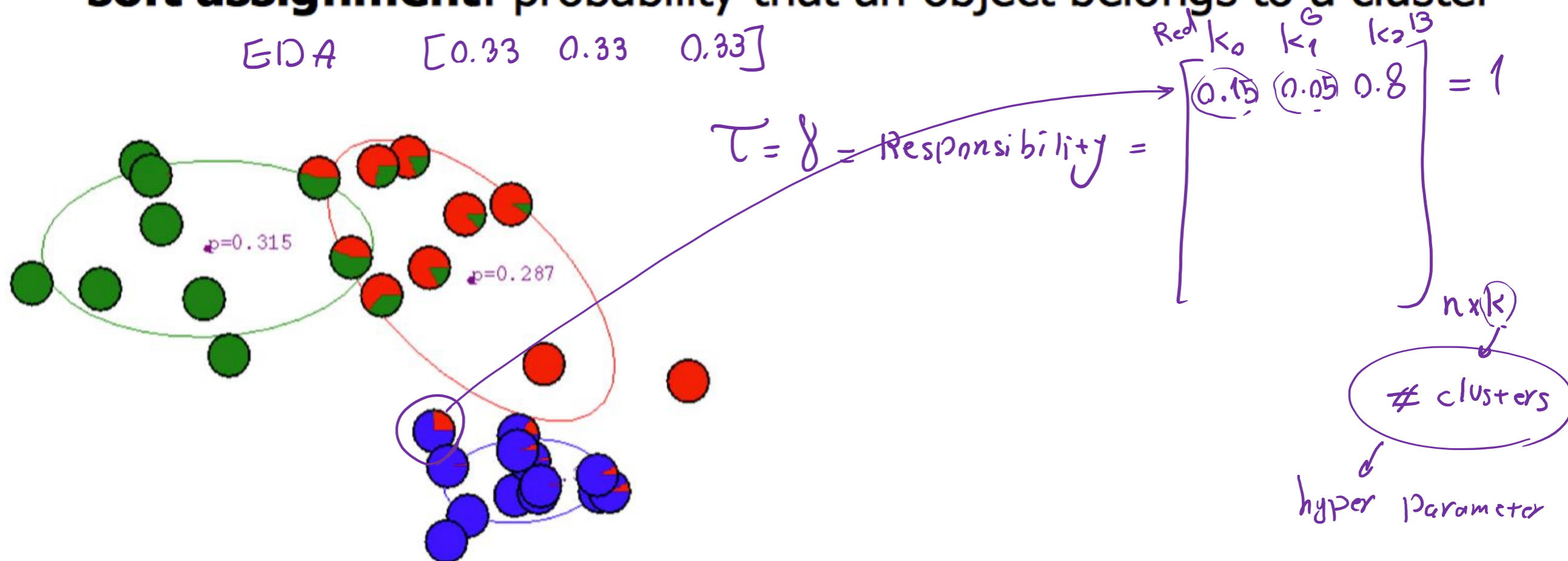
- **K-means**

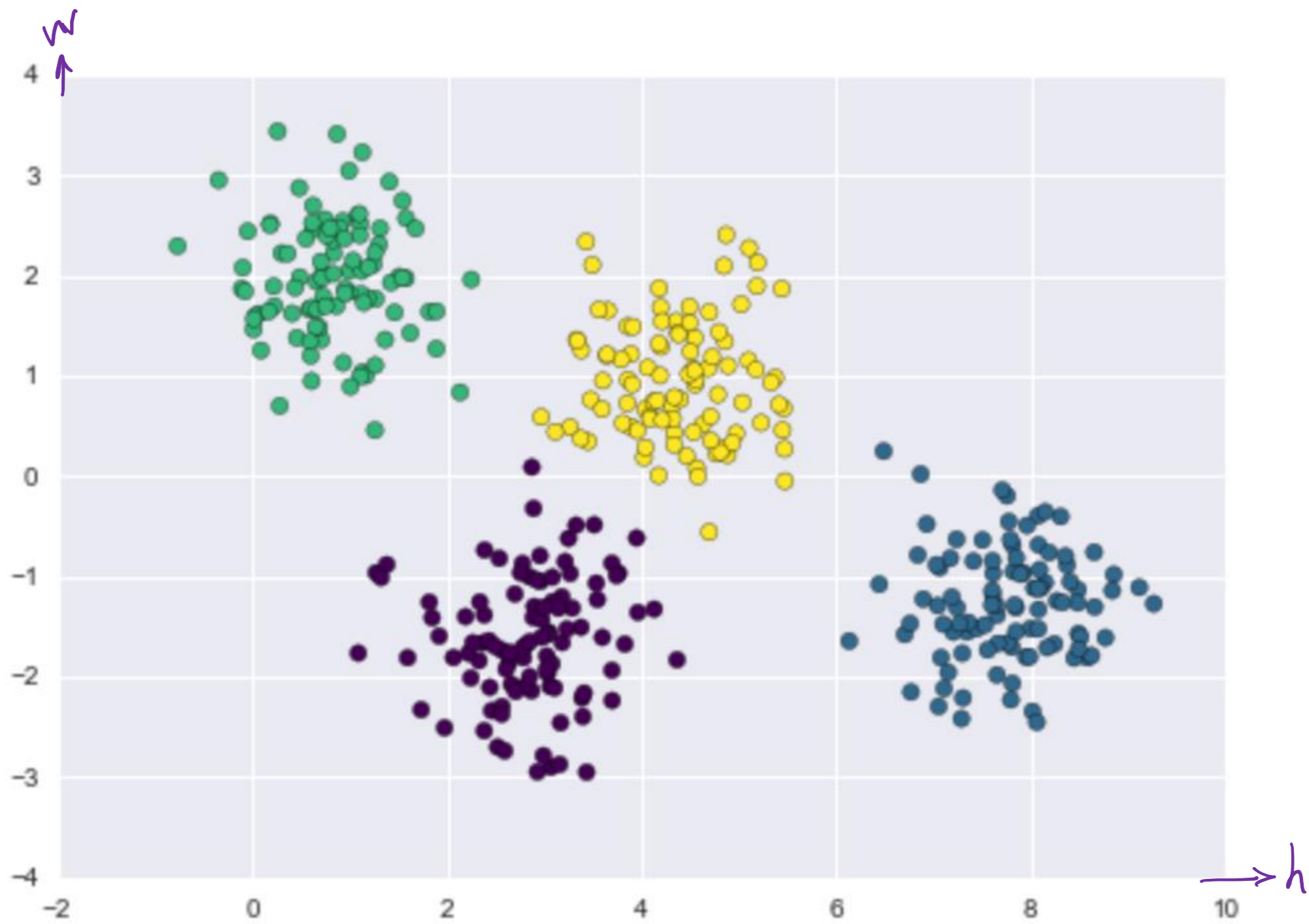
- **hard assignment:** each object belongs to only one cluster

$$\theta_i \in \{\theta_1, \dots, \theta_K\}$$

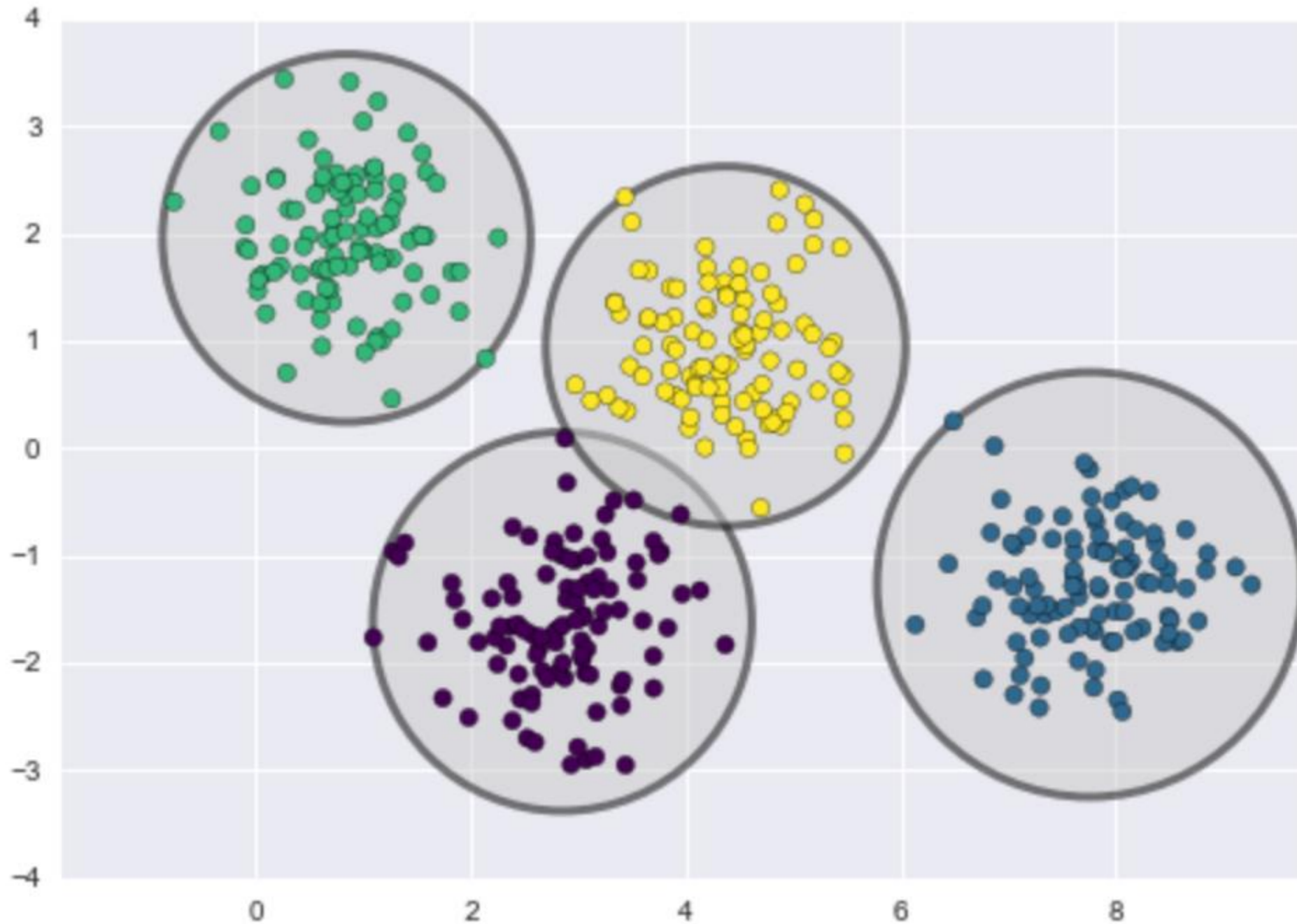
- **Mixture modeling**

- **soft assignment:** probability that an object belongs to a cluster

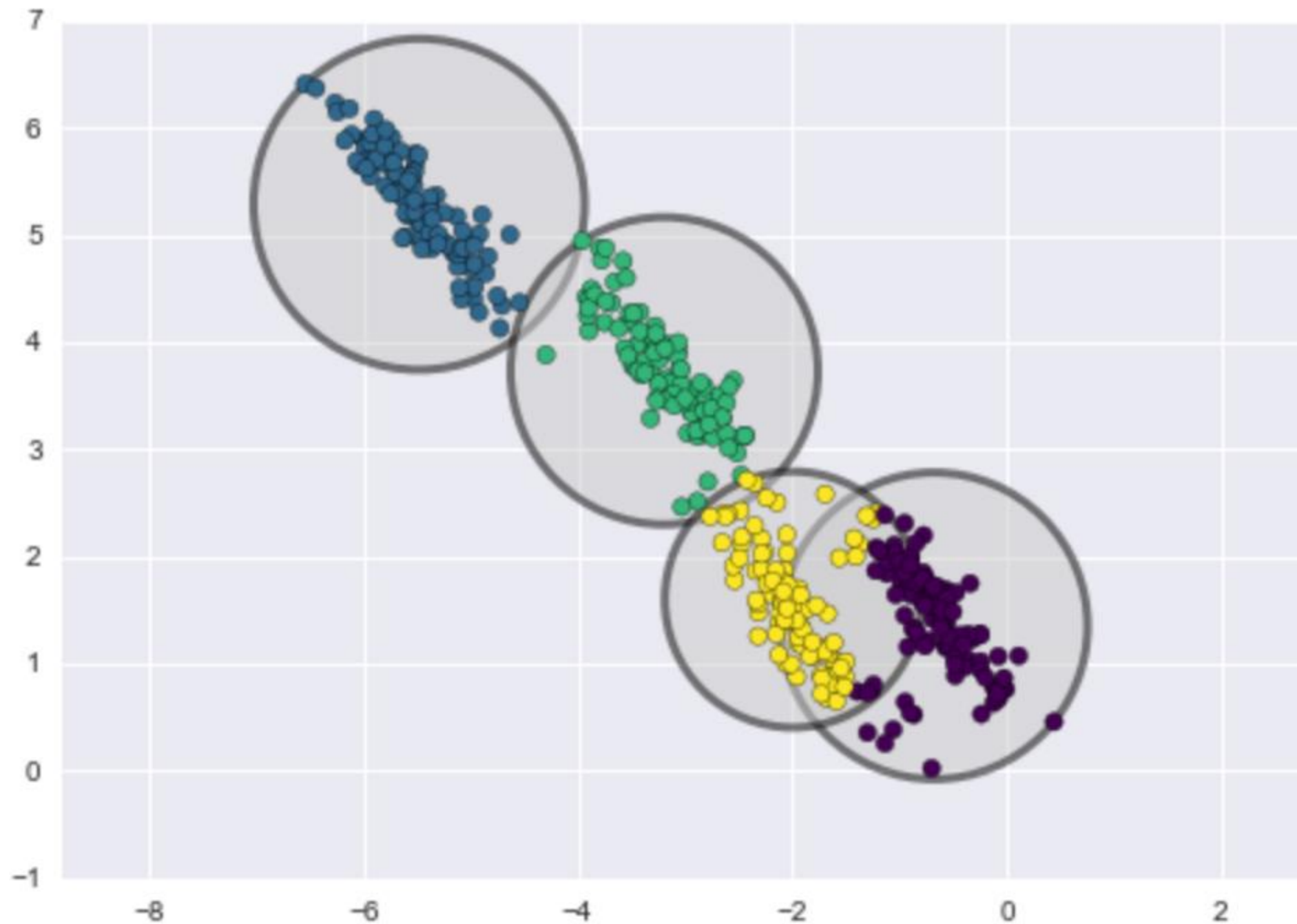




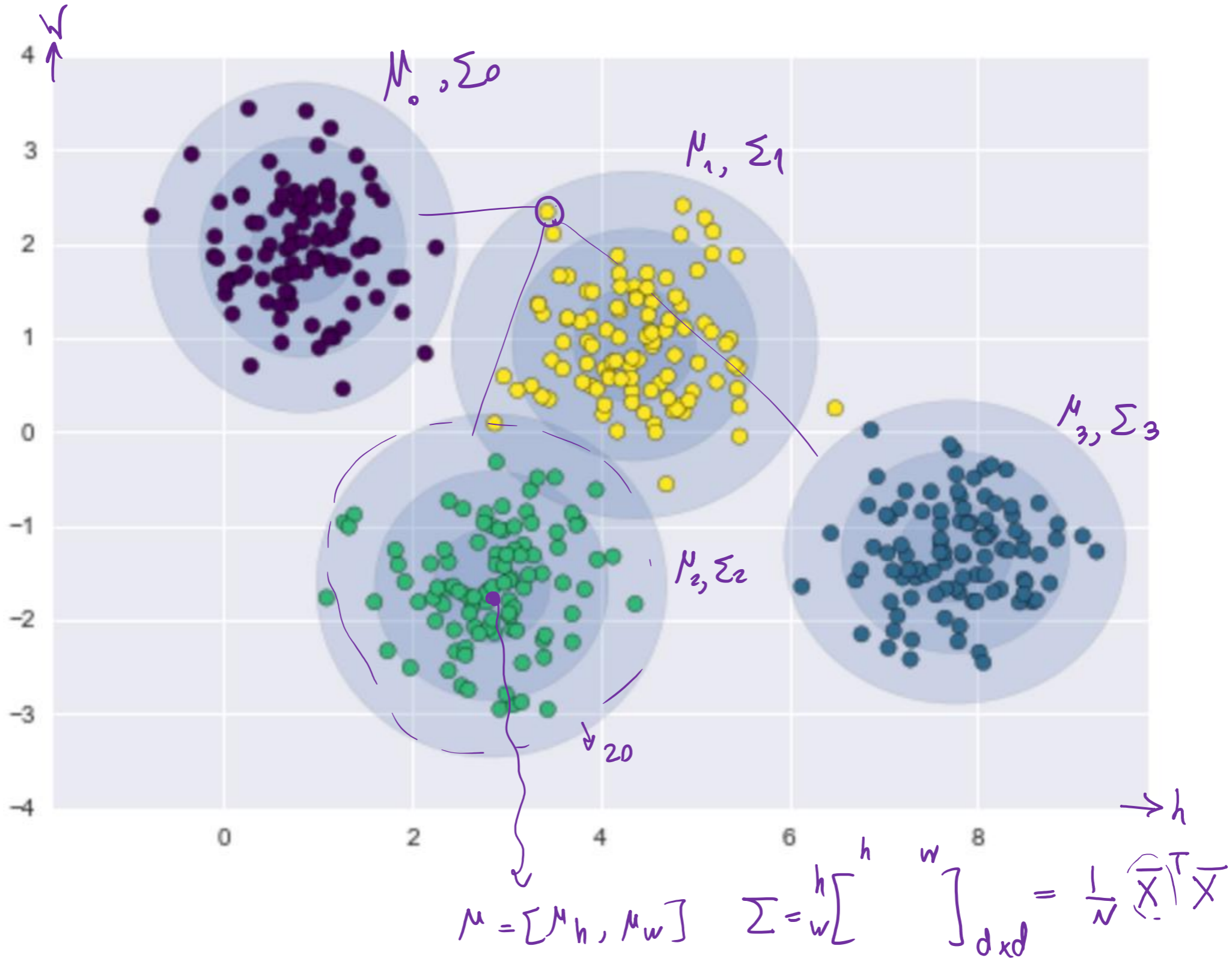
Let's run K-Means on the dataset

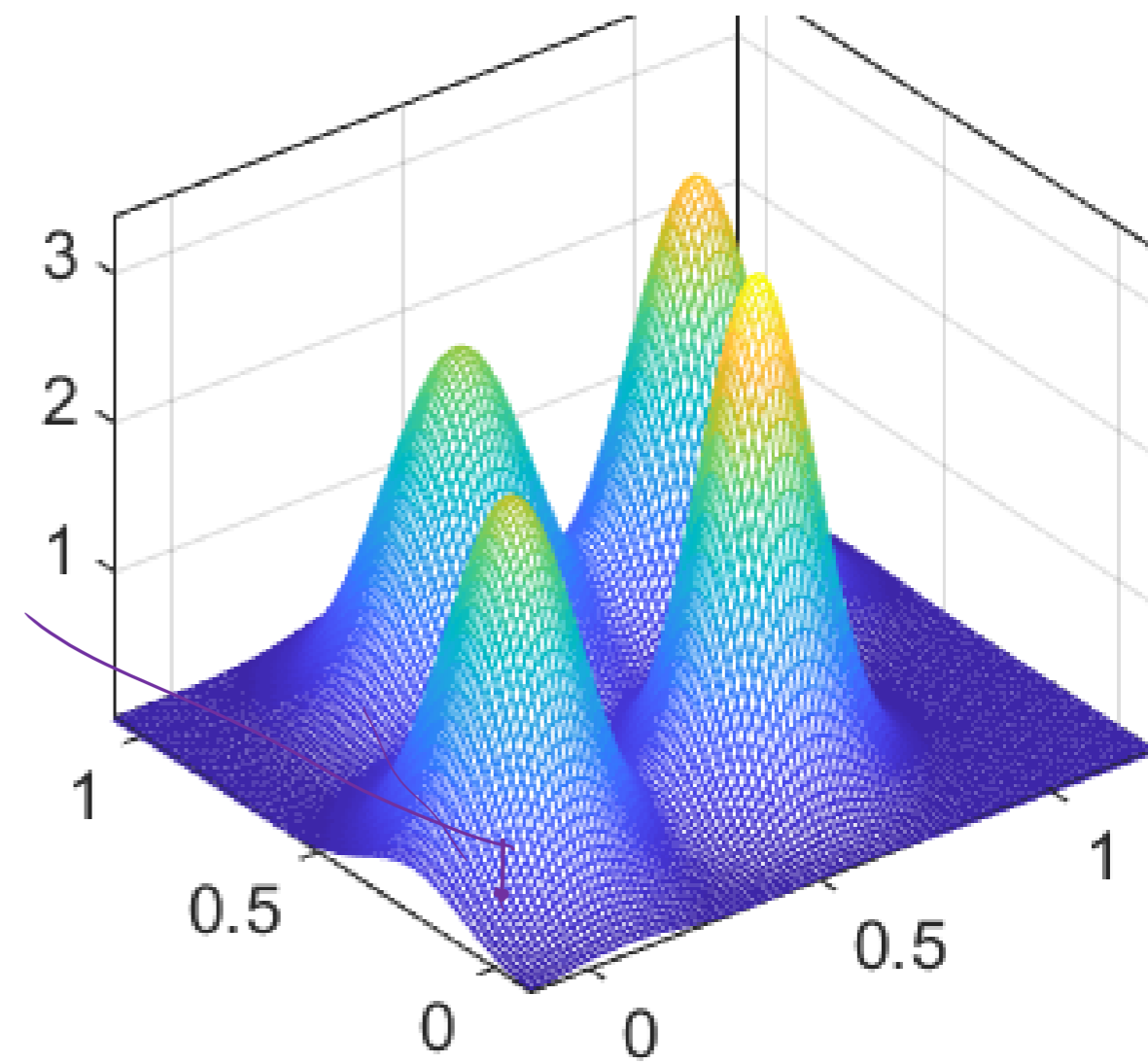
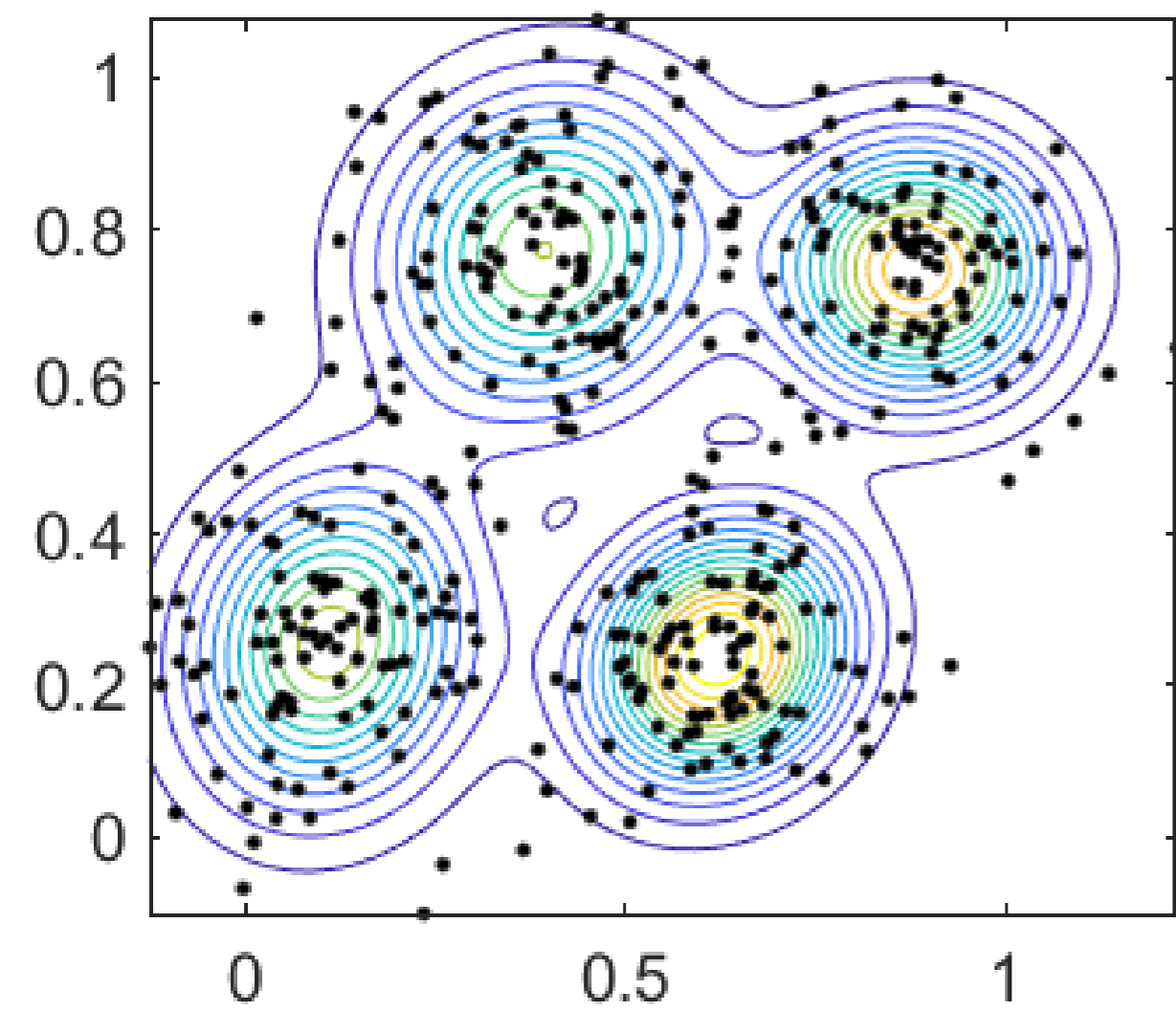


Let's generate a new dataset and run K-Means

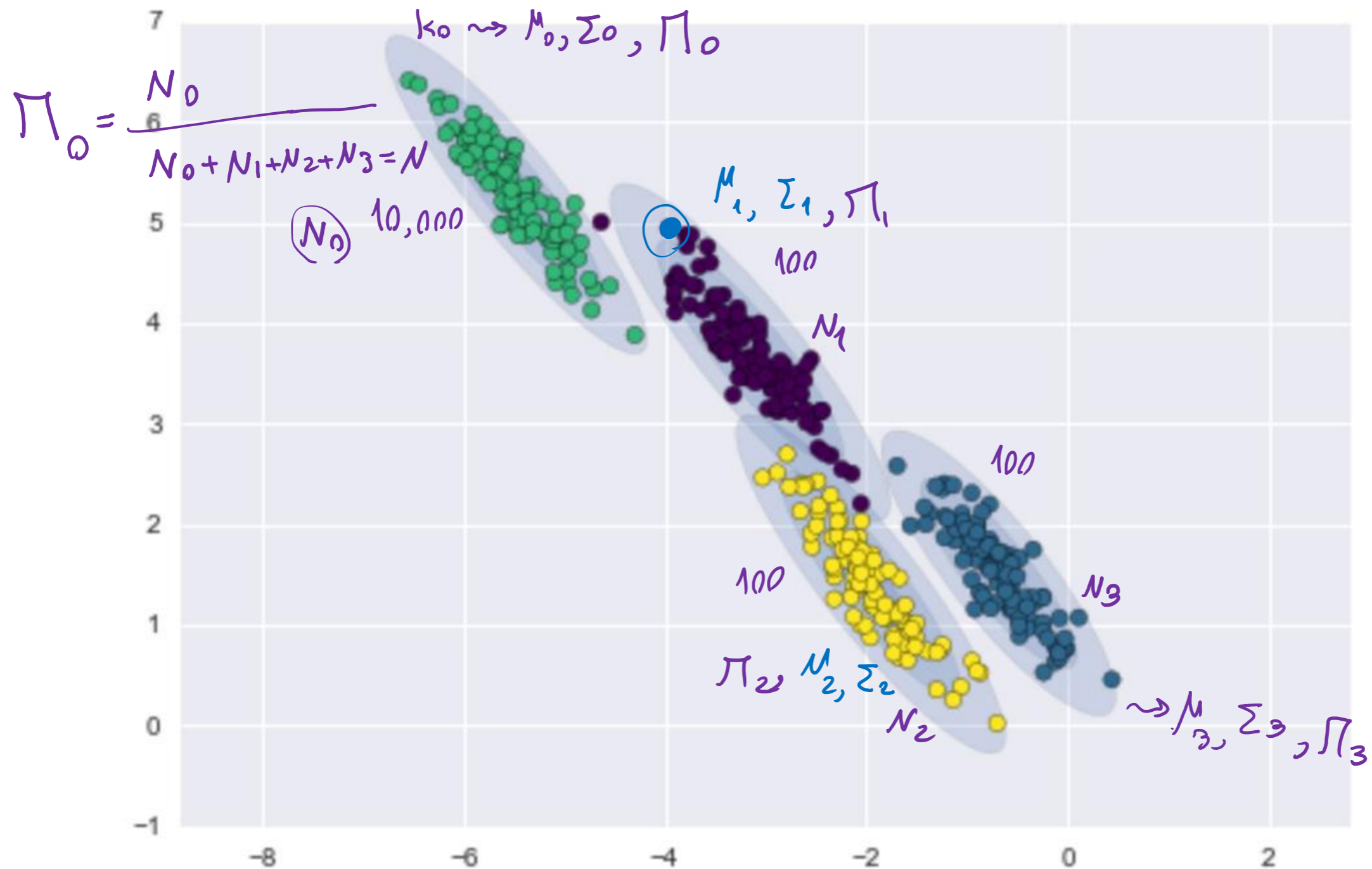


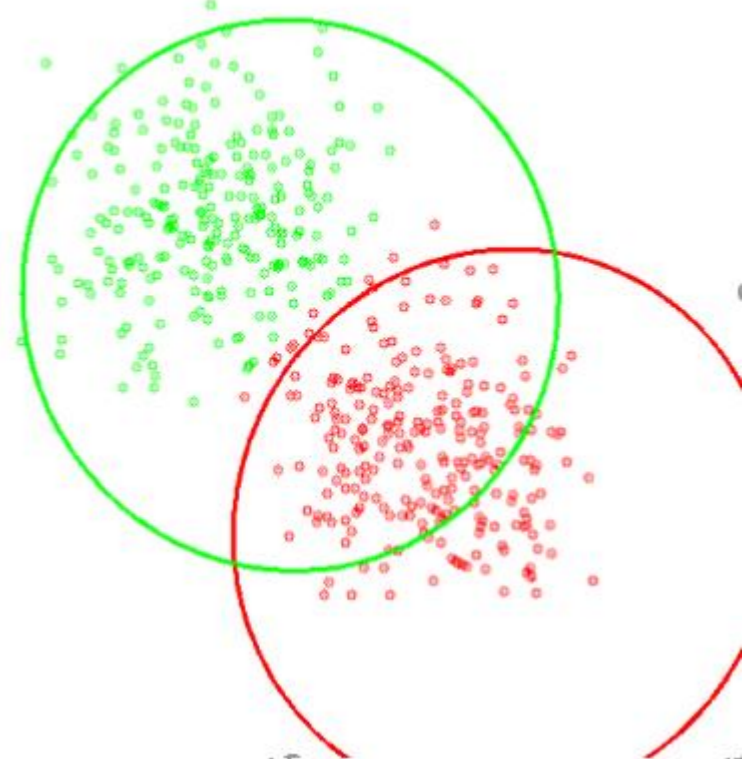
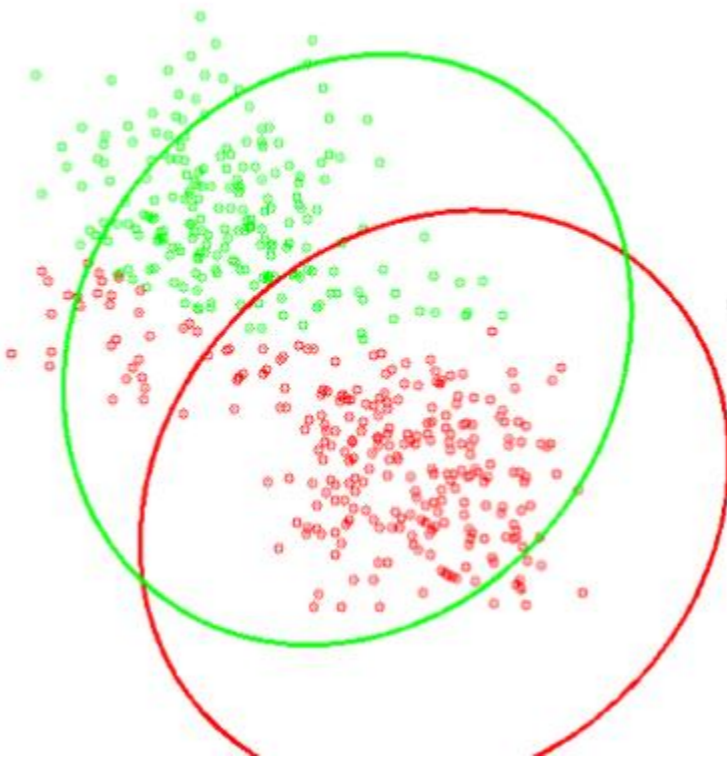
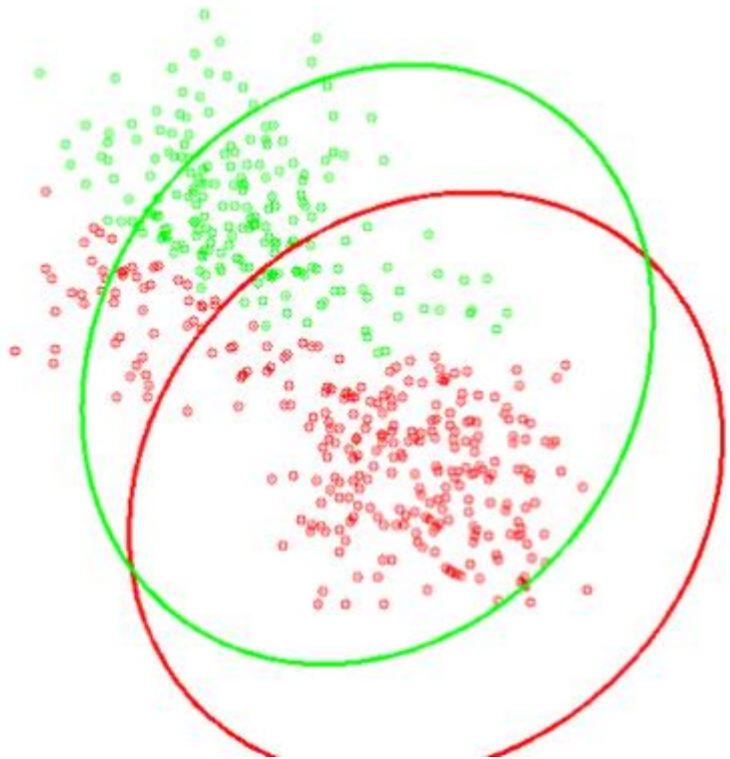
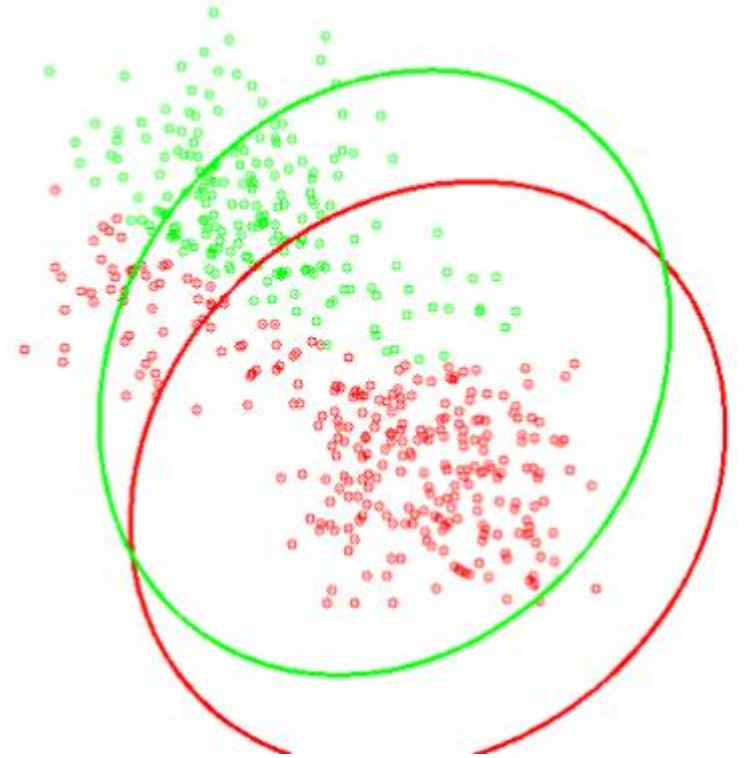
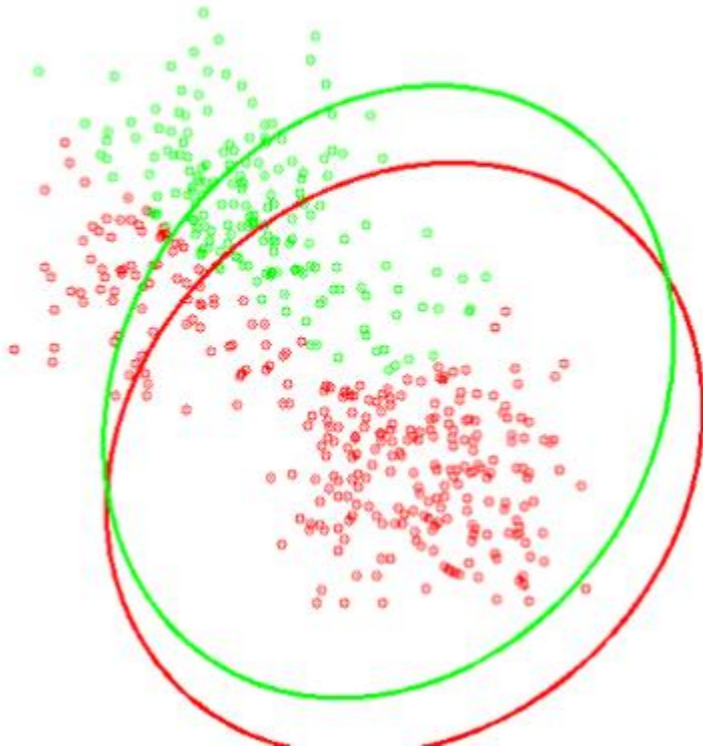
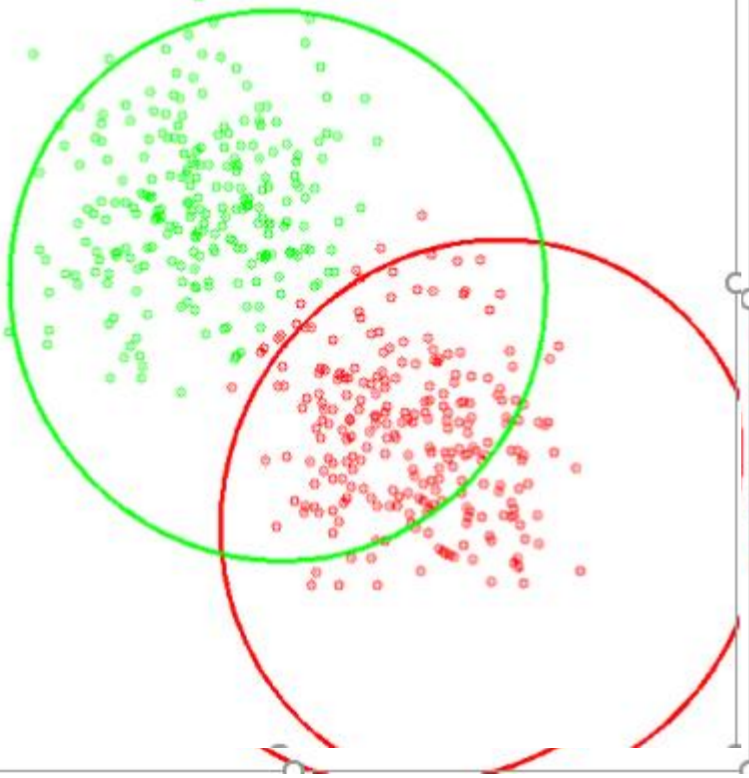
Let's run GMM on the first dataset

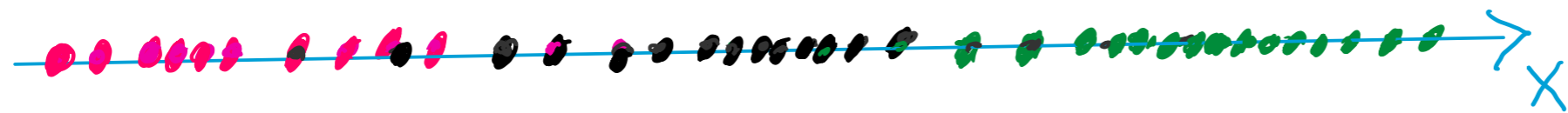
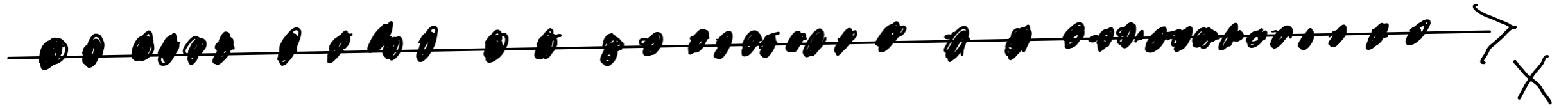




Let's do GMM on the second dataset

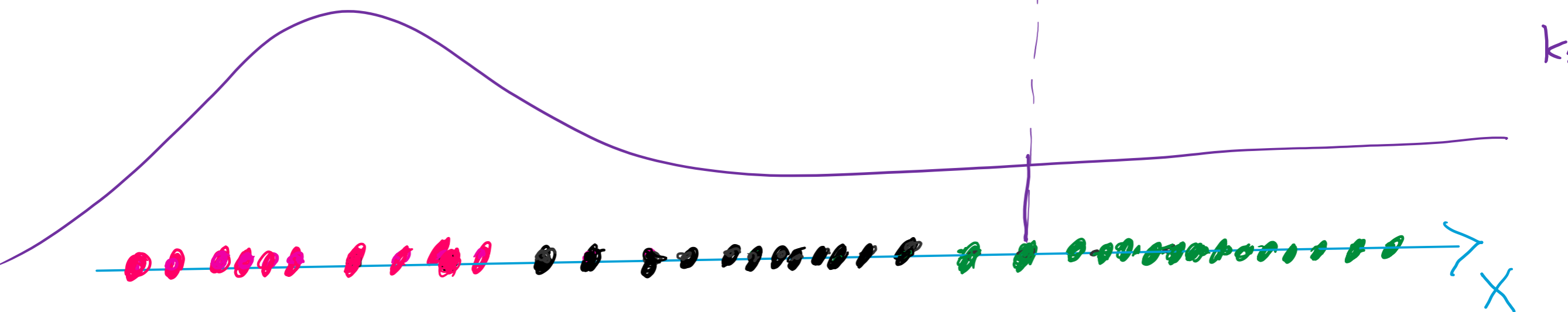
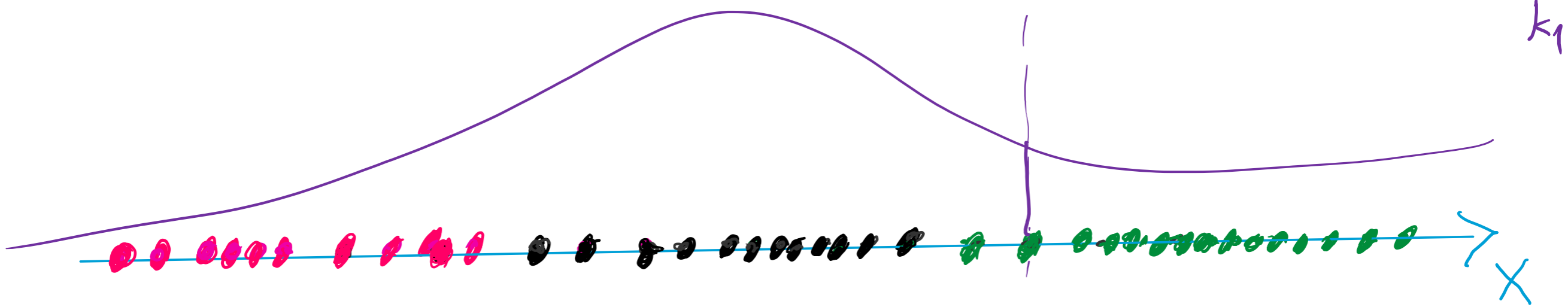
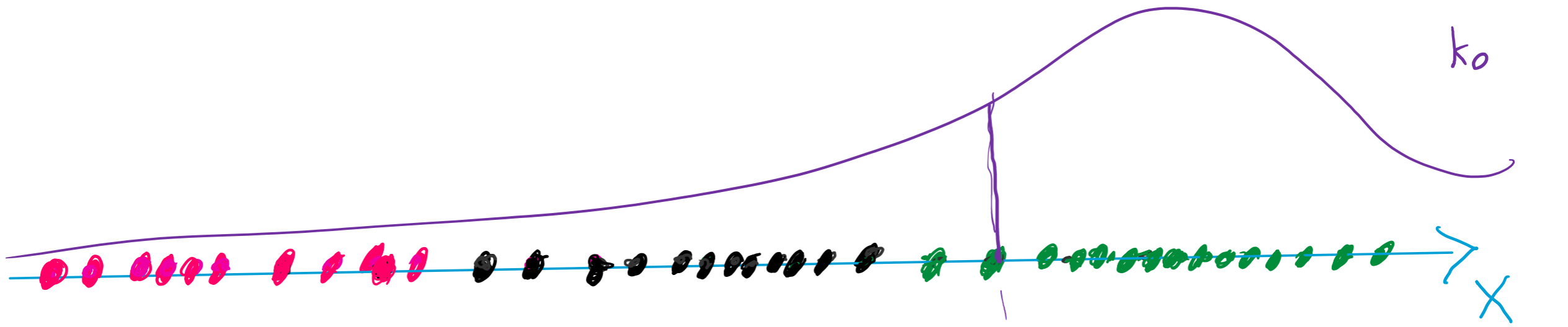
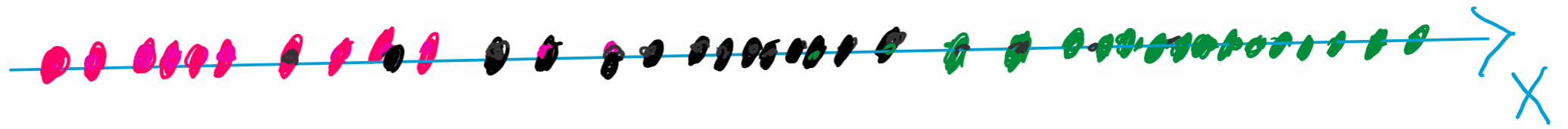






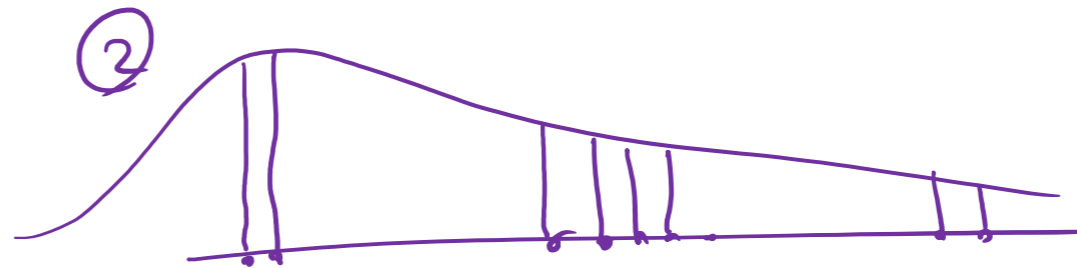
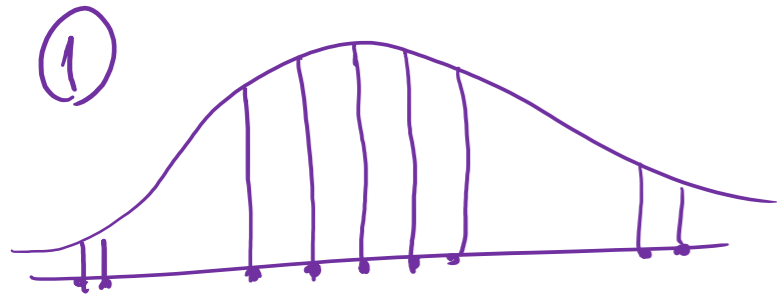
- cluster 1
- cluster 2
- cluster 3

Final Step



Gaussian Recap

$$f(x) = P(x) = (f(x|\theta)) = P(x|\theta) = (f(x|\mu, \sigma)) = N = N(x|\theta) = N(x|\mu, \sigma)$$



$$L(\theta|x) = \prod_{i=1}^N f(x^{(i)} | \mu, \sigma) \rightarrow \text{pdf} \rightsquigarrow \text{Gaussian distribution}$$

Likelihood function

$$l(\theta|x) = \log L(\theta|x) = \sum_{i=1}^N (\log f(x^{(i)} | \mu, \sigma))$$

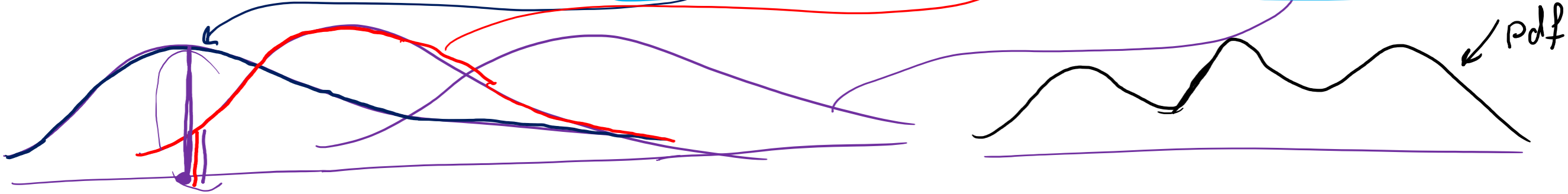
$$\frac{\partial l}{\partial \mu} = 0 \quad \frac{\partial l}{\partial \sigma} = 0$$

Mixture perspective – soft assignment

$$\pi_0 + \pi_1 + \pi_2 = 1 \quad \sum \pi_i = 1$$

Let's create a **SINGLE** pdf that combines all three Gaussians!!!!

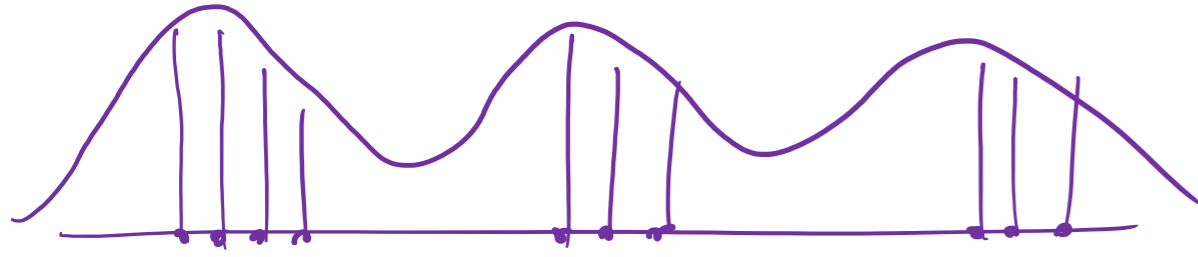
$$f(x) = f(x|\theta) = f(x|\mu, \sigma, \pi) = \frac{1}{3} f(x|\mu_0, \sigma_0) + \frac{1}{3} f(x|\mu_1, \sigma_1) + \frac{1}{3} f(x|\mu_2, \sigma_2)$$



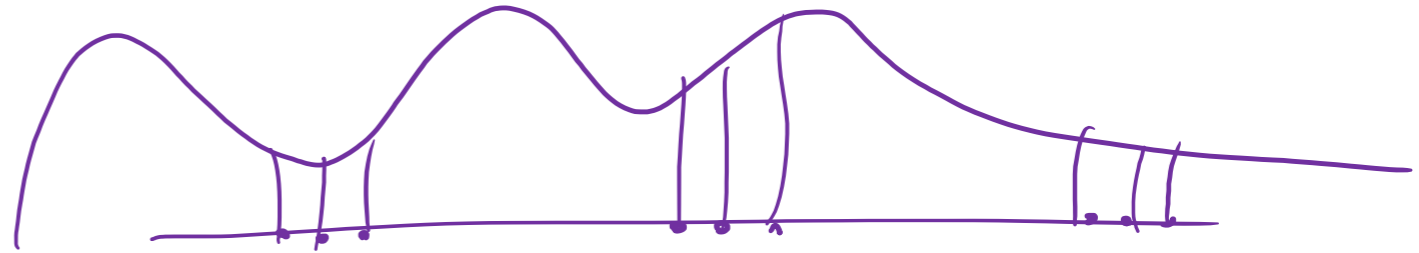
$$f(x) = f(x|\theta_0) + f(x|\theta_1) + f(x|\theta_2) \quad \Rightarrow \int f(x) dx = \underbrace{\int f(x|\theta_0) dx}_1 + \underbrace{\int f(x|\theta_1) dx}_1 + \underbrace{\int f(x|\theta_2) dx}_1$$

Constructing the log-likelihood of the mixture model

①



②



$$L(\theta | x) = L(x | \mu, \Sigma, \pi) = \prod_{i=1}^N f(x^{i3} | \mu, \Sigma, \pi)$$

Max

$$l(\theta | x) = \sum_{i=1}^N \log f(x^{i3} | \mu, \Sigma, \pi) = \sum_{i=1}^N \log \left[\pi_0 f_0(x | \mu_0, \Sigma_0) + \dots + \pi_k f_k(x | \mu_k, \Sigma_k) \right]$$

$$\frac{\partial l(\theta | x)}{\partial \pi} = 0$$

$$\frac{\partial l(\theta | x)}{\partial \Sigma} = 0$$

$$\frac{\partial l(\theta | x)}{\partial \mu} = 0$$

Using probability math notations: Hidden variables and responsibility

$K=3$

$$P(x) = P(x|\theta) = \underbrace{\pi_0}_{P(z_0)} \underbrace{f(x|\mu_0, \sigma_0)}_{P(x|z_0)} + \underbrace{\pi_1}_{P(z_1)} \underbrace{f(x|\mu_1, \sigma_1)}_{P(x|z_1)} + \pi_2 f(x|\mu_2, \sigma_2) + \dots + P(z_K)P(x|z_K)$$

\downarrow
 \downarrow
 \downarrow
 \downarrow
 \downarrow
 \downarrow
 N_0

$$P(x) = P(x|\theta) = P(z_0)P(x|z_0) + P(z_1)P(x|z_1) + P(z_2)P(x|z_2)$$

$$= P(x, z_0) + P(x, z_1) + P(x, z_2)$$

$$= \sum_K P(x, z_K) = P(x)$$

$\gamma = \tau = \text{Responsibility} =$

$k_0 \quad k_1 \quad k_2$
 $\begin{matrix} \circlearrowleft 0.7 \\ \circlearrowleft 0.2 \\ \circlearrowleft 0.1 \end{matrix} = 1$

$P(x|z) = N(x|\mu, \sigma)$

z_0

$P(x) = \frac{\pi_0 N_0}{a} + \frac{\pi_1 N_1}{b} + \frac{\pi_2 N_2}{c}$

$P(z_0|x^{\xi_13}) + P(z_1|x^{\xi_13}) + P(z_2|x^{\xi_13}) = 1$

$P(z_0|x^{\xi_13}) = \frac{P(z_0, x^{\xi_13})}{P(x^{\xi_13})} = \frac{P(x^{\xi_13}|z_0)P(z_0)}{P(x^{\xi_13}|z_0)P(z_0) + P(x^{\xi_13}|z_1)P(z_1) + P(x^{\xi_13}|z_2)P(z_2)}$

$\underbrace{P(x^{\xi_13}|z_0)}_{N_0} \underbrace{P(z_0)}_{\pi_0} \quad \underbrace{P(x^{\xi_13}|z_1)}_{N_1} \underbrace{P(z_1)}_{\pi_1} \quad \underbrace{P(x^{\xi_13}|z_2)}_{N_2} \underbrace{P(z_2)}_{\pi_2}$

$= \frac{a}{a+b+c}$

$P(z_1|x^{\xi_13}) = \frac{b}{a+b+c}$

$P(z_2|x^{\xi_13}) = \frac{c}{a+b+c}$

G Mixture Models

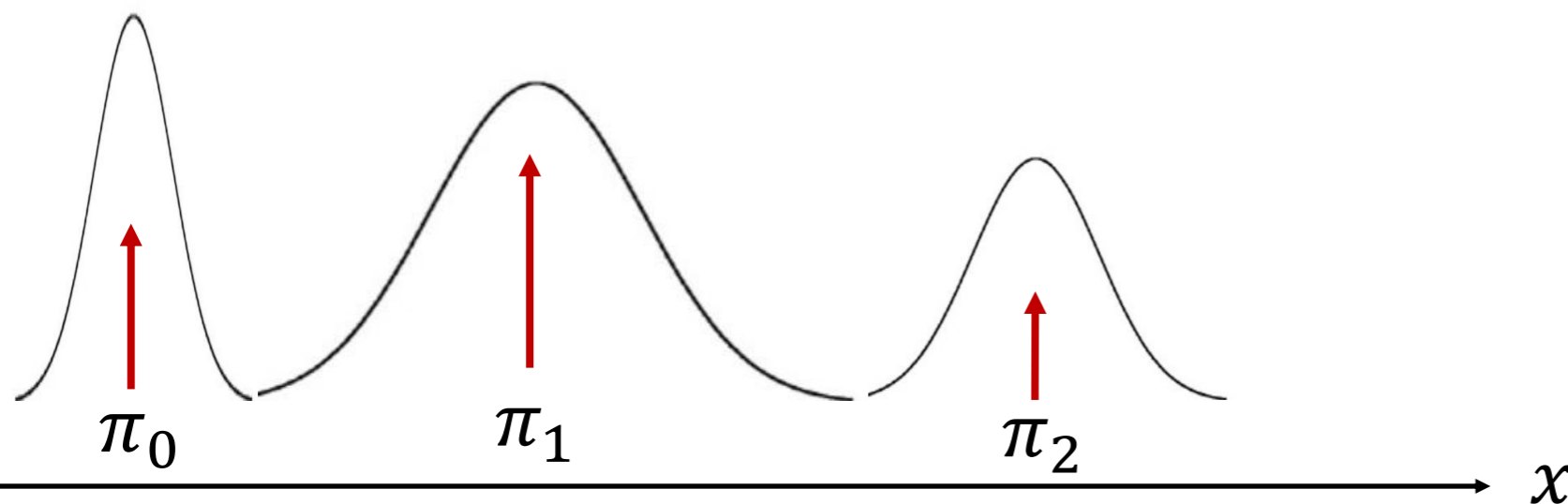
- Formally a Mixture Model is the weighted sum of a number of pdfs where the weights are determined by a distribution, π

$$p(x) = \pi_0 f_0(x) + \pi_1 f_1(x) + \pi_2 f_2(x) + \dots + \pi_k f_k(x)$$

where $\sum_{i=0}^k \pi_i = 1$

N_0

$$p(x) = \sum_{i=0}^k \pi_i f_i(x)$$



What is **f** in GMM?

Mixture Models are Generative

- Generative simply means dealing with joint probability $p(x, z)$

$$p(x) = \pi_0 f_0(x) + \pi_1 f_1(x) + \cdots + \pi_k f_k(x)$$

Let's say $f(\cdot)$ is a Gaussian distribution

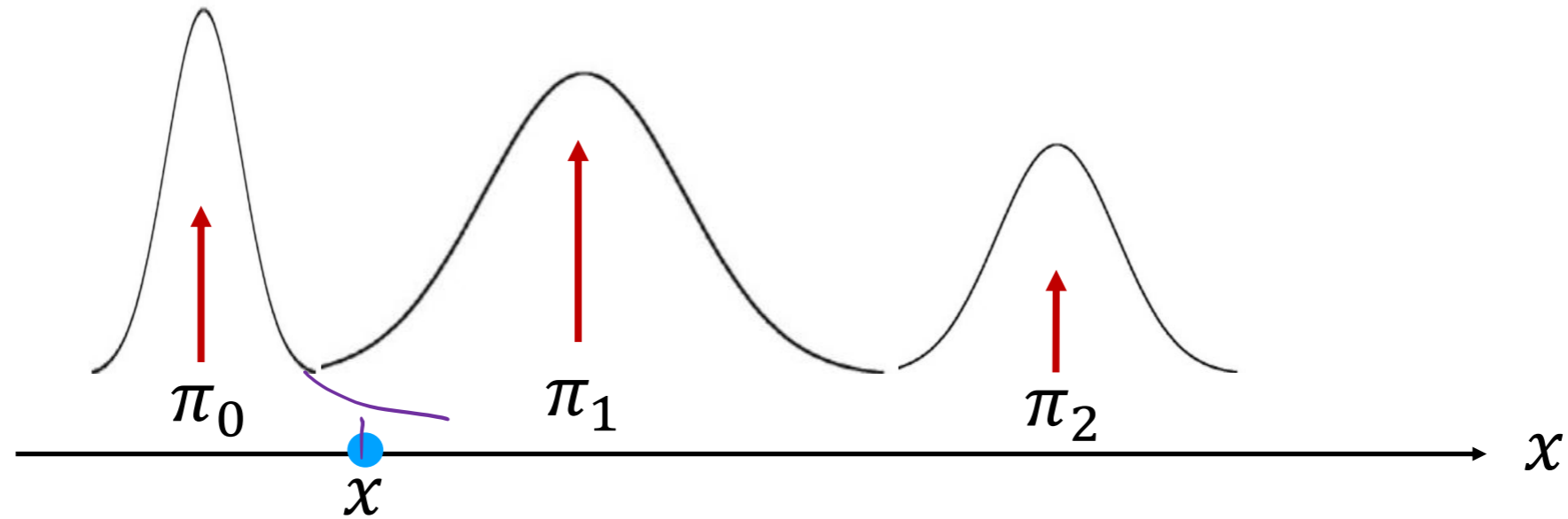
$$p(x) = \pi_0 N(X|\mu_0, \sigma_0) + \pi_1 N(X|\mu_1, \sigma_1) + \cdots + \pi_k N(X|\mu_k, \sigma_k)$$

$$p(x) = \sum_k N(x|\mu_k, \sigma_k) \pi_k$$

$$p(x) = \sum_k p(x|z_k) p(z_k) \quad z_k \text{ is component } k$$

$$p(x) = \sum_k p(x, z_k)$$

What is soft assignment?



$$\gamma = \begin{bmatrix} \circ & \circ & \circ \end{bmatrix}_{n \times k}$$

What is the probability of a datapoint x in each component?

How many components we have here? **3**

How many probability? **3**

What is the sum value of the 3 probabilities for each datapoint? **1**

Inferring Cluster Membership

- We have representations of the joint $p(x, z_{nk} | \theta)$ and the marginal, $p(x | \theta)$
- The conditional of $p(z_{nk} | x, \theta)$ can be derived using Bayes rule.
 - The **responsibility** that **a** mixture component takes for explaining an observation x .

$$\begin{aligned} \tau(z_k) = p(z_k | x) &= \frac{p(z_k) p(x | z_k)}{\sum_{j=1}^K p(z_j) p(x | z_j)} \\ &= \frac{\pi_k N(x | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x | \mu_j, \Sigma_j)} \end{aligned}$$

Handwritten annotations:

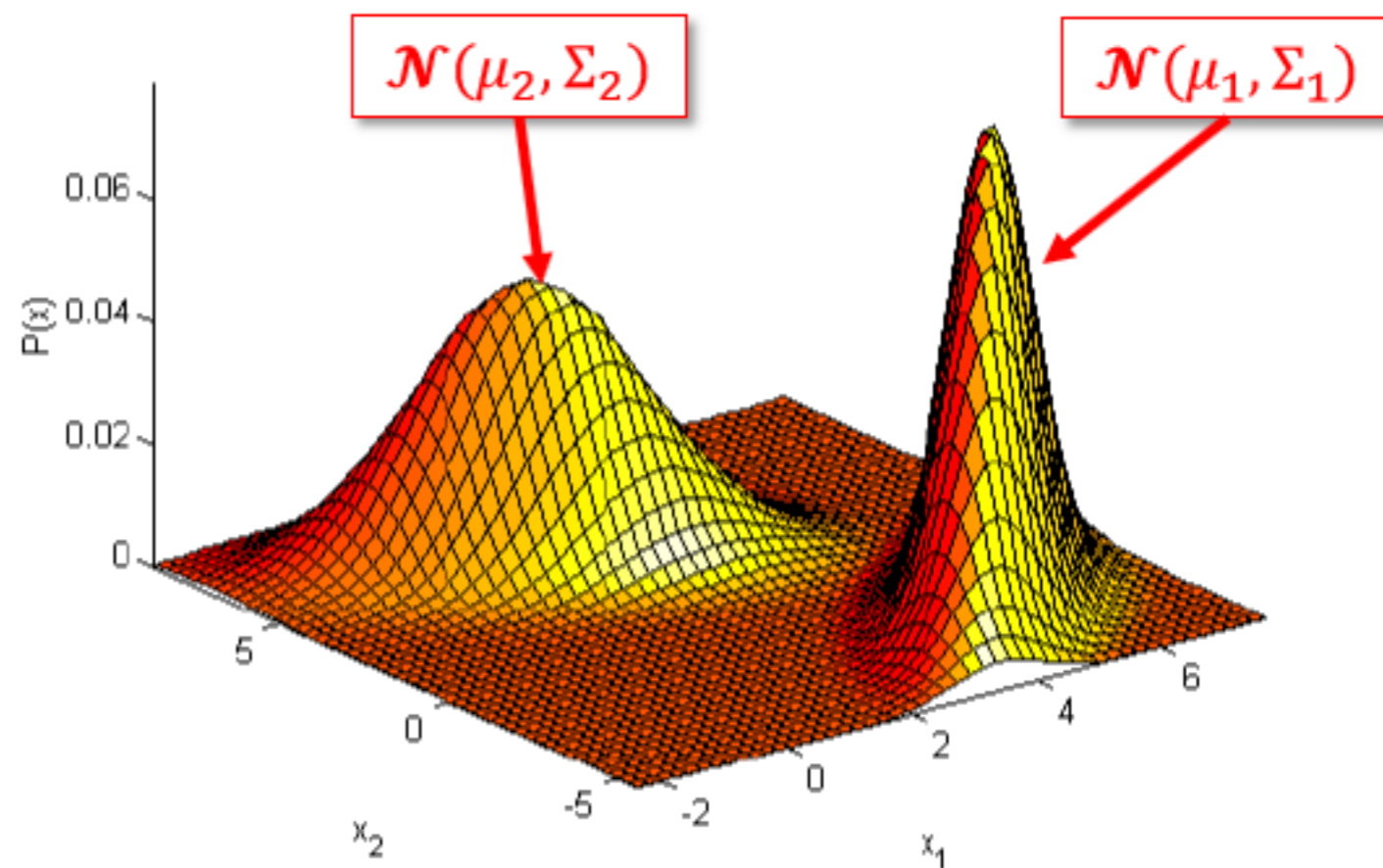
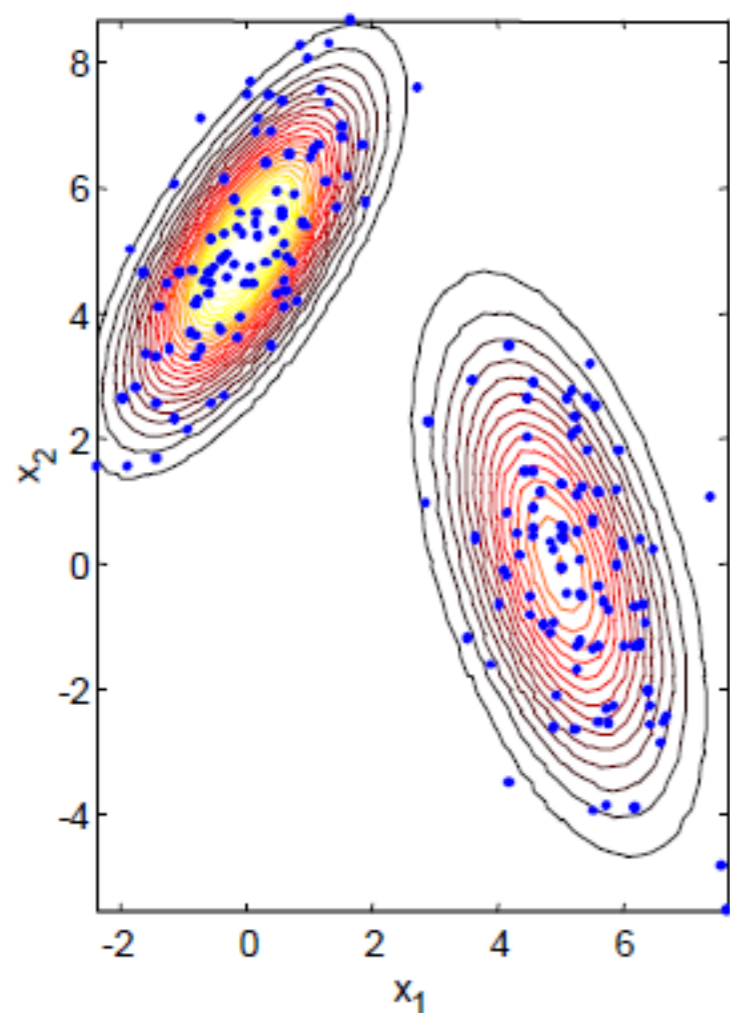
- A purple oval circles the left side of the equation: $\tau(z_k) = p(z_k | x)$.
- A purple oval circles the numerator of the second line: $\pi_k N(x | \mu_k, \Sigma_k)$. An arrow points from this oval to the label z_0 .
- A purple oval circles the denominator of the second line: $\sum_{j=1}^K \pi_j N(x | \mu_j, \Sigma_j)$. An arrow points from this oval to the label "Mixture model".
- The word "pdf" is written in purple below the denominator.
- A purple wavy line is drawn at the bottom of the page.

Why having “Latent variable”

- A variable can be unobserved (latent) because:
 - it is an imaginary quantity meant to provide some simplified and abstractive view of the data generation process.
 - e.g., speech recognition models, mixture models (soft clustering)...
 - it is a real-world object and/or phenomena, but difficult or impossible to measure
 - e.g., the temperature of a star, causes of a disease, evolutionary ancestors ...
 - it is a real-world object and/or phenomena, but sometimes wasn't measured, because of faulty sensors, etc.
- Discrete latent variables can be used to partition/cluster data into sub-groups.
- Continuous latent variables (factors) can be used for dimensionality reduction (factor analysis, etc).

How about GMM for multimodal distribution?

- What if we know the data consists of a few Gaussians
- What if we want to fit parametric models



Gaussian Mixture Model

- A density model $p(X)$ may be multi-modal: model it as a mixture of uni-modal distributions (e.g. Gaussians)

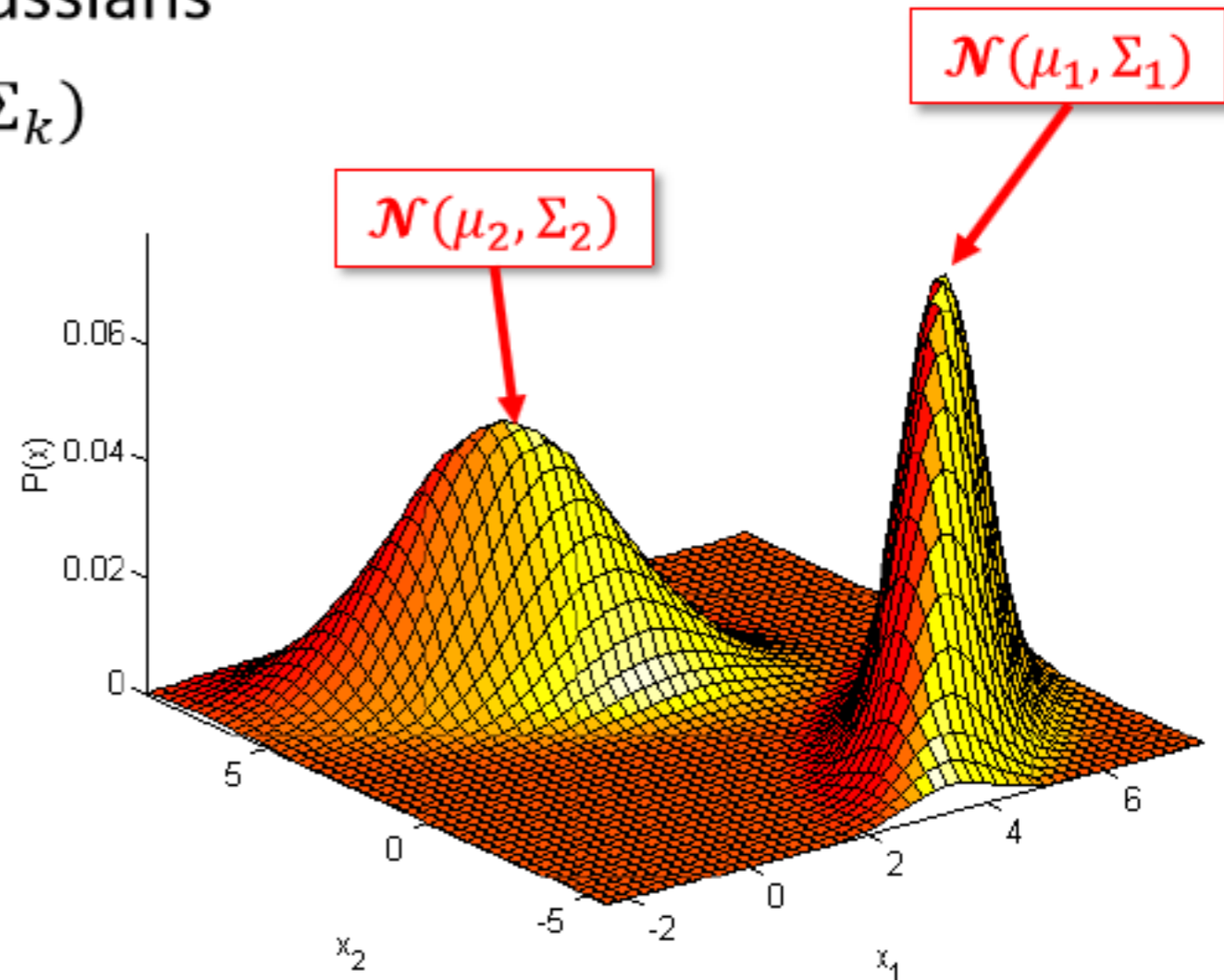
- Consider a mixture of K Gaussians

- $p(X) = \sum_{k=1}^K \pi_k \mathcal{N}(X|\mu_k, \Sigma_k)$

mixing
proportion

mixture
Component

- Learn $\pi_k \in (0,1), \mu_k, \Sigma_k$;



What are GMM parameters?

Mean μ_k

Variance σ_k

Size π_k

Marginal probability distribution

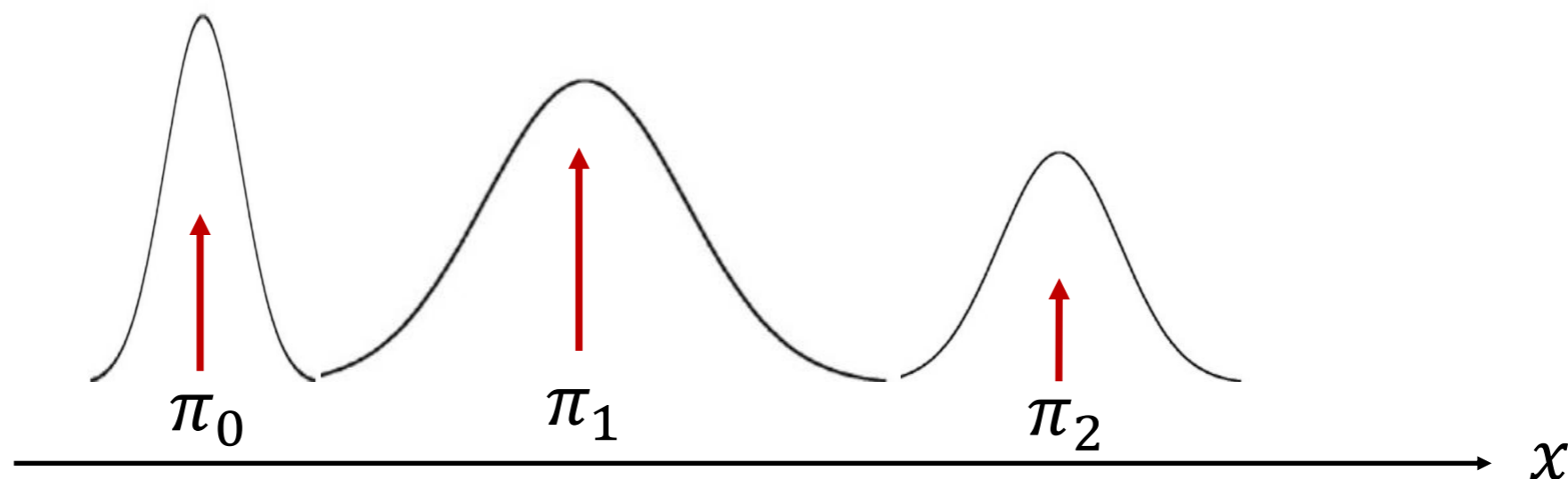
$$p(x|\theta) = \sum_k p(x, z_k|\theta) = \sum_k \underbrace{p(x|z_k, \theta)}_{\substack{N(\mu_k, \sigma_k) \\ f_k(x)}} \underbrace{p(z_k|\theta)}_{\pi_k} = \sum_k N(x|\mu_k, \sigma_k)\pi_k$$

$$p(z_k|\theta) = \pi_k$$

Select a mixture component with probability π

$$p(x|z_k, \theta) = N(x|\mu_k, \sigma_k)$$

Sample from that component's Gaussian



Parameters' definition

- Purpose: GMM is a clustering algorithm derived from probabilistic theory that uses soft-assignment, meaning that data points have probability of being associated/generated from K gaussians/clusters. This is as opposed to K -means where data points definitively are either from a cluster or they're not.
- Gaussian Parameters
 - μ : Mean of each gaussian, can be compared to the K -means cluster centers
 - Σ : Covariance matrix of each gaussian, which represents how dimensions vary between each other. If it's the covariance of a specific dimension with itself, it is just the standard deviation of that variable/dimension. This is in a $D \times D$ matrix (for each gaussian) and every element $\Sigma_{i,j}$ represents the covariance of dimension i with j . If you assume the dimensions are independent, then only the diagonals are non-zero.
 - z_{nk} : Latent variable which isn't explicitly known, but tells us which gaussian each datapoint was generated from. z is binary, it either's 1 (point x came from gaussian k) or 0 (point x did not come from gaussian k)
 - $p(z_{nk}) = \pi$: Mixing proportions/weights, which represent the fraction of data points that are generated from/associated with each gaussian. These sum to 1.
- $N(X_n | \mu_k, \Sigma_k)$: This term is the probability of some data point X_n occurring based on the assumption that it is generated from gaussian k . Mathematically, this is equal to the likelihood ($p(x|z, \mu_k, \Sigma_k)$). Multiply this with the mixing weight π , and you get the joint distribution $p(x=X_n, z=k)$.
- $P(X)$: If you sum up all the $N(X_n | \mu_k, \Sigma_k) * \pi$ terms for each cluster, you get the probability of the entire data set occurring.
- $\gamma(z_{nk})$ or $P(z_{nk} | x_n)$: We call this term the "responsibility." It is the probability of z for "A" data point x , meaning that this is probability that point n is generated from gaussian k normalized by $P(X)$, the probability of the entire data set occurring.

Well, we don't know π_k, μ_k, Σ_k
 What should we do?

We use a method called "Maximum Likelihood Estimation" (MLE) to solve the problem.

$$L(\theta | x) = \prod_{i=1}^N p(x^{(i)} | \theta) \quad \text{max} \quad \ell(\theta | x) = \sum \log p(x^{(i)} | \theta)$$

pdf $p(x) = p(x|\theta) = \sum_k p(x, z_k | \theta) = \sum_k p(z_k | \theta) p(x | z_k, \theta) = \sum_{k=0}^K \pi_k N(x | \mu_k, \Sigma_k)$

Let's identify a likelihood function, why?

Because we use likelihood function to optimize the probabilistic model parameters!

$$\arg \max p(x|\theta) = p(x|\pi, \mu, \Sigma) = \prod_{n=1}^N p(x^{\{n\}} | \theta) = \prod_{n=1}^N \sum_{k=0}^K \pi_k N(x^{\{n\}} | \mu_k, \Sigma_k)$$

$$\arg \max p(x) = p(x|\pi, \mu, \Sigma) = \prod_{n=1}^N p(x^{\{n\}}|\theta) = \prod_{n=1}^N \sum_{k=0}^K \pi_k N(x^{\{n\}}|\mu_k, \Sigma_k)$$

$$\ln[p(x)] = \ln[p(x|\pi, \mu, \Sigma)]$$

- As usual: Identify a likelihood function

μ, Σ, π

$\mu_0 \leftarrow \mu_0 + \alpha$

$\Sigma_0 \leftarrow \dots$

$\pi_0 \leftarrow \dots$

$$\frac{\partial \ln}{\partial \theta} \ln p(x|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n|\mu_k, \Sigma_k) \right\}$$

pdf

$\mu_1 \leftarrow \mu_{1+\dots}$

- And set partials to zero...

Maximum Likelihood of a GMM

- Optimization of means.

Objective function

$$\ln p(x|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$$

vb

$$\frac{\partial \ln p(x|\pi, \mu, \Sigma)}{\partial \mu_k} = \sum_{n=1}^N \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n | \mu_j, \Sigma_j)} \Sigma_k^{-1} (x_n - \mu_k) = 0$$

$$= \sum_{n=1}^N \tau(z_{nk}) \Sigma_k^{-1} (x_n - \mu_k) = 0$$

Concave
lb

$$\mu_k = \frac{\sum_{n=1}^N \tau(z_{nk}) x_n}{\sum_{n=1}^N \tau(z_{nk})}$$

Maximum Likelihood of a GMM

- Optimization of covariance

$$\ln p(x|\pi, \mu, \Sigma) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k N(x_n | \mu_k, \Sigma_k) \right\}$$

$$\Sigma_k = \frac{1}{\sum_{n=1}^N \tau(z_{nk})} \sum_{n=1}^N \tau(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T$$

Maximum Likelihood of a GMM

- Optimization of mixing term

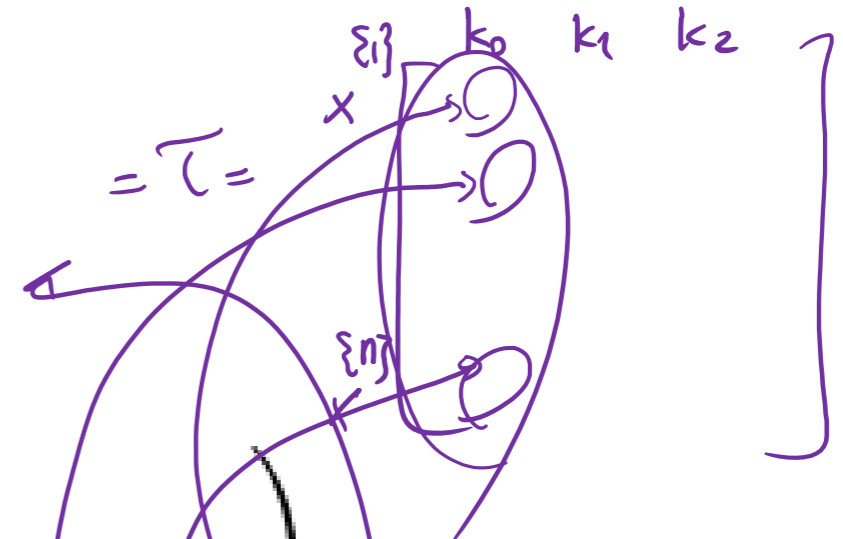
$$\pi_0 = \dots$$

$$\ln p(x|\pi, \mu, \Sigma) + \lambda \left(\sum_{k=1}^K \pi_k - 1 \right)$$

$$0 = \sum_{n=1}^N \frac{N(x_n|\mu_k, \Sigma_k)}{\sum_j \pi_j N(x_n|\mu_j, \Sigma_j)} + \lambda$$

$$\pi_k = \frac{\sum_{n=1}^N \tau(z_{nk})}{N}$$

$$= \frac{\sum_{n=1}^N \tau(z_{n0})}{N} = \frac{\tau(z_{10}) + \tau(z_{20}) + \dots + \tau(z_{N0})}{N}$$



MLE of a GMM

$$\mu_k = \frac{\sum_{n=1}^N \tau(z_{nk}) x_n}{N_k}$$

$$\Sigma_k = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk}) (x_n - \mu_k)(x_n - \mu_k)^T$$

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

$$N_k = \sum_{n=1}^N \tau(z_{nk})$$

Outline

- Overview
- Gaussian Mixture Model
- The Expectation-Maximization Algorithm ←

EM for GMMs

- E-step: Evaluate the Responsibilities

$$\tau(z_k) = \frac{\pi_k N(x_n | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j N(x_n | \mu_j, \Sigma_j)}$$

EM for GMMs

- M-Step: Re-estimate Parameters

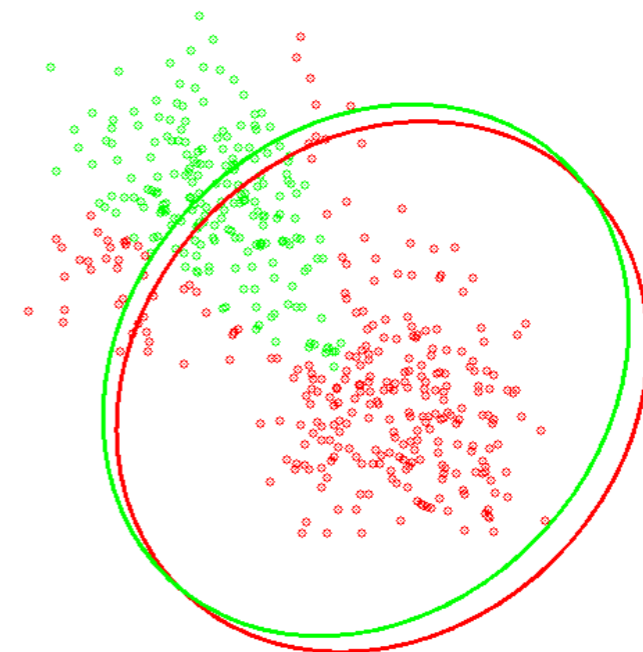
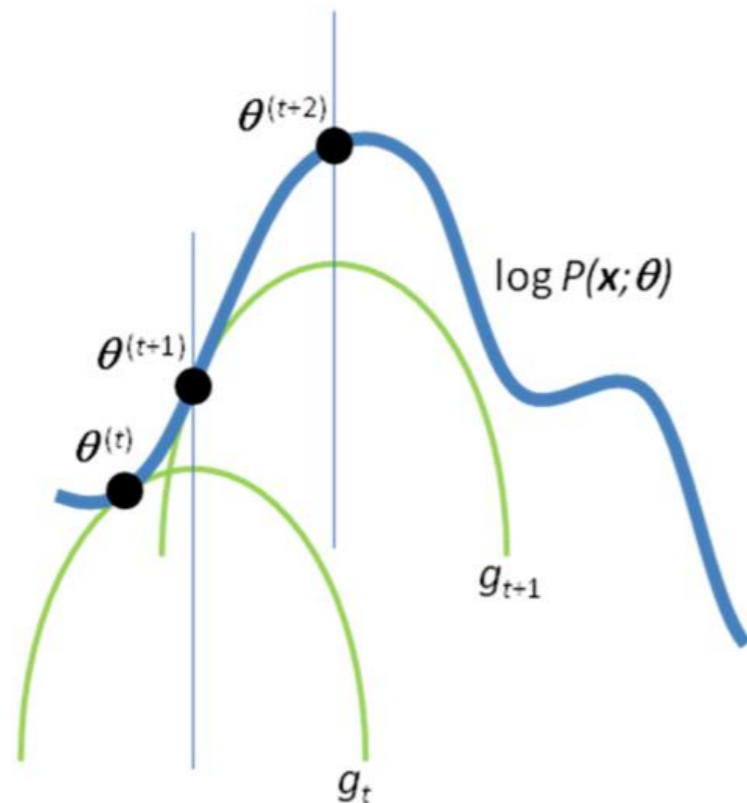
$$\mu_k^{new} = \frac{\sum_{n=1}^N \tau(z_{nk}) x_n}{N_k}$$

$$\Sigma_k^{new} = \frac{1}{N_k} \sum_{n=1}^N \tau(z_{nk}) (x_n - \mu_k^{new})(x_n - \mu_k^{new})^T$$

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

Expectation Maximization

- Expectation Maximization (EM) is a general algorithm to deal with hidden variables.
- Two steps:
 - E-Step: Fill-in hidden values using inference
 - M-Step: Apply standard MLE method to estimate parameters
- EM always converges to a local minimum of the likelihood.

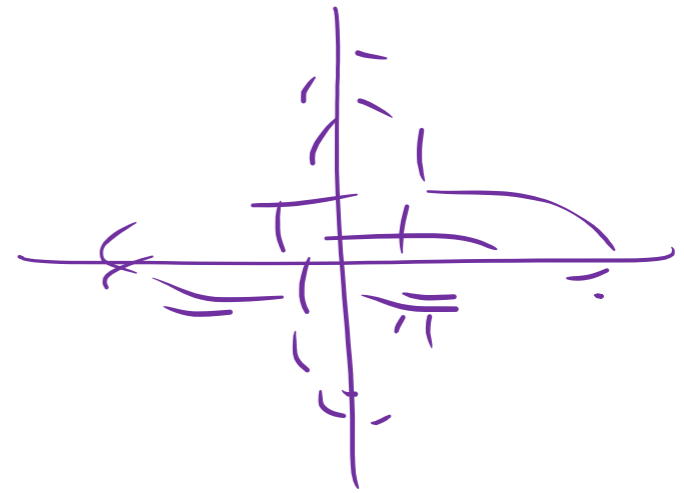


EM for Gaussian Mixture Model: Example

$$A = \begin{bmatrix} a & 0 \\ 0 & b \end{bmatrix} \quad A^{-1} = \begin{bmatrix} \frac{1}{a} & 0 \\ 0 & \frac{1}{b} \end{bmatrix}$$

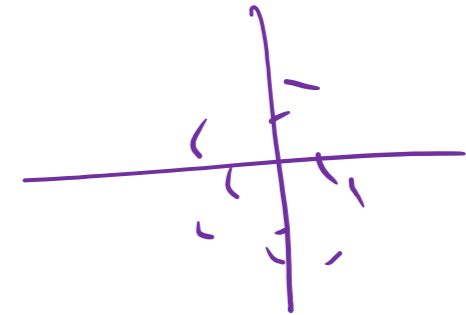
covariance_type="diag"

$$\begin{bmatrix} \sigma_h^2 & 0 \\ 0 & \sigma_w^2 \end{bmatrix}$$



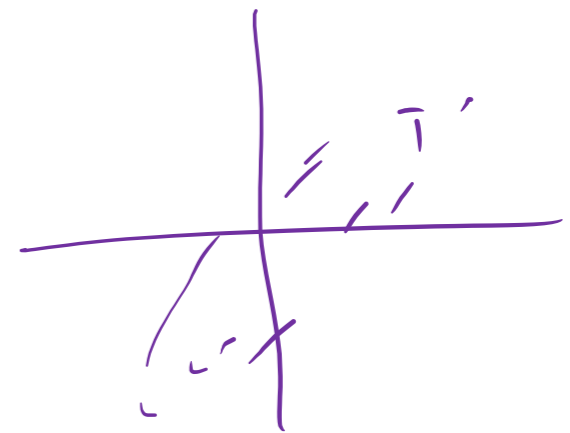
covariance_type="spherical"

$$\begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

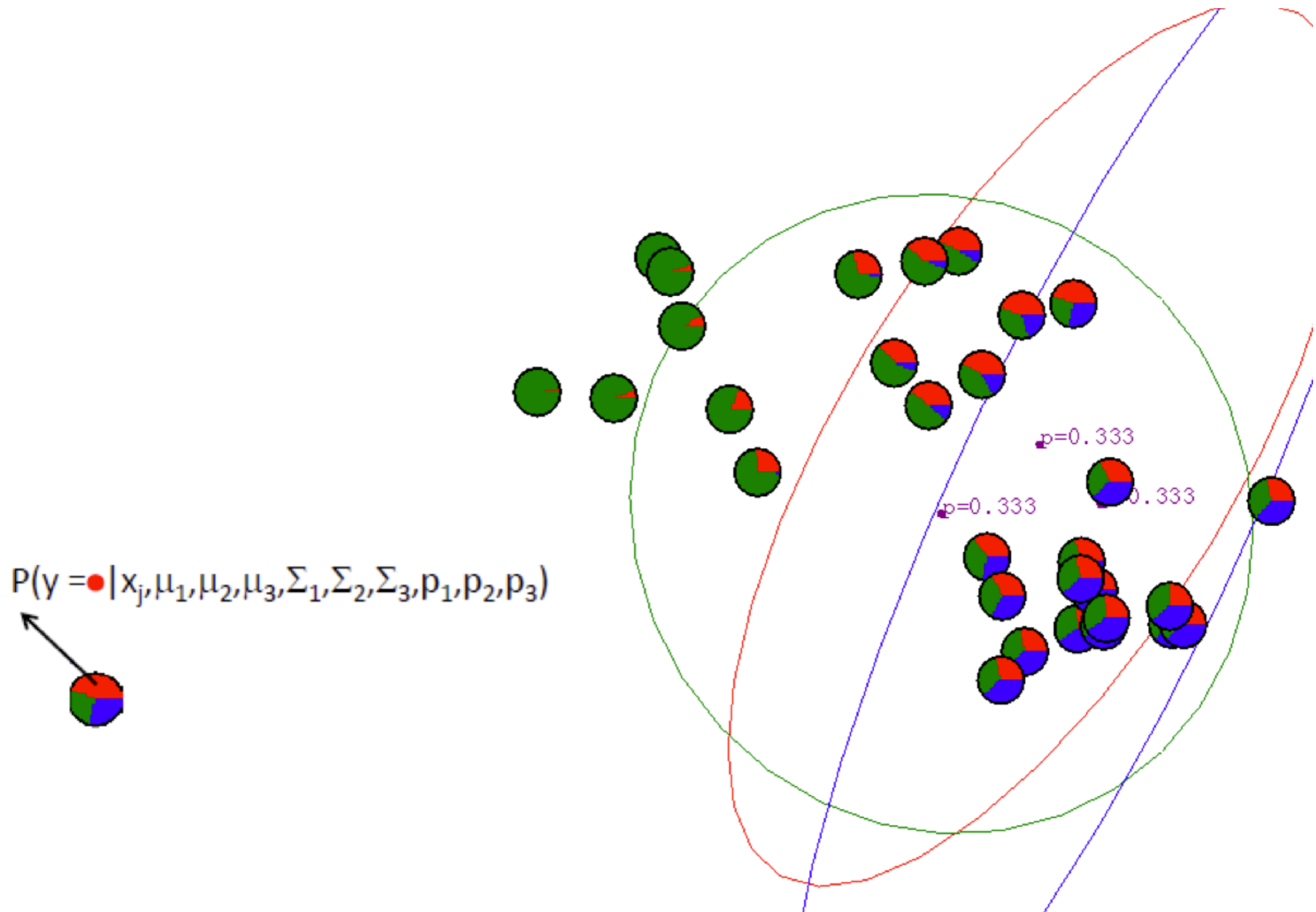


covariance_type="full"

$$\begin{bmatrix} \sigma_h^2 & \sigma_{hw} \\ \sigma_{wh} & \sigma_w^2 \end{bmatrix}$$

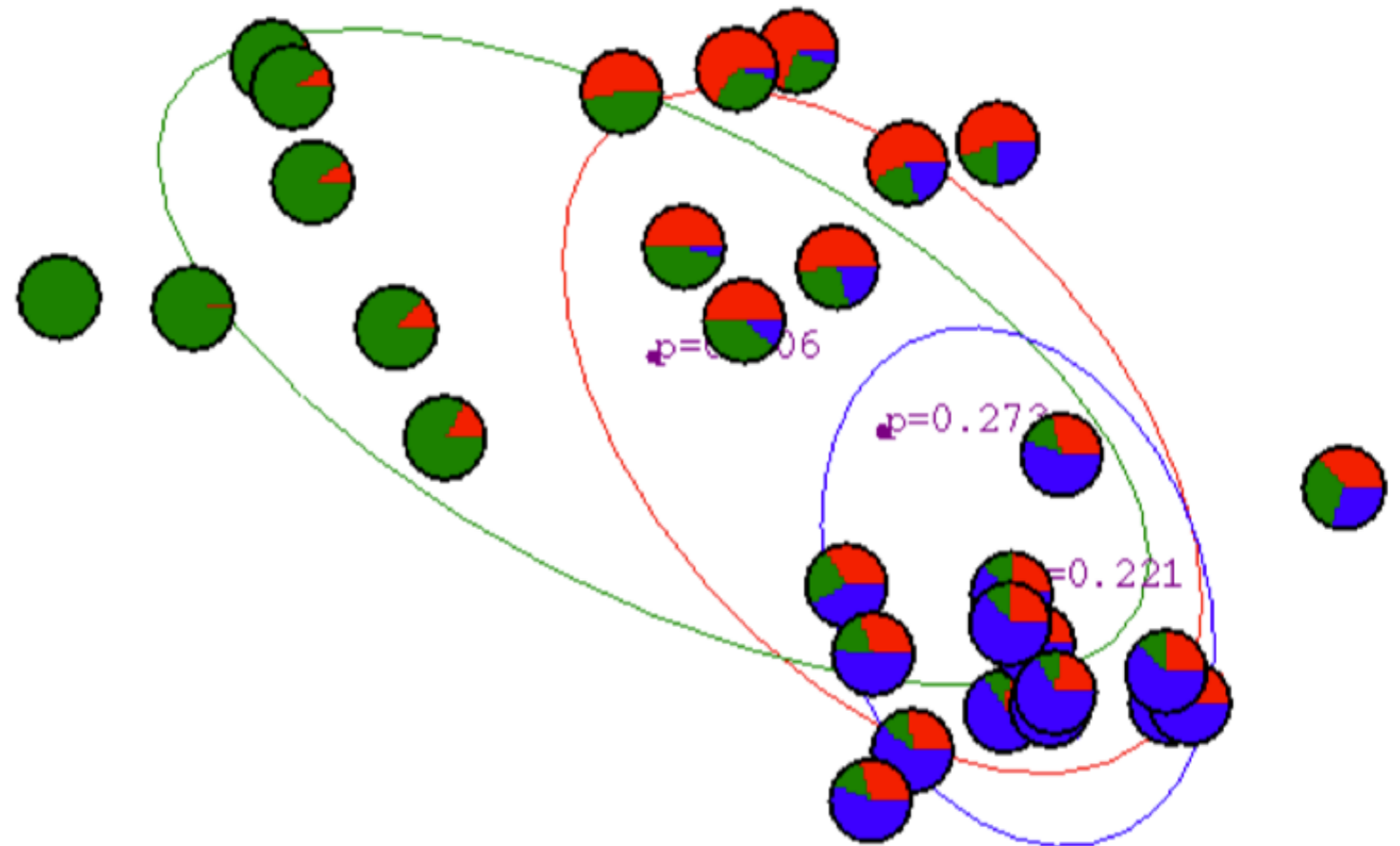


EM for Gaussian Mixture Model:



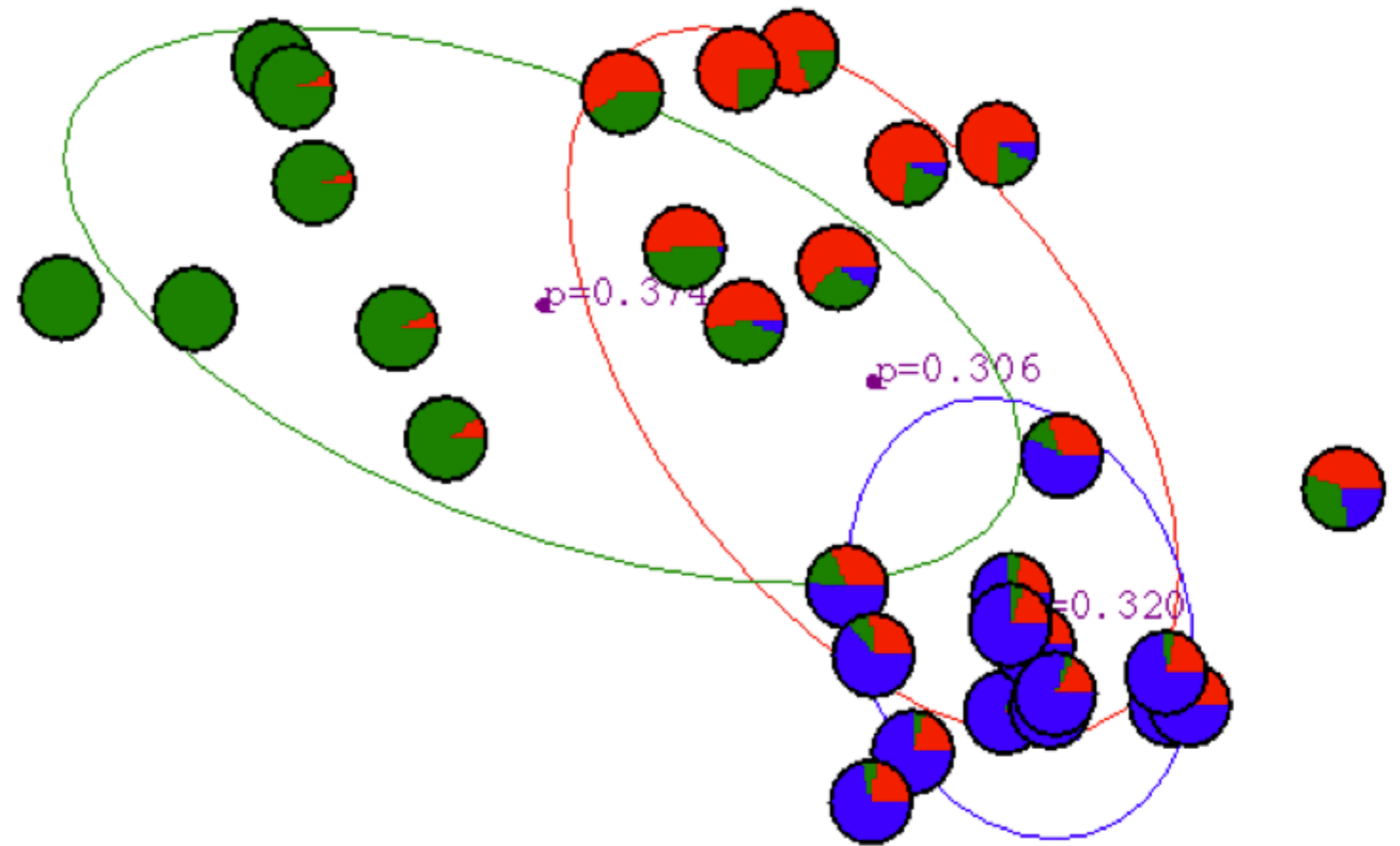
EM for Gaussian Mixture Model: Example

After 1st iteration



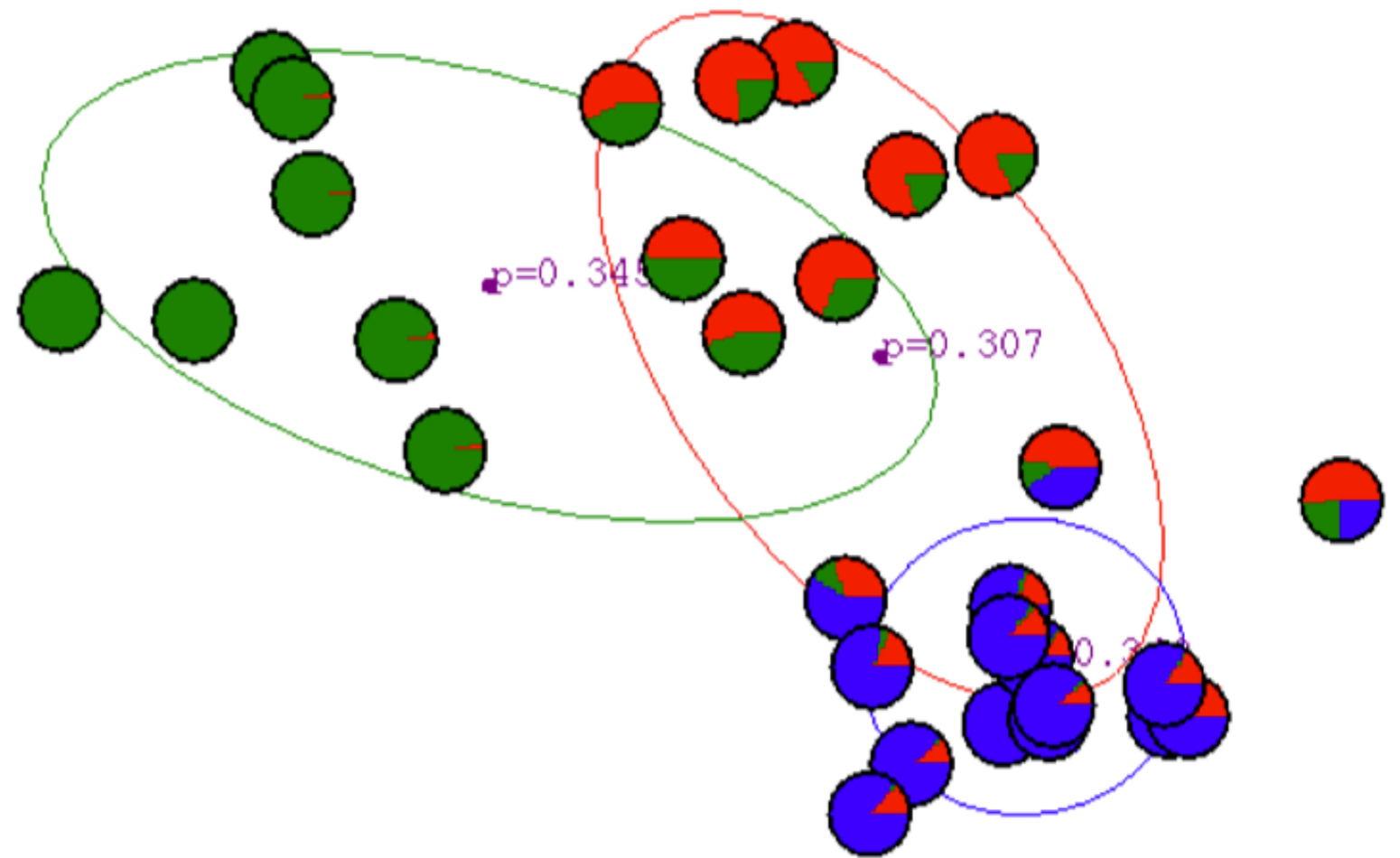
EM for Gaussian Mixture Model: Example

After 2nd iteration



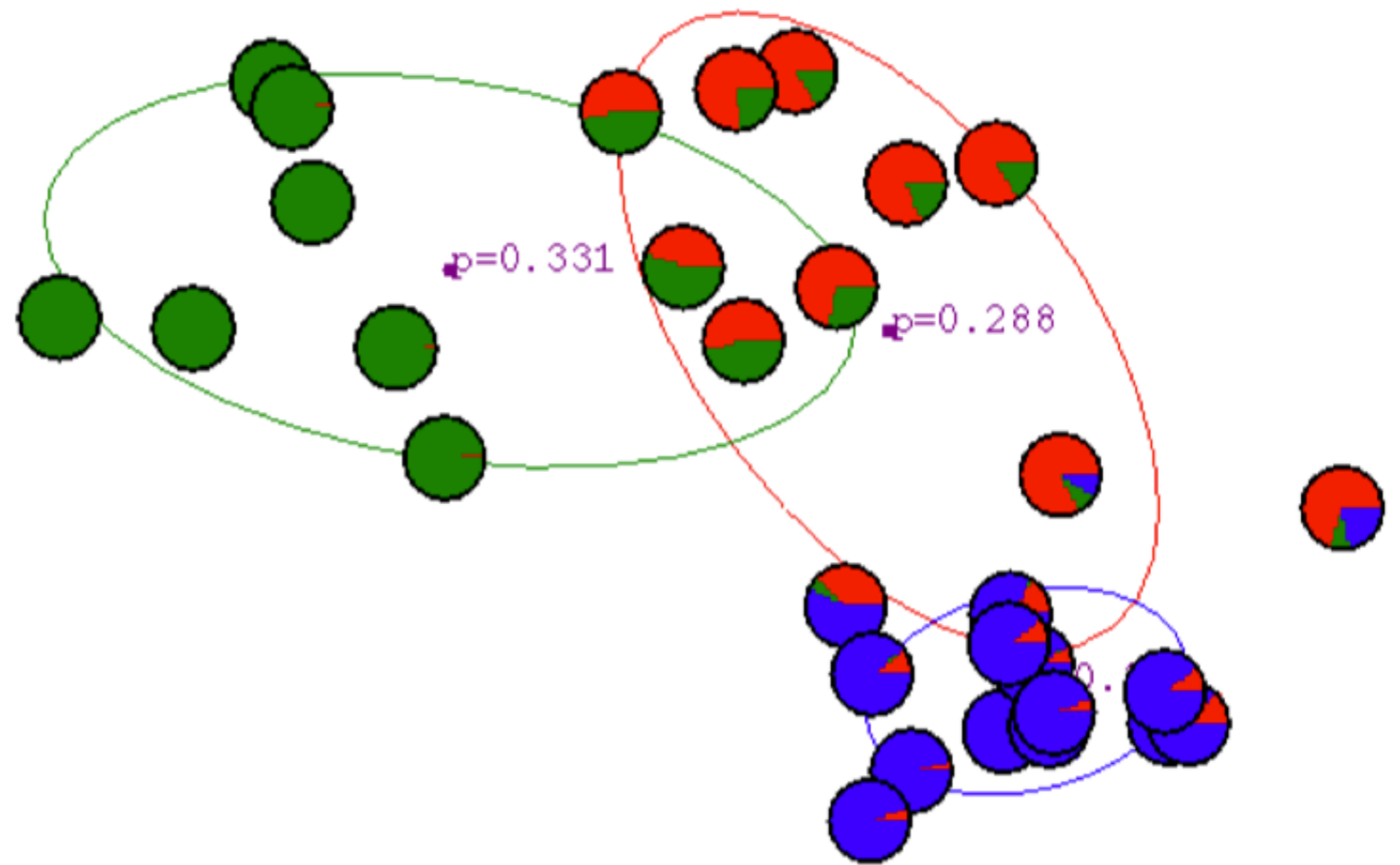
EM for Gaussian Mixture Model: Example

After 3rd iteration



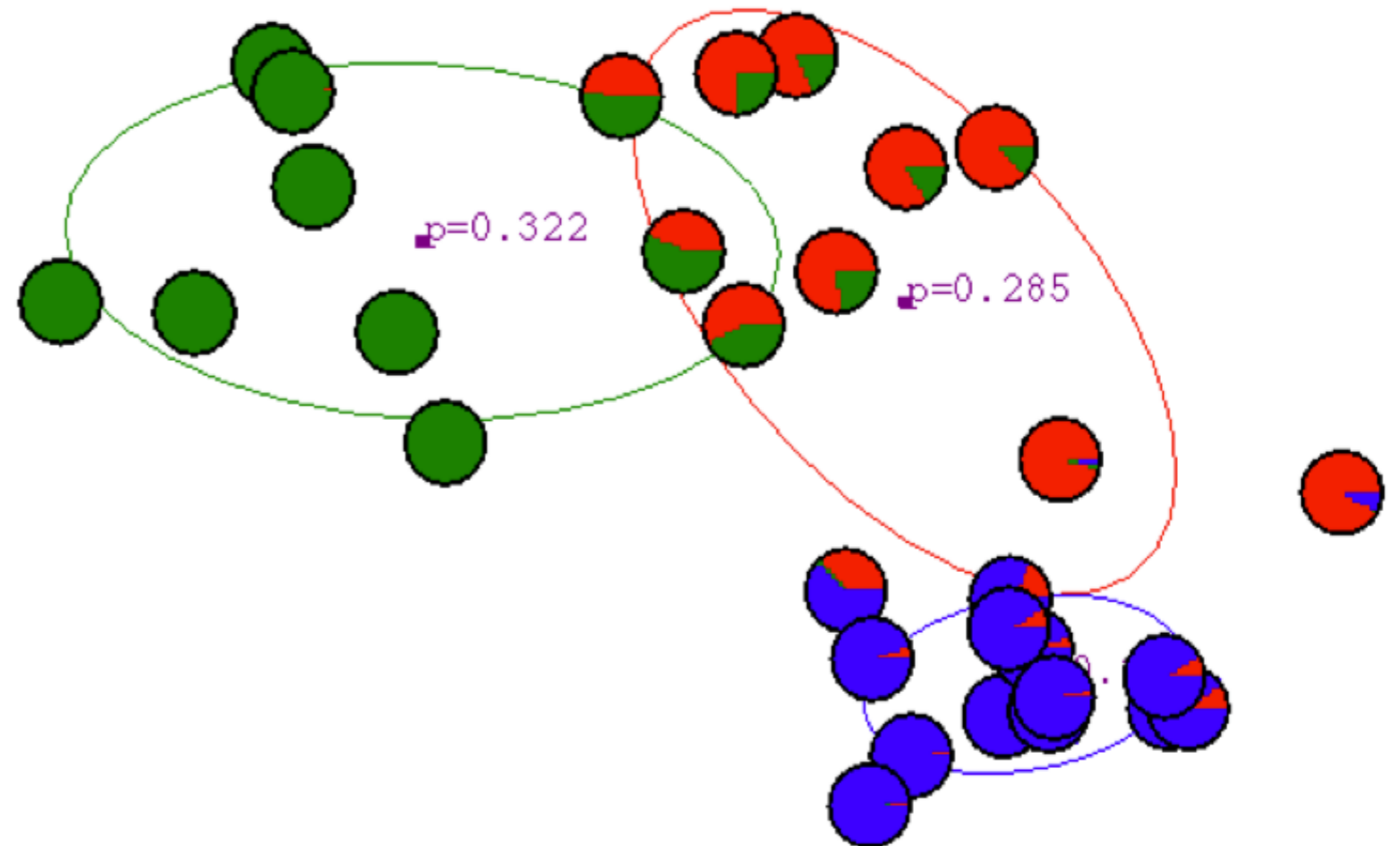
EM for Gaussian Mixture Model: Example

After 4th iteration



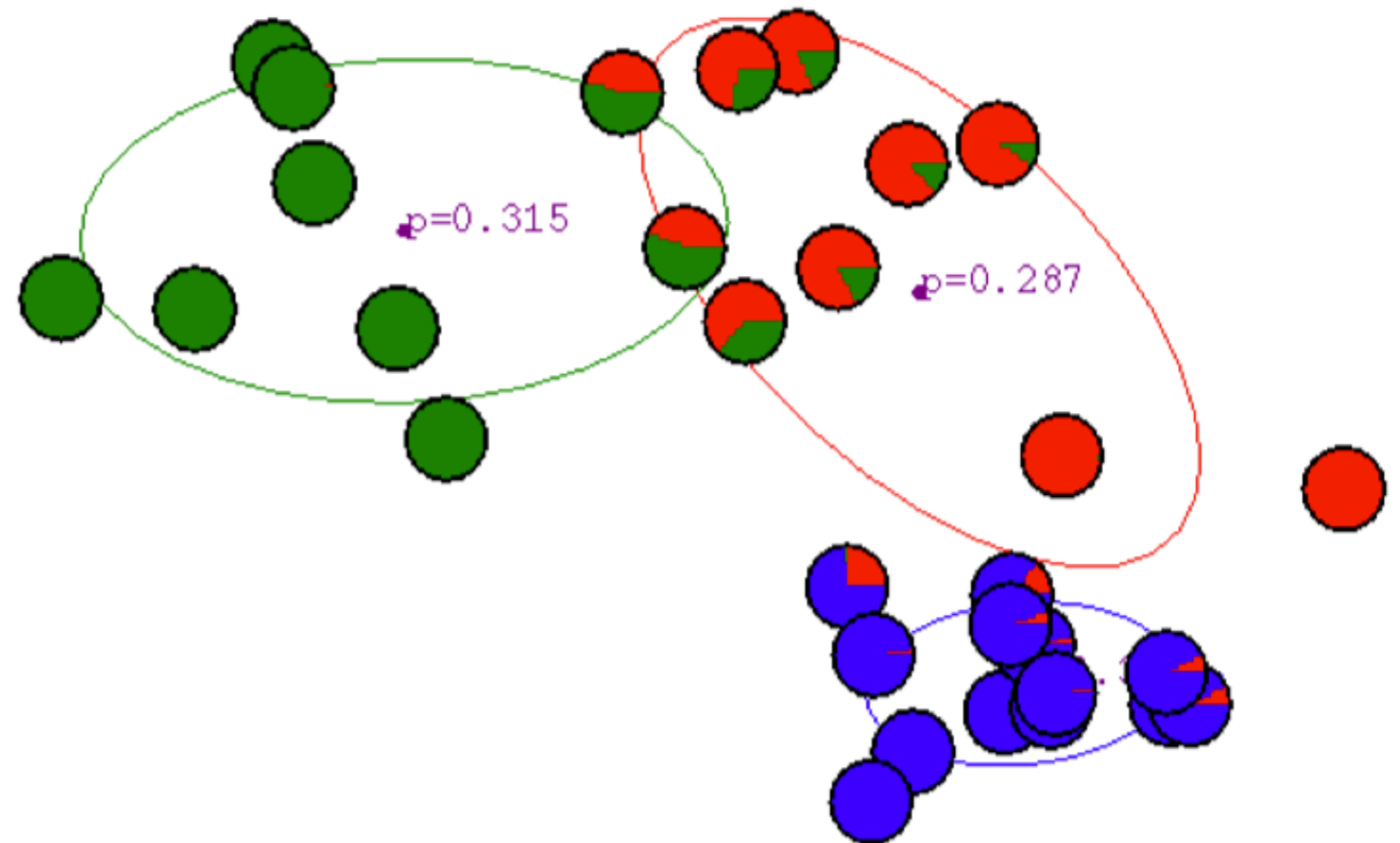
EM for Gaussian Mixture Model: Example

After 5th iteration



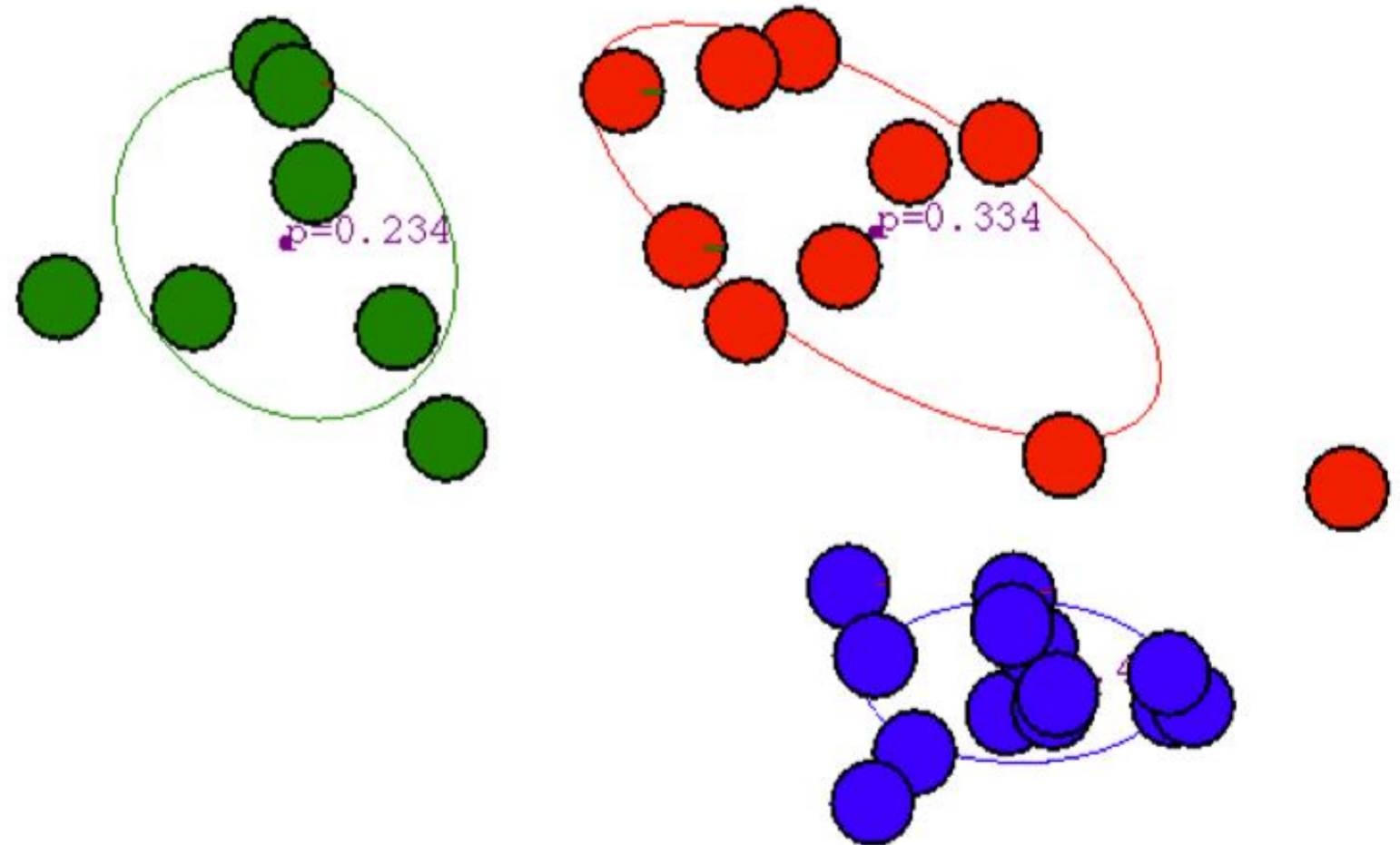
EM for Gaussian Mixture Model: Example

After 6th iteration



EM for Gaussian Mixture Model: Example

After 20th iteration



EM Algorithm for GMM (matrix form)

Given a Gaussian mixture model, the goal is to maximize the likelihood function with respect to the parameters comprising the means and covariances of the components and the mixing coefficients).

1. Initialize the means μ_j , covariances Σ_j and mixing coefficients π_j , and evaluate the initial value of the log likelihood.
2. **E step.** Evaluate the responsibilities using the current parameter values

$$\gamma(z_k) = \frac{\pi_k \mathcal{N}(\mathbf{x} | \mu_k, \Sigma_k)}{\sum_{j=1}^K \pi_j \mathcal{N}(\mathbf{x} | \mu_j, \Sigma_j)}$$

EM Algorithm for GMM (matrix form)

3. **M step.** Re-estimate the parameters using the current responsibilities

$$\mu_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) \mathbf{x}_n}{\sum_{n=1}^N \gamma(z_{nk})}$$

$$\Sigma_k = \frac{\sum_{n=1}^N \gamma(z_{nk}) (\mathbf{x}_n - \mu_k)(\mathbf{x}_n - \mu_k)^T}{\sum_{n=1}^N \gamma(z_{nk})}$$

$$\pi_k = \frac{1}{N} \sum_{n=1}^N \gamma(z_{nk})$$

4. Evaluate log likelihood

$$\ln p(\mathbf{X} | \mu, \Sigma, \pi) = \sum_{n=1}^N \ln \left\{ \sum_{k=1}^K \pi_k \mathbf{N}(\mathbf{x}_n | \mu_k, \Sigma_k) \right\}$$

If there is no convergence, return to step 2.

Relationship to K-means

- K-means makes **hard** decisions.
 - Each data point gets assigned to a single cluster.
- GMM/EM makes **soft** decisions.
 - Each data point can yield a posterior $p(z|x)$
- K-means is a special case of EM.

General form of EM

- Given a joint distribution over observed and latent variables: $p(X, Z|\theta)$
- Want to maximize: $p(X|\theta)$

1. Initialize parameters θ^{old}

2. E Step: Evaluate: $p(Z|X, \theta^{old})$

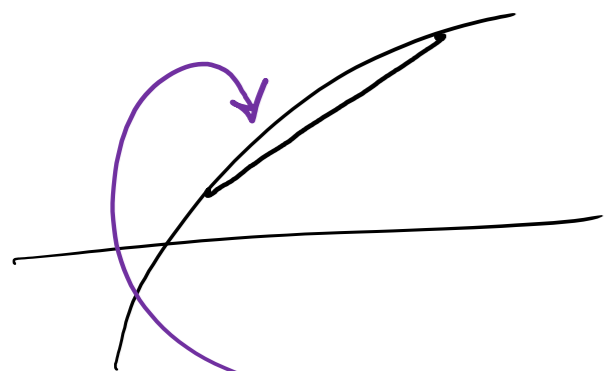
3. M-Step: Re-estimate parameters (based on expectation of complete-data log likelihood)

$$\theta^{new} = \operatorname{argmax}_{\theta} \sum_Z p(Z_k|X, \theta^{old}) \ln p(X, z_k|\theta) = \operatorname{argmax}_{\theta} \operatorname{Exp}[\log(p(x, z_k|\theta))]$$

1. Check for convergence of params or likelihood

$$P(x) = P(x|\theta) = \pi_0 N_0 + \pi_1 N_1 + \pi_2 N_2 = \underbrace{P(z_0) P(x|z_0) + P(z_1) P(x|z_1) + P(z_2) P(x|z_2)}$$

$$\ln P(x) = \ln \sum_k \underbrace{P(x, z_k)}_{P(x, z_0)} = \ln \sum_k \underbrace{q(z_k|x)}_{\text{pdf}} \underbrace{\frac{P(x, z_k)}{q(z_k|x)}}_{\substack{g(x) \\ \text{pdf}}} = \ln \underbrace{\left(\sum P(x) g(x) \right)}_{E[g(x)]}$$



$$\ln(E[g(x)]) \geq E[\log g(x)]$$

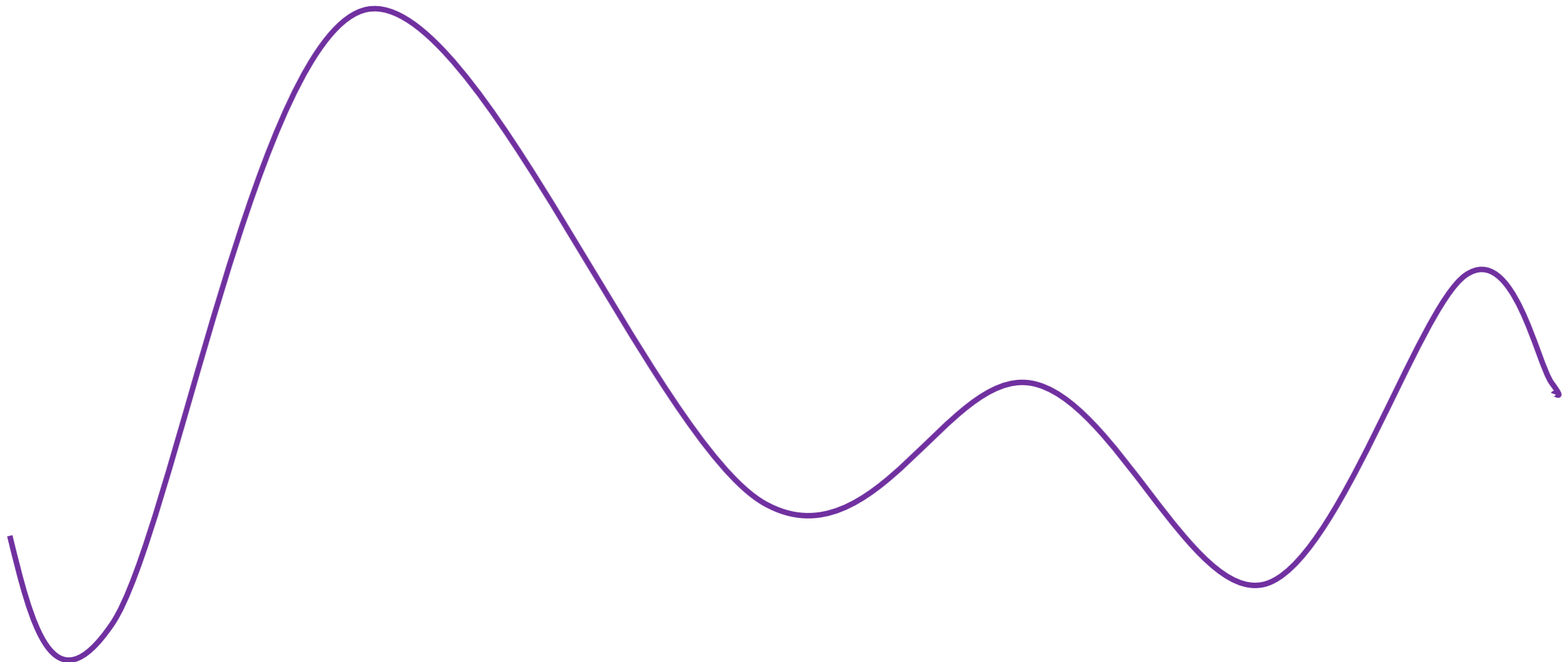
$$\geq \sum P(x) \log g(x)$$

$$\geq \sum q(z_k|x) \log \frac{P(x, z_k)}{q(z_k|x)}$$

$$l(\theta|x) = \log p(x|\theta) = \log \sum_k p(x, z_k|\theta) \geq \sum_k q(z_k|x) \log \frac{p(x, z_k|\theta)}{q(z_k|x)}$$

$$\log \left(\sum_k p(x, z_k|\theta) \right) = \log \left(\sum_z \underbrace{p(x|\theta, z_k)}_{\mathcal{N}_k} * \underbrace{p(z_k|\theta)}_{\pi_k} \right)$$

$$= \log(N(x|\mu_0, \Sigma_0) * \pi_0 + \dots + N(x|\mu_k, \Sigma_k) * \pi_0)$$



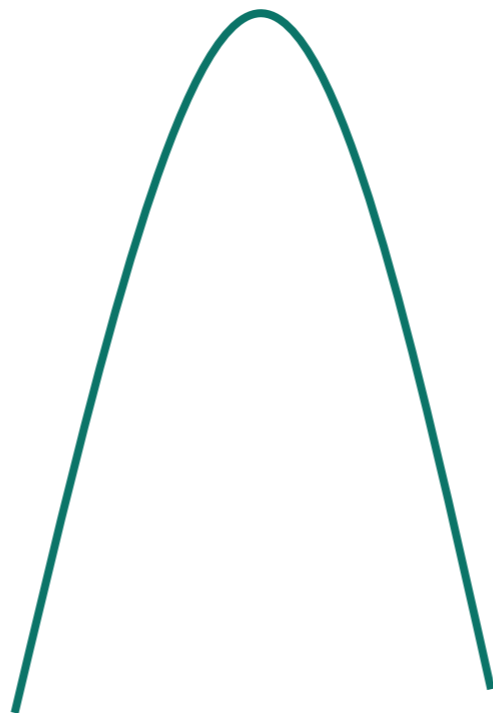
$$l(\theta|x) = \log p(x|\theta) = \log \sum_k p(x, z_k|\theta) \geq \sum_k \underbrace{q(z_k|x)}_{C_k} \log \frac{p(x, z_k|\theta)}{\underbrace{q(z_k|x)}_{C_k}}$$

$q(z_k|x) = C_k \Rightarrow$ It is given to us

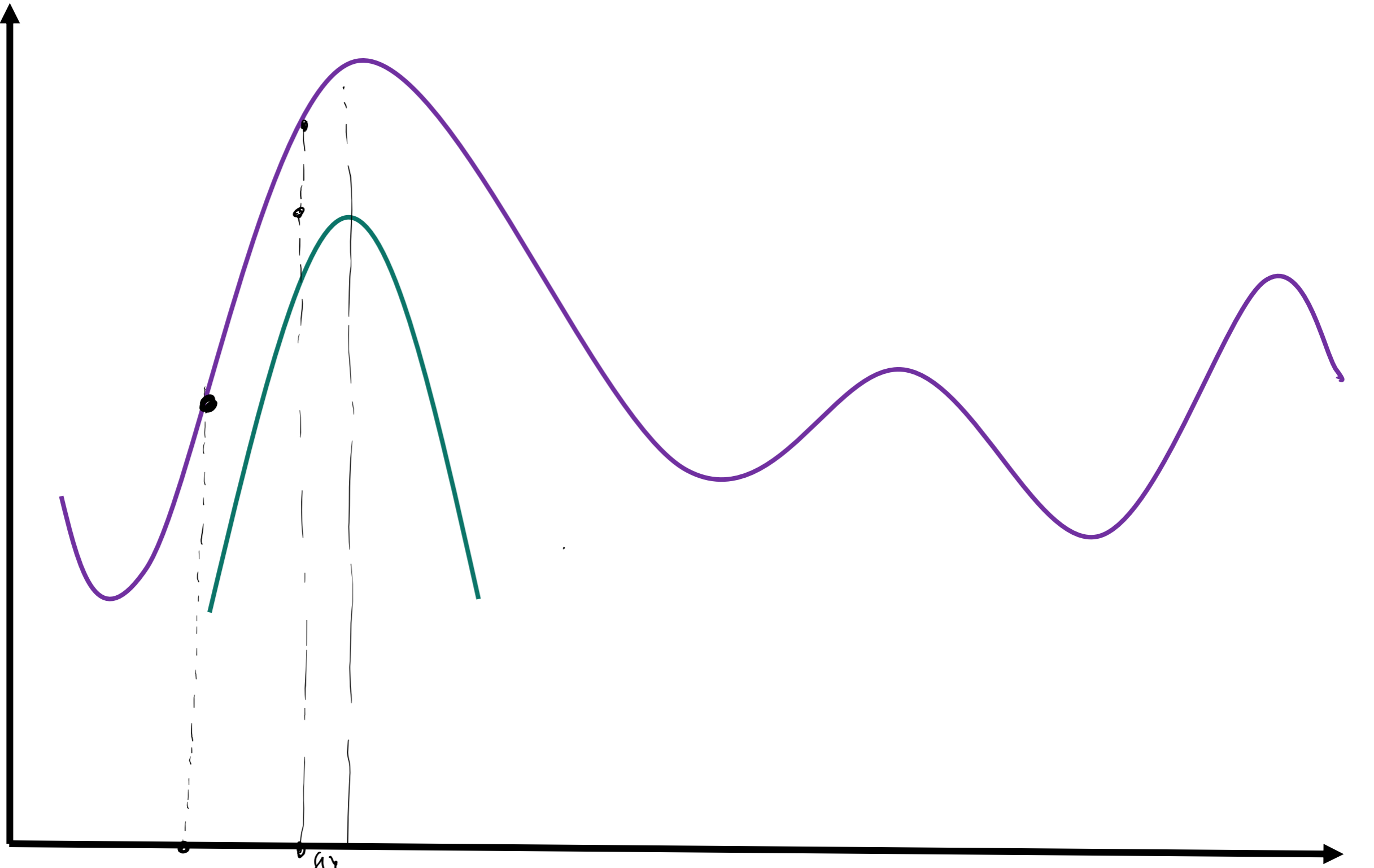
$$\sum_k q(z_k|x) \log \frac{p(x, z_k|\theta)}{q(z_k|x)} =$$

$\rightarrow p(x|z_k) p(z_k)$

$$C_0 \log \left(\frac{1}{C_0} * N(x|\mu_0, \Sigma_0) * \pi_0 \right) + \dots + C_k \log \left(\frac{1}{C_k} * N(x|\mu_k, \Sigma_k) * \pi_k \right)$$



$l(x|\theta)$



$\hat{\theta}$
 $\hat{\theta}_0$
 $\Theta = \{\mu, \sigma, \pi\}$

θ

$$l(\theta|x) = \log p(x|\theta) = \log \sum_z p(x, z|\theta) \stackrel{ub}{=} \sum_z q(z|x) \log \frac{p(x, z|\theta)}{q(z|x)} \quad \frac{a}{b}$$

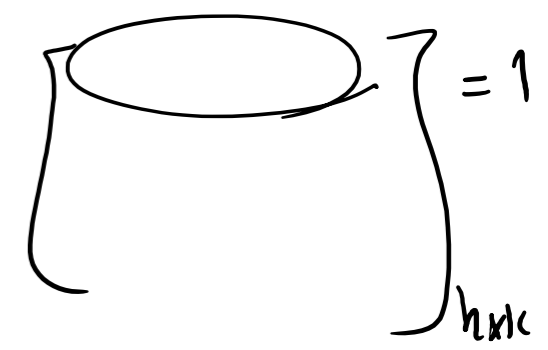
$$q(z|x) = p(z_k | x, \theta^{old})$$

$$\Rightarrow p(x, z|\theta) = p(z_k | x, \theta) * p(x|\theta) = a$$

$$lb = \sum p(z_k | x, \theta^{old}) \left[\log p(z_k | x, \theta) + \log p(x|\theta) - \log p(z_k | x, \theta^{old}) \right]$$

$$lb = \sum_z p(z_k | x, \theta^{old}) * (\log p(x|\theta)) = \log p(x|\theta) \left(\sum p(z_k | x, \theta^{old}) \right)$$

$$lb = \log p(x|\theta)$$



$$l(\theta|x) = \sum_{n=1}^N \sum_{k \in \{0,1\}} p(z_k|x_n, \theta_{old}) \ln [p(x_n, z_k|\theta)]$$

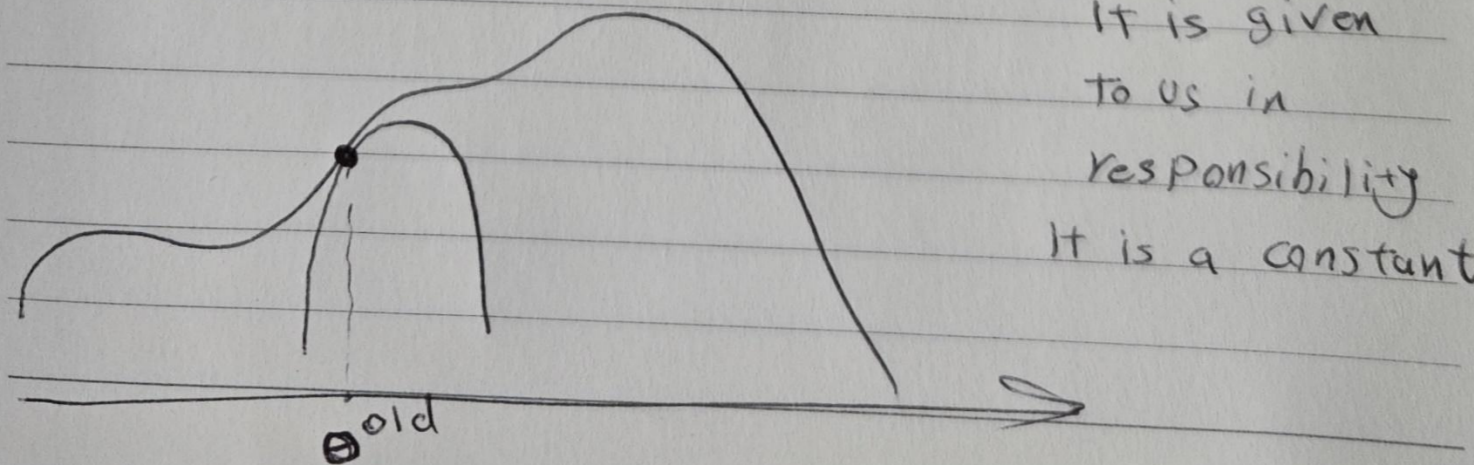
$$\log \sum_k p(x, z_k|\theta) \geq \sum_k q(z_k|x) \log \frac{p(x, z_k|\theta)}{q(z_k|x)}$$

$$\text{If } \Rightarrow q(z_k|x) = p(z_k|x, \theta^{old})$$

$$\text{lower bound} = \sum_k p(z_k|x, \theta^{old}) \log \frac{p(x, z_k|\theta)}{p(z_k|x, \theta^{old})}$$

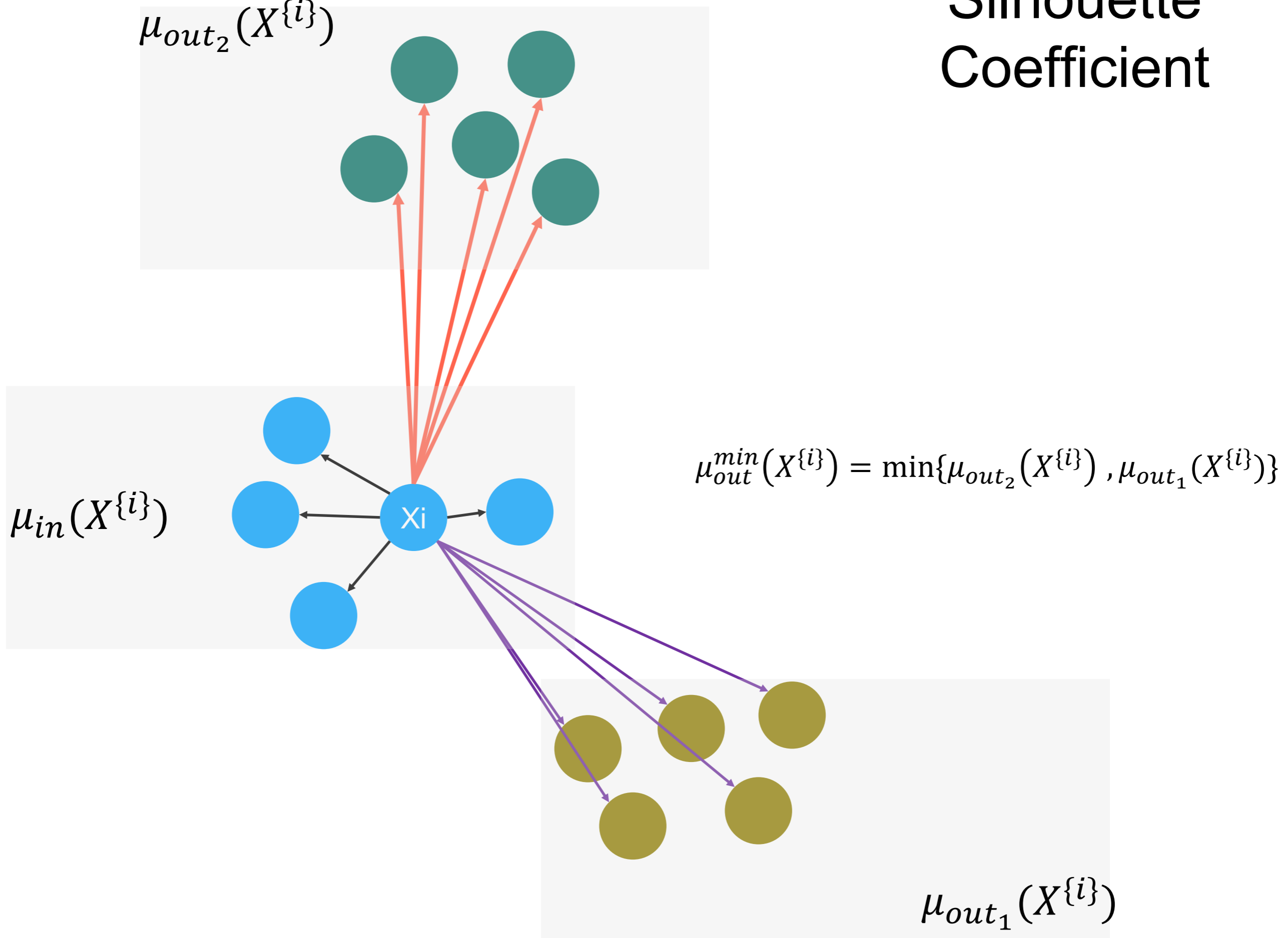
$$\text{lower bound} = \sum_k p(z_k|x, \theta^{old}) [\log p(x, z_k|\theta) - \log p(z_k|x, \theta^{old})]$$

$$\text{lower bound} = \sum_k p(z_k|x, \theta^{old}) \log p(x, z_k|\theta) + \underbrace{\left(\sum_k p(z_k|x, \theta^{old}) \log p(z_k|x, \theta^{old}) \right)}_{\text{constant}}$$



<https://www.dropbox.com/scl/fi/j0nmx654bbluf3zowp78/EM-maximization-and-equality.mp4?rlkey=wlua715r88kdtjdjoru7qc6f6&st=5h7rbts6&dl=0>

Silhouette Coefficient



Silhouette Coefficient

Define the silhouette coefficient of a point \mathbf{x}_i as

$$s_i = \frac{\mu_{out}^{\min}(x^{\{i\}}) - \mu_{in}(x^{\{i\}})}{\max\{\mu_{out}^{\min}(x^{\{i\}}), \mu_{in}(x^{\{i\}})\}}$$

where $\mu_{in}(x^{\{i\}})$ is the mean distance from $x^{\{i\}}$ to points in its own cluster $y^{\{i\}}$

$$\mu_{in}(x^{\{i\}}) = \frac{\sum_{\mathbf{x}_j \in C_{\hat{y}_i}, j \neq i} \delta(x^{\{i\}}, x^{\{j\}})}{n_{y^{\{i\}}} - 1}$$

and $\mu_{out}^{\min}(x^{\{i\}})$ is the mean of the distances from $x^{\{i\}}$ to points in the closest cluster:

$$\mu_{out}^{\min}(x^{\{i\}}) = \min_{j \neq y^{\{i\}}} \left\{ \frac{\sum_{\mathbf{y} \in C_j} \delta(x^{\{i\}}, \mathbf{y})}{n_j} \right\}$$

The Silhouette Coefficient for clustering C: $SC = \frac{1}{n} \sum_{i=1}^n s_i$.

SC close to 1 implies a good clustering (Points are close to their own clusters but far from other clusters)

Take-Home Messages

- The generative process of Gaussian Mixture Model
- Inferring cluster membership based on a learned GMM
- The general idea of Expectation-Maximization
- Expectation-Maximization for GMM
- Silhouette Coefficient

