

# Gaussian Mixture Models

---

Machine Learning 1 — Lecture 13

19<sup>th</sup> June 2023

Robert Peharz

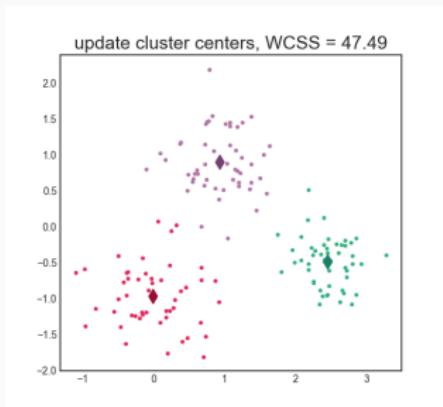
Institute of Theoretical Computer Science  
Graz University of Technology

## K-Means Problem

$$\begin{aligned} \min_{\mu^1, \dots, \mu^K, Z} \quad & \sum_{i=1}^N \sum_{k=1}^K z_{i,k} \|x^{(i)} - \mu^{(k)}\|^2 \\ \text{s.t.} \quad & z_{i,k} \in \{0, 1\} \\ & \sum_k z_{i,k} = 1 \quad i = 1, \dots, N \end{aligned}$$

## K-Means Algorithm

- assign data points to cluster centers
- set cluster centers as means of assigned data points
- repeat



Today, we discuss **Gaussian mixture models (GMMs)** which can be seen as a probabilistic version of k-means. The k-means algorithm corresponds to the **expectation-maximization (EM)** algorithm for learning GMMs.

Main purpose of GMMs is **density estimation**, i.e., approximating the true data generating distribution  $p_{true}$ .

Assume  $N$  samples  $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$  drawn i.i.d. (**independently and identically distributed**) from a true distribution with density  $p_{true}(\mathbf{x})$ . In **density estimation**, we want to learn/estimate a **model distribution**  $p_{model}(\mathbf{x})$  such that

$$p_{model} \approx p_{true}$$

In the following, we will denote  $p_{model}$  simply as  $p$ .

# Applications of Density Estimation

Say we have found a model with  $p_{model} \approx p_{true}$ . Generally, knowing  $p_{true}$  allows us—in principle—to answer any question we might have about the data.

The model can be used for

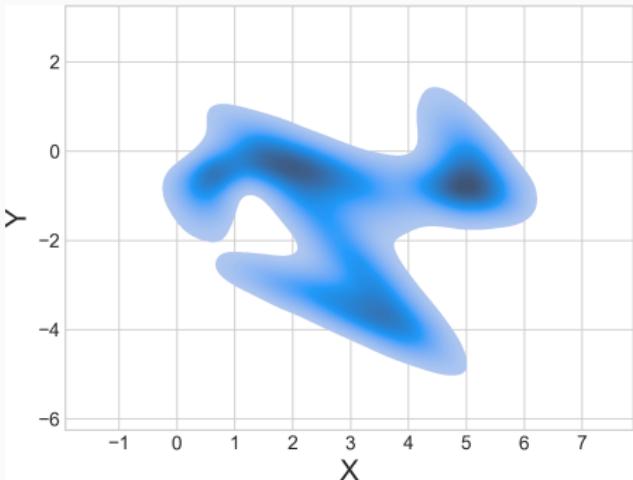
- **outlier detection** (monitor whether  $p_{model}(\mathbf{x})$  is low)
- **generative modeling**: sample “new data” ( $\mathbf{x} \sim p_{model}$ )
- prediction via **conditional distributions**

$$p_{model}(\mathbf{x}_q | \mathbf{x}_e) = \frac{p_{model}(\mathbf{x}_q, \mathbf{x}_e)}{p_{model}(\mathbf{x}_e)}$$

where  $\mathbf{x}_e$  is evidence for some variables  $\mathbf{X}_e$  and  $\mathbf{x}_q$  are query variables.

### Regression

Say you want to predict  $Y$  from  $X$ , and you would know the true distribution  $p_{true}(y|x)$  producing i.i.d. samples of  $X, Y$ .

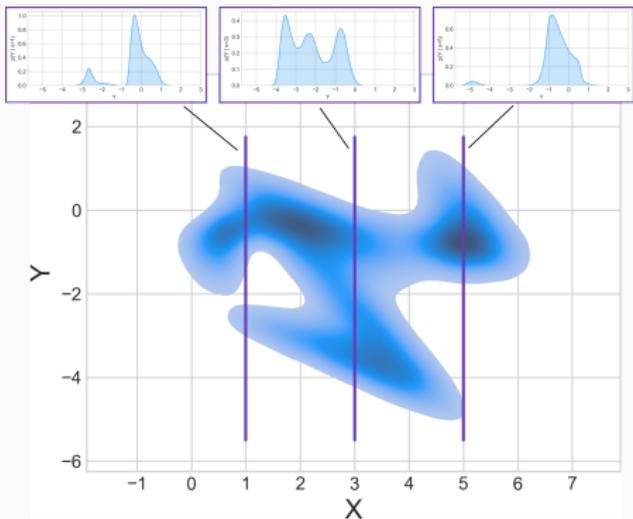


In practice, we almost never know  $p_{true}$ , so this is a very strong assumption.

You might compute the **conditional distribution**

$$p_{true}(y | x) = \frac{p_{true}(y, x)}{\underbrace{p_{true}(x)}_{\int p_{true}(x, y) dy}}$$

for all values of  $x$ .

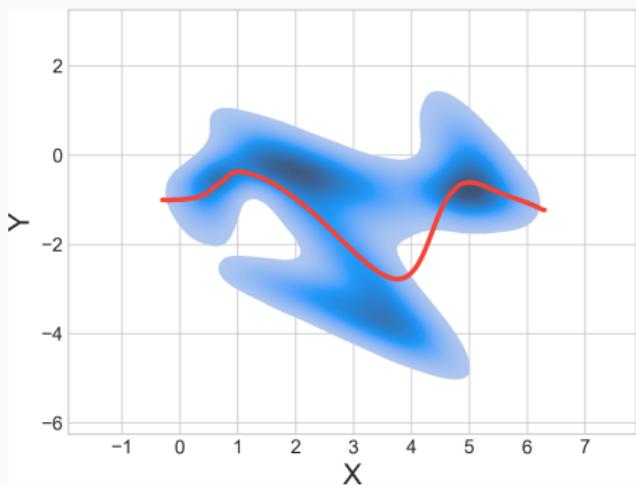


Then you can construct the  
**true regression function**

$$f(x) := \mathbb{E}_{p_{true}(y|x)}[Y|x]$$

which is just the **expectation**  
of  $Y$  given  $x$

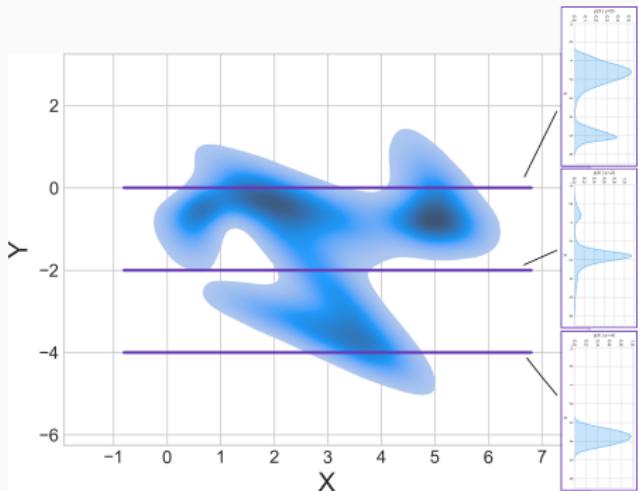
No other function has lower  
**true squared loss.**



Of course, when we now replace  $p_{true}$  with  $p_{model}$ , this is not true. In practice, directly learning  $f(x)$  will deliver better results.

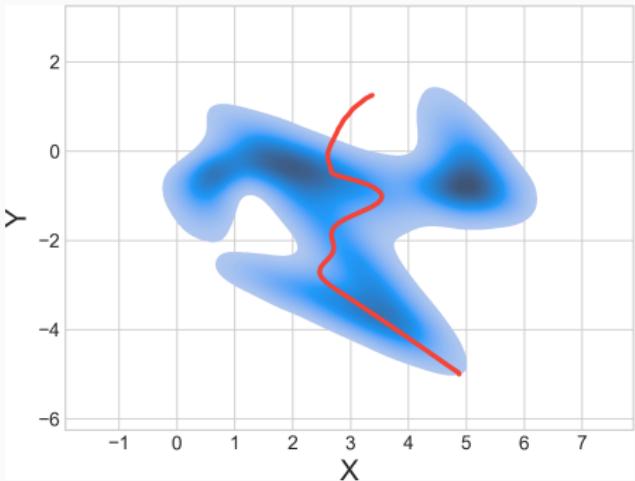
Of course, that  $Y$  is output and  $X$  is input was arbitrary, so the whole “trick” works also in the other direction:

$$p_{true}(x | y) = \frac{p_{true}(y, x)}{\underbrace{p_{true}(y)}_{\int p_{true}(x, y) dx}}$$



Of course, that  $Y$  is output and  $X$  is input was arbitrary, so the whole “trick” works also in the other direction:

$$f(y) := \mathbb{E}_{p_{true}(x|y)}[X | y]$$



## Maximum Likelihood Principle

---

## Setting

Say we have a dataset  $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$  drawn i.i.d. from some unknown distribution  $p_{true}(\mathbf{x})$ .

Further, assume we have a **parametric model**  $p(\mathbf{x}; \boldsymbol{\theta})$ , i.e., a **family of distributions** or a **model class** described by a parameter vector  $\boldsymbol{\theta}$ , e.g., Gaussian, Poisson, Bernoulli, etc.

How do we fit  $p(\mathbf{x}; \boldsymbol{\theta})$  to  $p_{true}(\mathbf{x})$ ?

Since the samples are drawn i.i.d., the probability density of the whole dataset  $\mathcal{D}$  is

$$p(\mathcal{D}; \theta) = \prod_{i=1}^N p(x^{(i)}; \theta)$$

Note: Generally, multiple random variables  $X_1, X_2, \dots, X_D$  are **statistically independent** if their distribution factorizes:

$$p(X_1, X_2, \dots, X_D) = p(X_1) \times p(X_2) \times \cdots \times p(X_D)$$

The probability density  $p(\mathcal{D}; \theta)$  factorizes due to the i.i.d. assumption.

The quantity

$$p(\mathcal{D}; \theta) = \prod_{i=1}^N p(x^{(i)}; \theta)$$

depends on the dataset  $\mathcal{D}$  and the parameters  $\theta$ .

- When interpreted as function of  $\mathcal{D}$ : probability density of  $\mathcal{D}$
- When interpreted as function of  $\theta$ : **likelihood** of  $\theta$

**Maximum Likelihood Estimator**: Parameters  $\theta^*$  which maximize the likelihood

$$\theta^* = \arg \max_{\theta} p(\mathcal{D}; \theta) = \arg \max_{\theta} \prod_{i=1}^N p(x^{(i)}; \theta),$$

i.e., the parameters  $\theta^*$  under which the data is most likely.

Instead of likelihood, we can maximize the **log-likelihood**:

$$\mathcal{L}(\theta) := \log \prod_{i=1}^N p(\mathbf{x}^{(i)}; \theta) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}; \theta)$$

$$\theta^* = \arg \max_{\theta} \mathcal{L}(\theta)$$

The log-likelihood is measured in **bits** (logarithm of base 2) or **nats** (logarithm of base  $e$ ). Since the log is a **strictly increasing** function, the log-likelihood has exactly the same maxima as the likelihood, i.e. maximizing log-likelihood is equivalent to maximizing likelihood.

## Why the Log?

Optimizing the log-likelihood (sum of terms) is usually easier than optimizing the likelihood (product of factors). Moreover, when using numerical optimizers, the log-likelihood is **numerically stable**, while the likelihood, a product of many small factors, becomes quickly numerical zero.

Why is maximizing the (log-)likelihood meaningful?

In density estimation, we want to find parameters  $\theta$  such that

$$p(\mathbf{x}; \theta) \approx p_{\text{true}}(\mathbf{x})$$

To quantify the approximation quality (" $\approx$ ") we require a distance or **divergence measure** for comparing  $p(\mathbf{x}; \theta)$  and  $p_{\text{true}}(\mathbf{x})$ . One of the most widely divergences is the **Kullback-Leibler divergence**:

$$KL(p_{\text{true}} || p) = \mathbb{E}_{p_{\text{true}}} \left[ \log \frac{p_{\text{true}}(\mathbf{x})}{p(\mathbf{x}; \theta)} \right] = \int p_{\text{true}}(\mathbf{x}) \log \frac{p_{\text{true}}(\mathbf{x})}{p(\mathbf{x}; \theta)} d\mathbf{x}$$

It holds that

- $KL(p_{\text{true}} || p) \geq 0$  and
- $KL(p_{\text{true}} || p) = 0$  iff  $p_{\text{true}} \equiv p$  almost everywhere.

# Kullback-Leibler Divergence and Log-likelihood

Note that the KL divergence can be written as

$$\begin{aligned} KL(p_{true} || p) &= \int p_{true}(\mathbf{x}) \log \frac{p_{true}(\mathbf{x})}{p(\mathbf{x}; \theta)} d\mathbf{x} \\ &= \int p_{true}(\mathbf{x}) [\log p_{true}(\mathbf{x}) - \log p(\mathbf{x}; \theta)] d\mathbf{x} \\ &= \mathbb{E}_{p_{true}} [\log p_{true}(\mathbf{x}) - \log p(\mathbf{x}; \theta)] \\ &= \underbrace{\mathbb{E}_{p_{true}} [\log p_{true}(\mathbf{x})]}_{\text{neg. entropy } -H[p_{true}], \text{const.}} - \mathbb{E}_{p_{true}} [\log p(\mathbf{x}; \theta)] \end{aligned}$$

The second term,  $-\mathbb{E}_{p_{true}} [\log p(\mathbf{x}; \theta)]$ , can be estimated with Monte Carlo, which is proportional to the log-likelihood:

$$-\mathbb{E}_{p_{true}} [\log p(\mathbf{x}; \theta)] \approx -\frac{1}{N} \overbrace{\sum_{i=1}^N \log p(\mathbf{x}^{(i)}; \theta)}^{\text{log-likelihood } \mathcal{L}(\theta)}$$

## Kullback-Leibler Divergence and Log-likelihood cont'd

Thus, maximum likelihood (ML) is well-justified for density approximation:

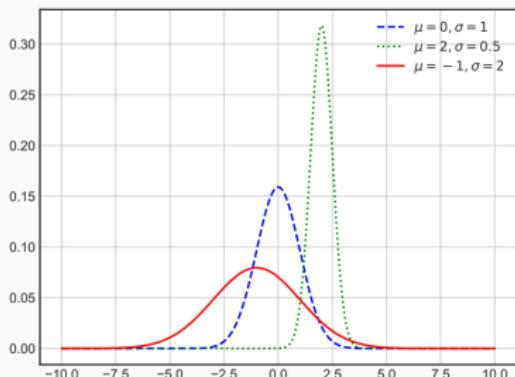
- maximizing log-likelihood is equivalent to minimizing a Monte Carlo estimate of the KL divergence  $KL(p_{true} || p)$
- if  $p_{true}$  is contained in the model class  $\{p(\mathbf{x}; \boldsymbol{\theta})\}_{\boldsymbol{\theta}}$ , then the maximum likelihood distribution  $p(\mathbf{x}; \boldsymbol{\theta}^*)$  will converge to  $p_{true}$  when  $N \rightarrow \infty$
- if  $p_{true}$  is not contained in the model class  $\{p(\mathbf{x}; \boldsymbol{\theta})\}_{\boldsymbol{\theta}}$ ,  $p(\mathbf{x}; \boldsymbol{\theta}^*)$  will converge to the best approximation of  $p_{true}$  in KL sense

## Gaussian Distribution, Normal Distribution

Let  $p$  be a probability density defined as

$$p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2}$$

with parameters **mean**  $\mu$  and **standard deviation**  $\sigma > 0$ . A distribution with density  $p$  is called **Gaussian distribution** or **normal distribution**.



Carl Friedrich Gauss 1777–1855

Image: wikipedia

## Maximum Likelihood for Gaussians

- Assume i.i.d. dataset  $\mathcal{D} = \{x^{(1)}, x^{(2)}, \dots, x^{(N)}\}$
- The Gaussian pdf is  $p(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$
- Thus, the log-likelihood is

$$\mathcal{L}(\mu, \sigma) = \log p(\mathcal{D}; \mu, \sigma) = \sum_{i=1}^N -\log(\sqrt{2\pi}\sigma) - \frac{1}{2} \left( \frac{x^{(i)} - \mu}{\sigma} \right)^2$$

- We want to maximize  $\mathcal{L}(\mu, \sigma)$  w.r.t.  $\mu$  and  $\sigma > 0$
- It turns out, that we can solve the problem step-wise: first find a maximizer for  $\mu$ , then for  $\sigma$

## Maximum Likelihood for Gaussians: $\mu^*$

The log-likelihood is a negative quadratic function in  $\mu$ :

$$\mathcal{L}(\mu, \sigma) = \sum_{i=1}^N -\log(\sqrt{2\pi} \sigma) - \frac{1}{2} \left( \frac{x^{(i)} - \mu}{\sigma} \right)^2$$

Derive after  $\mu$  and set to 0:

$$\frac{\partial \mathcal{L}(\mu, \sigma)}{\partial \mu} = \sum_{i=1}^N -\frac{1}{2} 2 \left( \frac{x^{(i)} - \mu}{\sigma} \right) \left( -\frac{1}{\sigma} \right) \stackrel{!}{=} 0$$

Multiplying by  $\sigma^2$  and bringing  $\mu$  on the other side delivers the **average** as maximum likelihood (ML) solution:

$$\sum_{i=1}^N x^{(i)} = \sum_{i=1}^N \mu^* = N\mu^* \quad \Rightarrow \quad \mu^* = \frac{1}{N} \sum_{i=1}^N x^{(i)}$$

Note that this is the ML solution for  $\mu$  for **any**  $\sigma$ . Thus we can already assume  $\mu^*$  when computing  $\sigma^*$ .

## Maximum Likelihood for Gaussians: $\sigma^*$

$$\mathcal{L}(\mu^*, \sigma) = \sum_{i=1}^N -\log(\sqrt{2\pi} \sigma) - \frac{1}{2} \left( \frac{x^{(i)} - \mu^*}{\sigma} \right)^2$$

Deriving after  $\sigma$ :

$$\begin{aligned}\frac{\partial \mathcal{L}(\mu, \sigma)}{\partial \sigma} &= \sum_{i=1}^N -\frac{\sqrt{2\pi}}{\sqrt{2\pi} \sigma} - \frac{1}{2} 2 \left( \frac{x^{(i)} - \mu^*}{\sigma} \right) (-1) \left( \frac{x^{(i)} - \mu^*}{\sigma^2} \right) \\ &= \sum_{i=1}^N -\frac{1}{\sigma} + \left( \frac{(x^{(i)} - \mu^*)^2}{\sigma^3} \right) \stackrel{!}{=} 0\end{aligned}$$

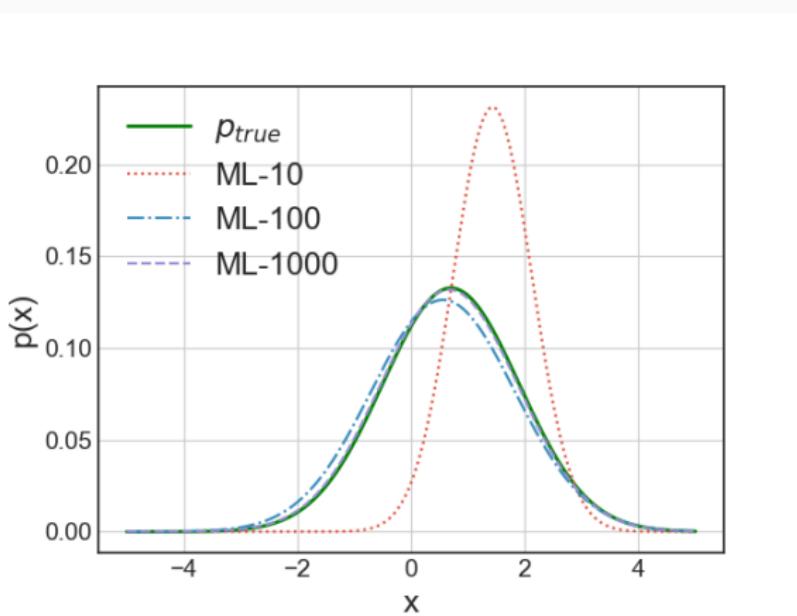
Multiplying by  $\sigma^3$  and resolving after  $\sigma$  delivers the **empirical standard deviation** as ML solution:

$$\sigma^* = \sqrt{\frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu^*)^2}$$

# Maximum Likelihood for Gaussians

Example

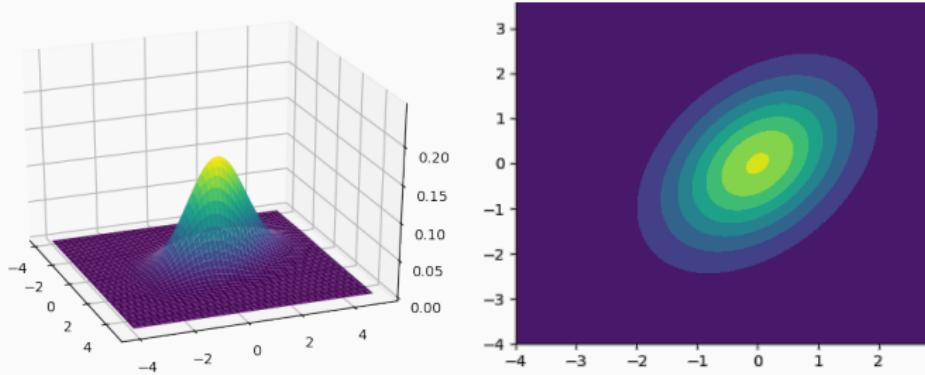
Drawing three datasets with  $N = 10$ ,  $N = 100$  and  $N = 1000$  from a ground-truth Gaussian with  $\mu_{true} = 0.7$  and  $\sigma_{true} = 1.2$ . With increased sample size, the ML fit approaches the ground-truth:



Let  $p$  be defined as

$$p(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^D |\boldsymbol{\Sigma}|}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1} (\mathbf{x}-\boldsymbol{\mu})}$$

with  $D$ -dimensional **mean vector**  $\boldsymbol{\mu}$  and  $D \times D$  positive definite **covariance matrix**  $\boldsymbol{\Sigma}$ . A distribution with density  $p$  is called **multivariate Gaussian distribution**. It generalizes Gaussians to higher dimensions.



# Maximum Likelihood for Multivariate Gaussians

- The ML solution for 1d Gaussians...

$$\mu^* = \frac{1}{N} \sum_{i=1}^N x^{(i)} \quad \sigma^2* = \frac{1}{N} \sum_{i=1}^N (x^{(i)} - \mu^*)^2$$

- ... carries over to multivariate Gaussians:

$$\boldsymbol{\mu}^* = \frac{1}{N} \sum_{i=1}^N \mathbf{x}^{(i)} \quad \boldsymbol{\Sigma}^* = \frac{1}{N} \sum_{i=1}^N (\mathbf{x}^{(i)} - \boldsymbol{\mu}^*)(\mathbf{x}^{(i)} - \boldsymbol{\mu}^*)^T$$

# Gaussian Mixture Models

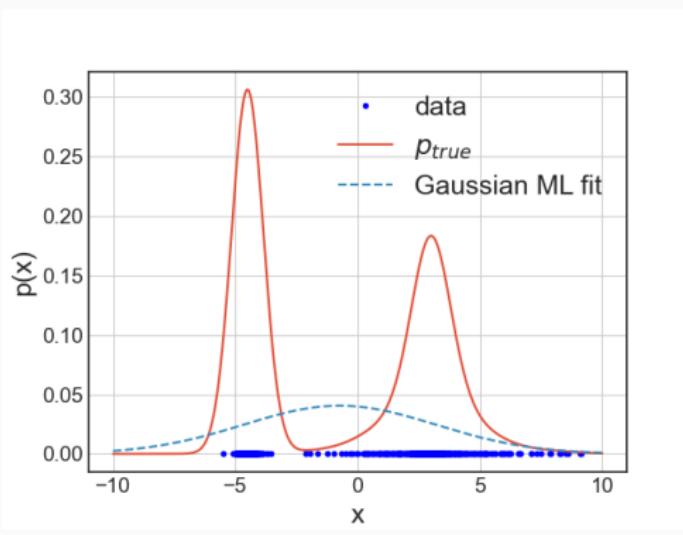
---

Are Gaussian models always a good choice?

## A Gaussian (Mis)fit

Example

Assume a true distribution  $p_{true}$  as shown below. We collect  $N = 1000$  samples and fit a Gaussian with maximum likelihood:



The fit is awful, since  $p_{true}$  is non-Gaussian. In particular, it is a **bimodal** distribution (**mode**: local maximum of density).

Let  $K \geq 1$  and

- $\mu_1, \mu_2, \dots, \mu_K$  be  $K$  mean parameters
- $\Sigma_1, \Sigma_2, \dots, \Sigma_K$  be  $K$  covariance matrices
- $w_1, w_2, \dots, w_K$  be  $K$  **mixture weights** with  $w_k \geq 0$  and  $\sum_k w_k = 1$

The **Gaussian mixture model (GMM)** parametrized by

$\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$  is given as

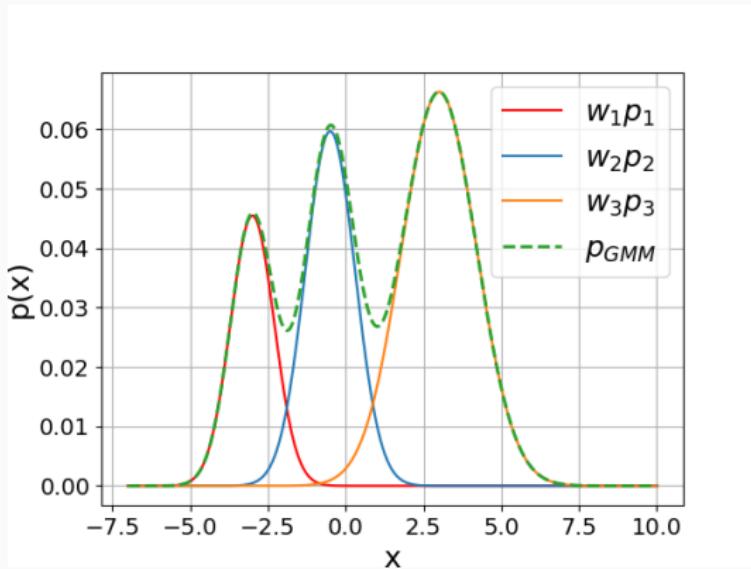
$$p_{GMM}(\mathbf{x}, \theta) = \sum_{k=1}^K w_k p(\mathbf{x}; \mu_k, \Sigma_k)$$

where  $p(\mathbf{x}; \mu_k, \Sigma_k)$  is the Gaussian density with parameters  $\mu_k, \Sigma_k$ .

That is, a GMM is a **convex combination** of  $K$  **Gaussian components**.

# Gaussian Mixture Model, 1 Dimension

Example



3 components:

$$w_1 = 0.2, w_2 = 0.3, w_3 = 0.5,$$

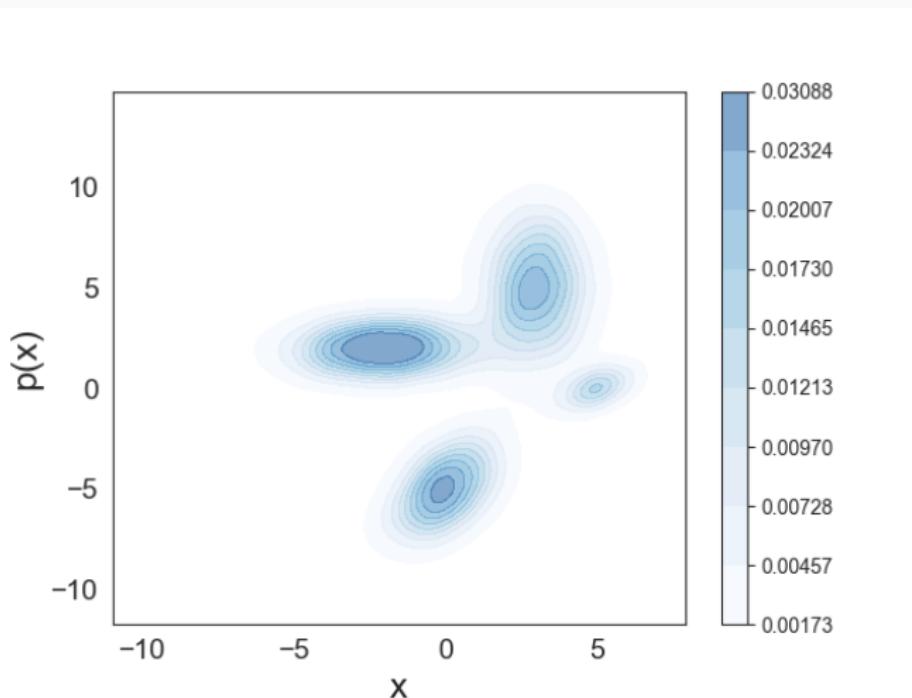
$$\mu_1 = -3, \mu_2 = -0.5, \mu_3 = 3,$$

$$\sigma_1 = 0.7, \sigma_2 = 0.8, \sigma_3 = 1.2$$

# Gaussian Mixture Model, 2 Dimensions

Example

A GMM with 4 components in 2 dimensions:



# Properties of GMMs

Any GMM specifies a proper probability density, since

- $p_{GMM}(\mathbf{x})$  is non-negative:
  - $w_k \geq 0$  and
  - $p(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \geq 0$
  - Hence,  $p_{GMM}(\mathbf{x}) = \sum_{k=1}^K w_k p(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \geq 0$
- $p_{GMM}$  is normalized:

$$\begin{aligned}\int p_{GMM}(\mathbf{x}) d\mathbf{x} &= \int \sum_{k=1}^K w_k p(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) d\mathbf{x} \\ &= \sum_{k=1}^K w_k \int p(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) d\mathbf{x} \\ &= \sum_{k=1}^K w_k 1 = 1\end{aligned}$$

## Properties of GMMs cont'd

GMMs are **universal approximators** of densities, that is, any density can be approximated arbitrarily well by a GMM. **This is true even when restricted Gaussian with only diagonal covariances.**

However, we don't know a big  $K$  needs to be and how to set the parameters.

# Expectation-Maximization

---

## GMM as Latent Variable Model

Note that the mixture weights  $w_k$  are non-negative and sum to one. Thus, we can interpret them as the distribution of a **categorical latent variable**  $Z$  taking values in  $\{1, \dots, K\}$ :

$$\begin{aligned} p_{GMM}(\mathbf{x}) &= \sum_{k=1}^K w_k p(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \\ &= \sum_{k=1}^K p(z=k) p(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \end{aligned}$$

Thus, a GMM actually defines a joint distribution over  $\mathbf{X}$  and  $Z$ :

$$p_{GMM}(\mathbf{x}, z = k) = \underbrace{p(z = k)}_{w_k} \underbrace{p(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}_{p(\mathbf{x} | z = k)}$$

and  $p_{GMM}(\mathbf{x})$  is actually just the **marginal distribution**

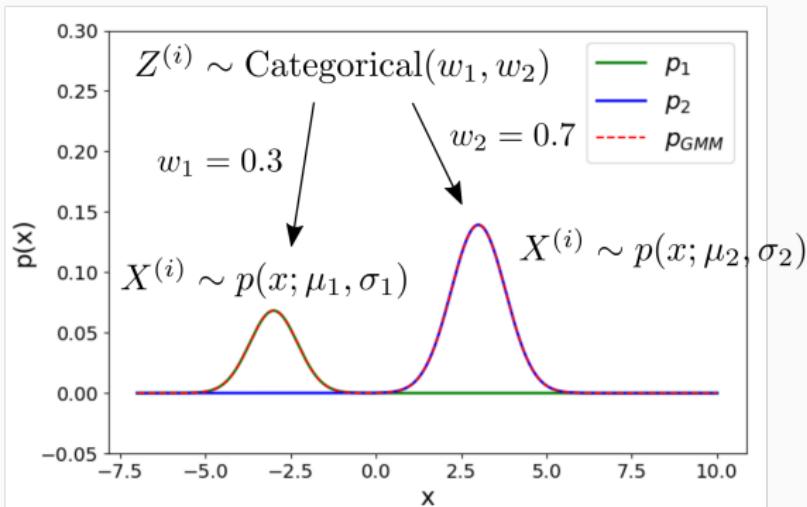
$$p_{GMM}(\mathbf{x}) = \sum_k p_{GMM}(\mathbf{x}, z = k)$$

Thus, a GMM can be seen as a **hierarchical generative model**:

For each  $i = 1 \dots N$ , generate  $Z^{(i)}, \mathbf{X}^{(i)}$  as follows:

$$Z^{(i)} \sim \text{Categorical}(w_1, \dots, w_K)$$

$$\mathbf{X}^{(i)} \sim p(\mathbf{x} | z^{(i)}) = p(\mathbf{x}; \boldsymbol{\mu}_{z^{(i)}}, \boldsymbol{\Sigma}_{z^{(i)}})$$



## Connection to K-Means

Recall, in k-means, we had  $z_{i,k} = 1$  iff  $\mathbf{x}^{(i)}$  belongs to  $k^{\text{th}}$  cluster.

The latent variables  $Z^{(i)}$  in GMMs have exactly the same meaning: they encode which Gaussian component has generated  $\mathbf{x}^{(i)}$ . Thus, GMMs can be seen as probabilistic clustering, or “soft clustering.”

**Hen and Egg Problem:** In GMMs, neither the parameters

$\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$  nor the latent variables

$\mathbf{Z} = \{Z^{(1)}, Z^{(2)}, \dots, Z^{(N)}\}$  are known. The problem would be easy, if either were known.

# The Hen and Egg Problem

If  $Z$  was known, ML estimation for  $\theta$  would be easy:

- $Z$  specify clusters:  $\mathcal{D}_k = \{\mathbf{x}^{(i)} \mid z^{(i)} = k\}$
- mean vectors  $\mu_k^*$  are estimated as the means of the  $k^{\text{th}}$  cluster

$$\mu_k^* = \frac{1}{|\mathcal{D}_k|} \sum_{\mathbf{x} \in \mathcal{D}_k} \mathbf{x}$$

- covariance matrices  $\Sigma_k^*$  are estimated as empirical covariances of  $\mathcal{D}_k$

$$\Sigma_k^* = \frac{1}{|\mathcal{D}_k|} \sum_{\mathbf{x} \in \mathcal{D}_k} (\mathbf{x} - \mu_k^*)(\mathbf{x} - \mu_k^*)^T$$

- the ML mixture weights are the empirical frequencies of  $Z^{(i)}$ :

$$w_k^* = \frac{1}{N} \sum_{i=1}^N \mathbb{1}(z^{(i)} = k)$$

## The Hen and Egg Problem cont'd

If  $\theta$  was known, we could infer  $Z$  via Bayes rule: \*

$$\begin{aligned} p(Z^{(i)} = k | \mathbf{x}^{(i)}) &= \frac{p(\mathbf{x}^{(i)} | Z^{(i)} = k) p(Z^{(i)} = k)}{p(\mathbf{x}^{(i)})} \\ &= \frac{p(\mathbf{x}^{(i)}; \boldsymbol{\mu}_k, \Sigma_k) w_k}{\sum_{k'} p(\mathbf{x}^{(i)}; \boldsymbol{\mu}_{k'}, \Sigma_{k'}) w_{k'}} =: \gamma_{i,k} \end{aligned}$$

Note that  $\gamma_{i,k} \geq 0$  and  $\sum_k \gamma_{i,k} = 1$ .

$\gamma_{i,k}$ : posterior over  $Z^{(i)}$  given  $\mathbf{x}^{(i)}$ , also called “cluster responsibilities.”

$\gamma_{i,k}$  can be interpreted as a fractional weight, i.e., how much the  $i^{\text{th}}$  sample belongs to the  $k^{\text{th}}$  component. Furthermore,  $\gamma_{i,k}$  can be seen as a “soft” version of the “hard” cluster assignments  $z_{i,k} \in \{0, 1\}$  in k-means.

\* Bayes rule in general:  $p(Y | X) = \frac{p(X | Y) p(Y)}{p(X)}$

Idea of the **Expectation-Maximization (EM)** algorithm:

**Repeat:**

- **E-step:** Compute cluster responsibilities  $\gamma_{i,k}$  using the current parameters  $\theta$
- **M-step:** Update  $\theta$  using maximum likelihood with samples weighted by  $\gamma_{i,k}$

# Expectation Maximization Algorithm for GMMs

**input:** Data  $\mathcal{D} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$

**output:** Learned GMM parameters  $\boldsymbol{\theta} = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$

Initialize  $\boldsymbol{\theta} = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$  randomly

**for** counter = 1 ... max number of iterations **do**

## E-Step

**for**  $i = 1 \dots N, k = 1 \dots K$  **do**

$$\text{let } \gamma_{i,k} \leftarrow \frac{p(\mathbf{x}^{(i)}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) w_k}{\sum_{k'} p(\mathbf{x}^{(i)}; \boldsymbol{\mu}_{k'}, \boldsymbol{\Sigma}_{k'}) w_{k'}} \quad \triangleright \text{cluster responsibilities}$$

**end**

## M-Step

**for**  $k = 1 \dots K$  **do**

$$\boldsymbol{\mu}_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} \mathbf{x}^{(i)}}{\sum_{i=1}^N \gamma_{i,k}} \quad \triangleright \text{weighted mean}$$

$$\boldsymbol{\Sigma}_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k} (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^T}{\sum_{i=1}^N \gamma_{i,k}} \quad \triangleright \text{weighted covariance}$$

$$w_k \leftarrow \frac{\sum_{i=1}^N \gamma_{i,k}}{N} \quad \triangleright \text{weighted empirical frequencies}$$

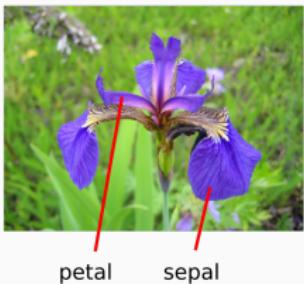
**end**

**end**

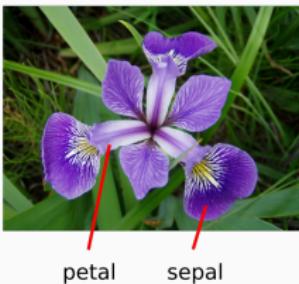
# Learning GMM with EM on Iris Data

Example

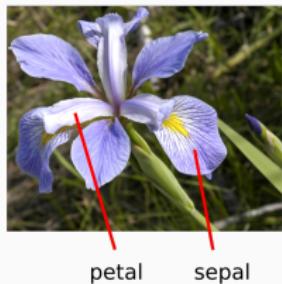
Iris Setosa



Iris Versicolor



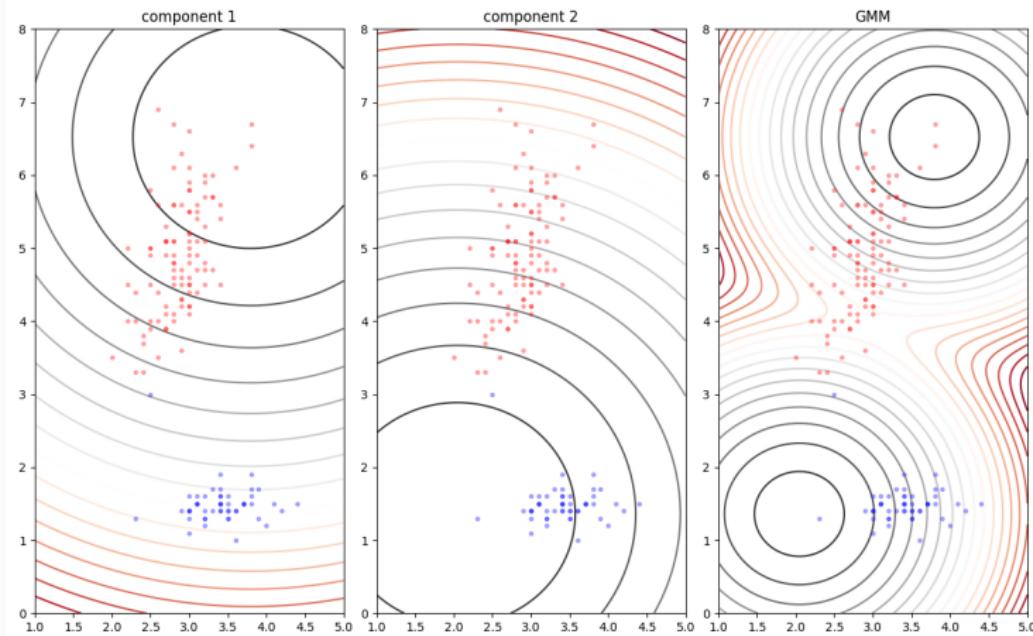
Iris Virginica



- $D = 2$  dimension
- $K = 2$  components
- Initialization  $w_1 = w_2 = 0.5$
- $\mu_1, \mu_2$ : uniformly at random from data range
- $\Sigma_1 = \Sigma_2 = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$

# Learning GMM with EM on Iris Data

Example

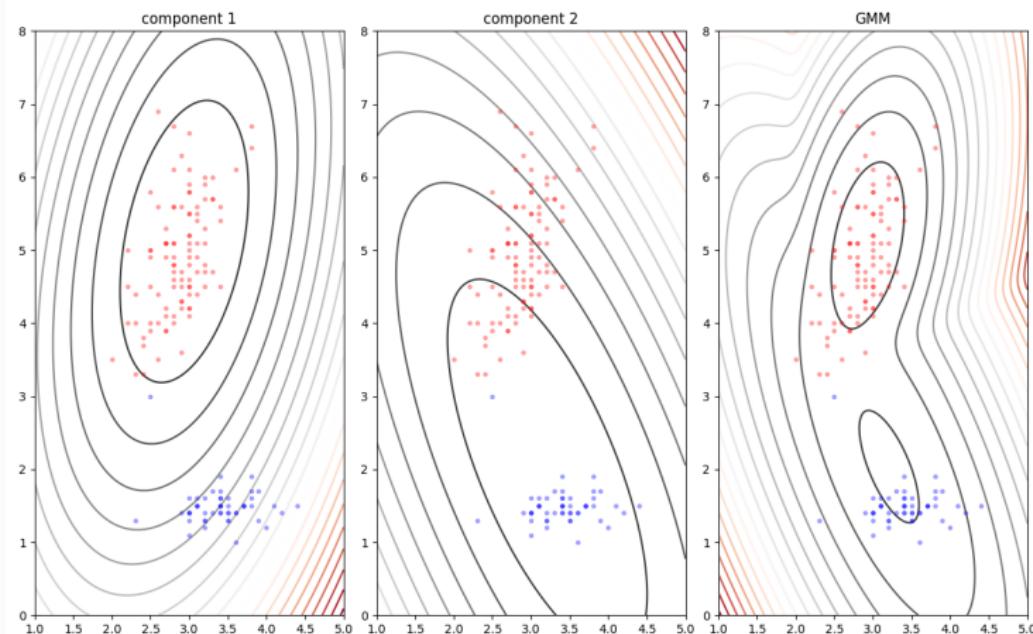


Init: log-likelihood =  $-604.86$  nats,  $w_1 = 0.5$ ,  $w_2 = 0.5$

The blue and red labels is the clustering by K-Means; this **not** used by EM

# Learning GMM with EM on Iris Data

Example

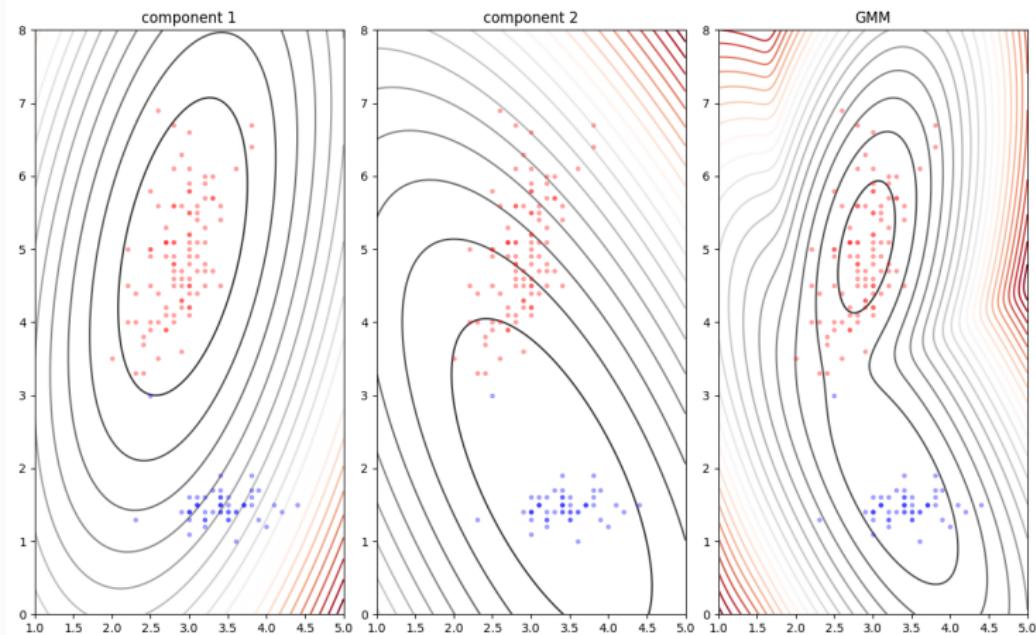


Iteration 1:  $\log\text{-likelihood} = -309.27 \text{ nats}$ ,  $w_1 = 0.56$ ,  $w_2 = 0.44$

The blue and red labels is the clustering by K-Means; this **not** used by EM

# Learning GMM with EM on Iris Data

Example

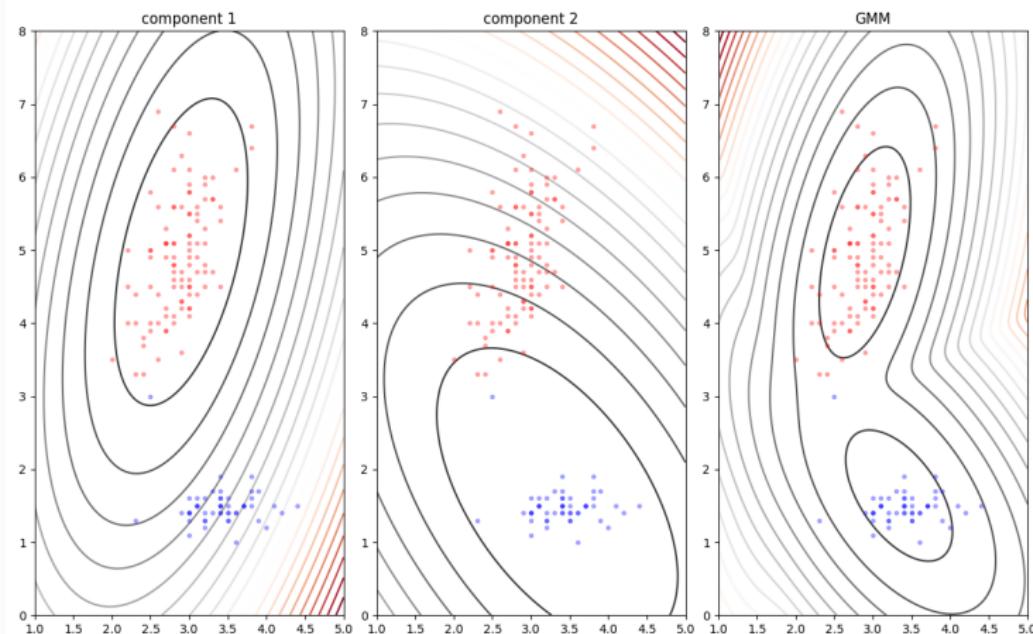


Iteration 5:  $\log\text{-likelihood} = -299.75$  nats,  $w_1 = 0.6$ ,  $w_2 = 0.4$

The blue and red labels is the clustering by K-Means; this **not** used by EM

# Learning GMM with EM on Iris Data

Example

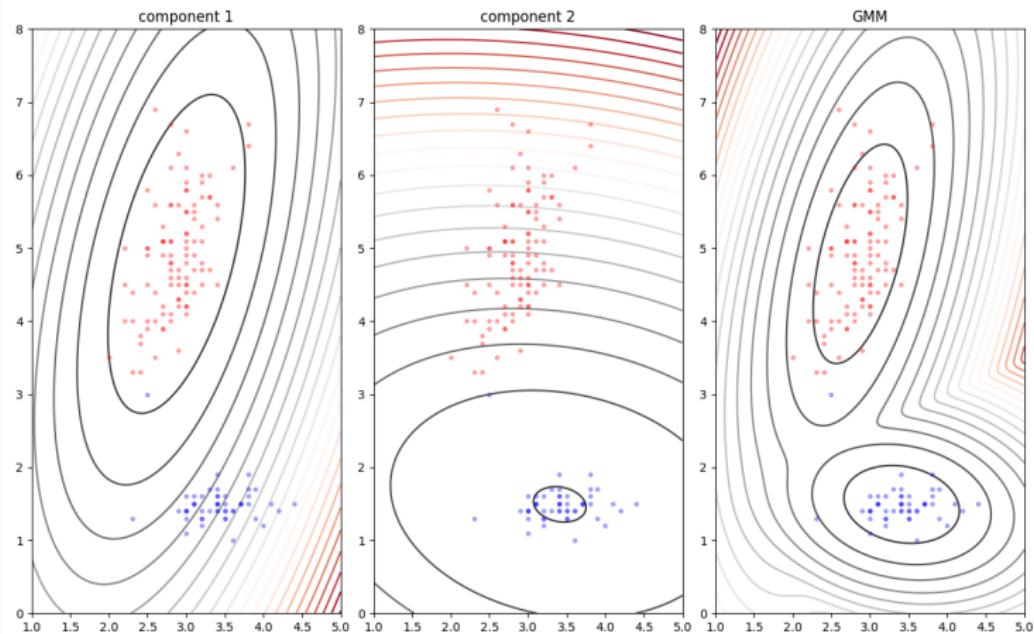


Iteration 10: log-likelihood =  $-282.92$  nats,  $w_1 = 0.64$ ,  $w_2 = 0.36$

The blue and red labels is the clustering by K-Means; this **not** used by EM

# Learning GMM with EM on Iris Data

Example

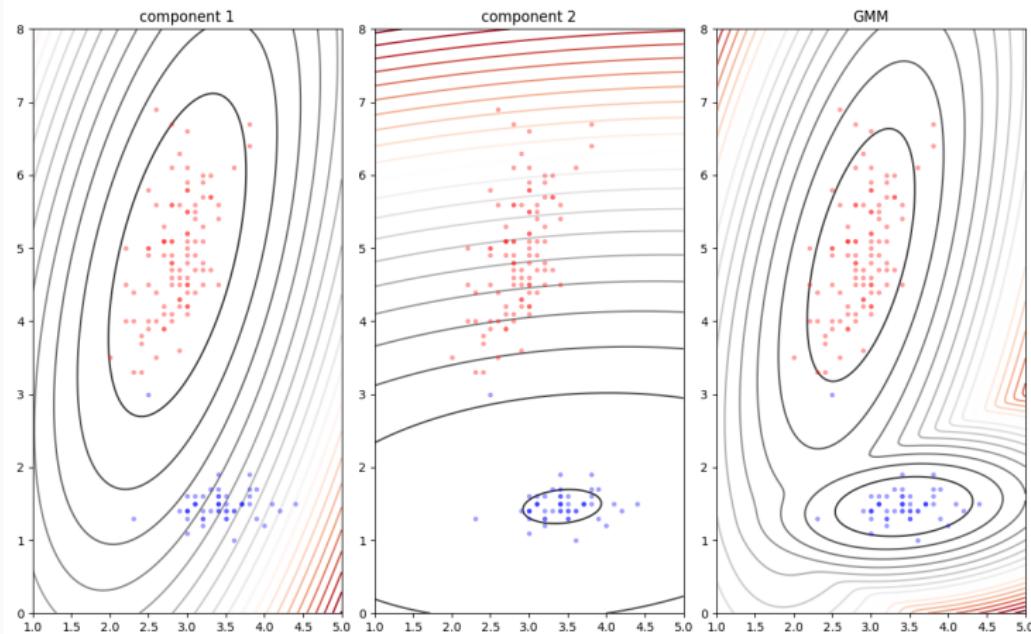


Iteration 15: log-likelihood =  $-248.68$  nats,  $w_1 = 0.66$ ,  $w_2 = 0.34$

The blue and red labels is the clustering by K-Means; this **not** used by EM

# Learning GMM with EM on Iris Data

Example

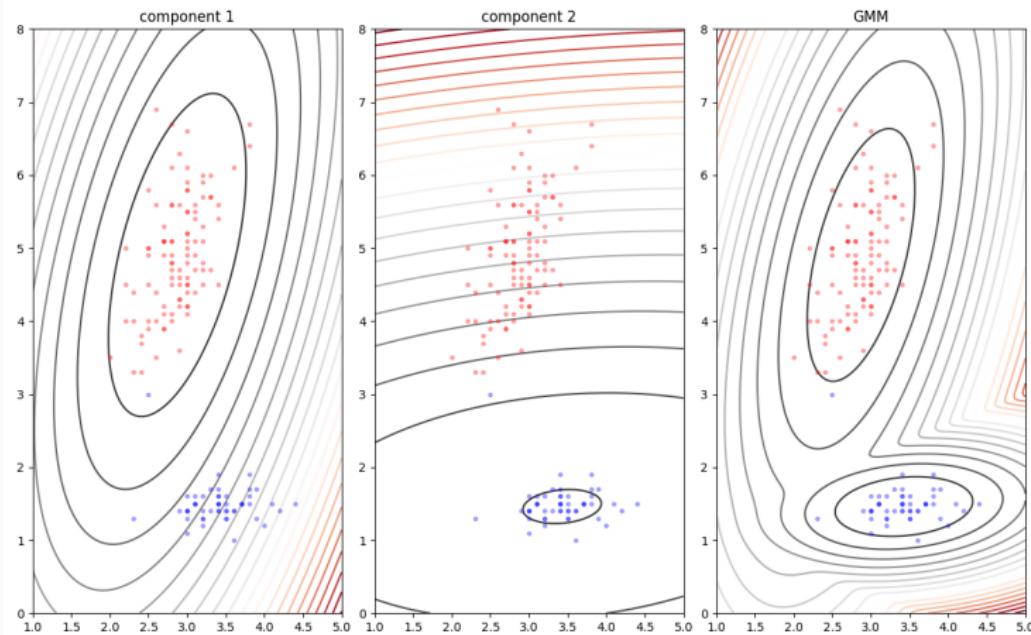


Iteration 20: log-likelihood =  $-237.35$  nats,  $w_1 = 0.67$ ,  $w_2 = 0.33$

The blue and red labels is the clustering by K-Means; this **not** used by EM

# Learning GMM with EM on Iris Data

Example



Iteration 25: log-likelihood =  $-237.35$  nats,  $w_1 = 0.67$ ,  $w_2 = 0.33$

The blue and red labels is the clustering by K-Means; this **not** used by EM

## Properties of the EM Algorithm

- **Monotonous in log-likelihood:** In each iteration, the GMM's log-likelihood is **guaranteed** to increase or stay the same (when converged)
- Thus, we can terminate the algorithm when improvement is small (stopping criterion)
- However, for GMMs, the log-likelihood is non-concave in the parameters, thus only a **local maximum** can be achieved in general
- EM tends to **increase log-likelihood rapidly in the beginning**; close to an optimum it has a rather slow theoretical convergence

- Gaussian distributions might be a bad fit
- Gaussian mixture models: universal approximators of densities
- Related to K-means: hard vs. soft clustering
- EM algorithm analogous to K-means algorithm
  - assign data points to clusters
  - update parameters
- EM algorithm monotonous in (log-)likelihood, but converges to local optimum only

# Please Evaluate

ML1 (VO) is nominated for an Award for Excellence in Teaching

