

# Gaussian Mixture Models

---

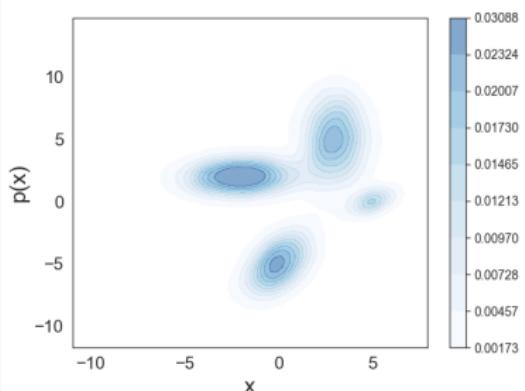
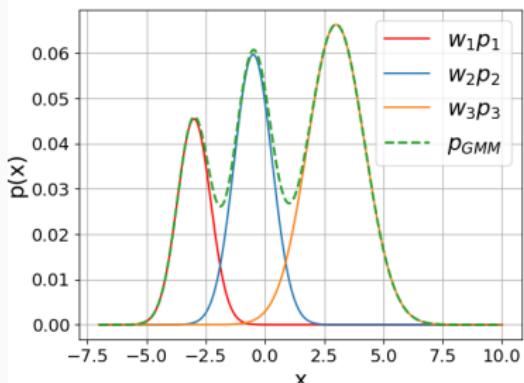
Probabilistic Decision Making — Lecture 8  
19<sup>th</sup> November 2025

Robert Peharz  
Institute of Machine Learning and Neural Computation  
Graz University of Technology

The **Gaussian mixture model (GMM)** parametrized by  $\theta = \{w_k, \mu_k, \Sigma_k\}_{k=1}^K$  is given as

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

where  $p(x | \mu_k, \Sigma_k)$  is a Gaussian with parameters  $\mu_k, \Sigma_k$ .



GMMs are **universal approximators** of densities, that is, any density can be approximated arbitrarily well by a GMM.

This is true even when:

- Gaussians are restricted to diagonal covariances,
- the diagonal is constant,
- and the mixture weights  $w_k$  are constant

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

Note the parallel to multi layer perceptrons (MLPs), which are universal approximators for continuous functions.

A simple way to “learn” a GMM is to place an **isotropic** (constant diagonal) Gaussian component on each data point.

- let  $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$  be a dataset
- the GMM defined as

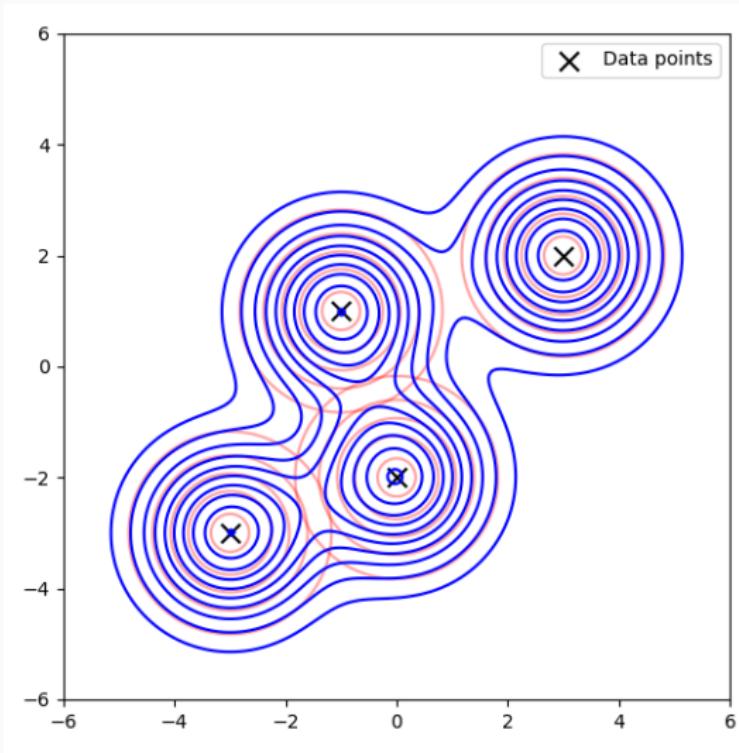
$$p_{\text{KDE}}(\mathbf{x}) = \sum_{i=1}^N \frac{1}{N} p_{\text{gauss}}(\mathbf{x} | \mathbf{x}^{(i)}, \sigma I)$$

is called **kernel density estimator (KDE)** for data  $\mathcal{D}$

- here,  $\sigma$  is some **bandwidth parameter**

# Kernel Density Estimator

Example



KDE is simple but a **consistent density estimator**. In particular, when

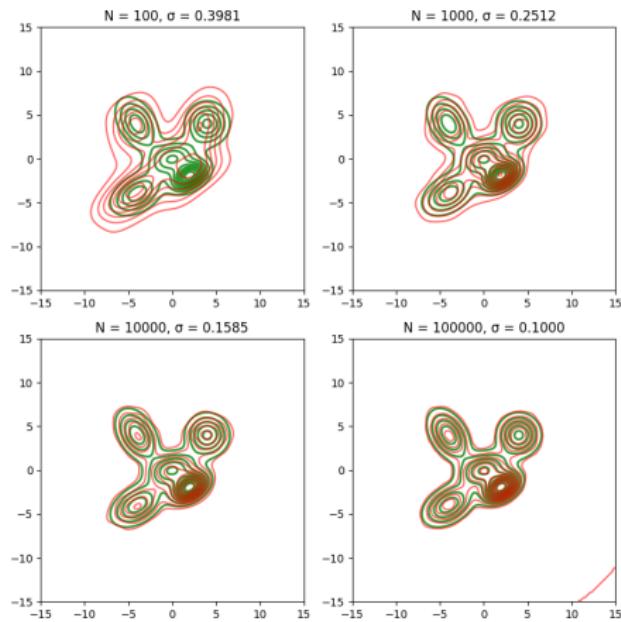
- number of data samples  $N \rightarrow \infty$
- bandwidth  $\sigma \rightarrow 0$
- but “ $\sigma$  goes to 0 slowly”, i.e.  $\sigma N \rightarrow \infty$ ,

then the KDE converges to  $p^*$ .

- this is also true for arbitrary bounded and symmetric density components, not only Gaussians
- theoretically optimal convergence rate for  $\sigma = N^{-0.2}$

# Consistency of KDE

## Example



We use a ground truth  $p^*$  (here a 5 component GMM) from which samples can be drawn (green). In red we show the KDE fit for a varying number of samples.

## GMMs as Latent Variable Models

---

## GMMs as Latent Variable Models

Recall the GMM definition:  $p(\mathbf{x}) = \sum_{k=1}^K w_k p(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$

- the weights  $w_k \geq 0$ ,  $\sum_k w_k = 1$  look like a **Categorical distribution** of some **latent** (unobserved) random variable  $Z$
- $p(\mathbf{x} | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$  might be interpreted as a conditional distribution  $p(\mathbf{x} | z = k)$

This yields a joint via the product rule:

$$p(\mathbf{x}, z) = \underbrace{p(z)}_{w_z} \underbrace{p(\mathbf{x} | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)}_{p(\mathbf{x} | z)}$$

Since  $Z$  is unobserved, we **marginalize**:

$$p(\mathbf{x}) = \sum_z p(\mathbf{x}, z) = \sum_z p(z) p(\mathbf{x} | z)$$

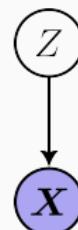
retrieving exactly the GMM definition above!

## GMMs as Latent Variable Models cont'd

The GMM might be understood as a joint over  $\mathbf{X}$  and a latent  $Z$ .

$$Z \sim \text{Categorical}(w_1, w_2, \dots, w_K)$$

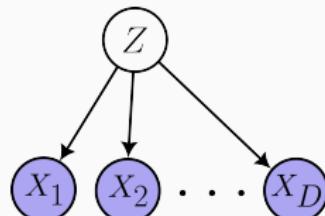
$$\mathbf{X} \sim \text{Gaussian}(\mu_Z, \Sigma_Z)$$



If the Gaussian components have diagonal covariances:

$$Z \sim \text{Categorical}(w_1, w_2, \dots, w_K)$$

$$X_d \sim \text{Gaussian}(\mu_{Zd}, \sigma_{Zd}), \quad d = 1 \dots D$$



# **Sampling GMMs**

**How can we draw samples from a GMM?**

# Sampling GMMs

## How can we draw samples from a GMM?

We combine two previous ideas:

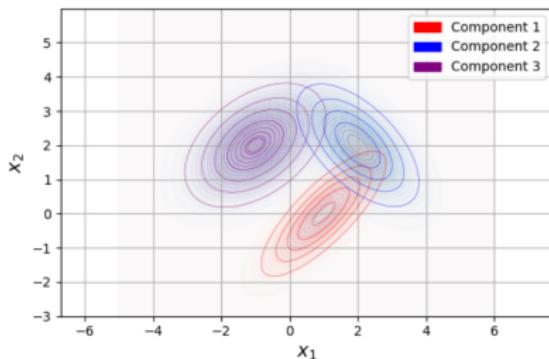
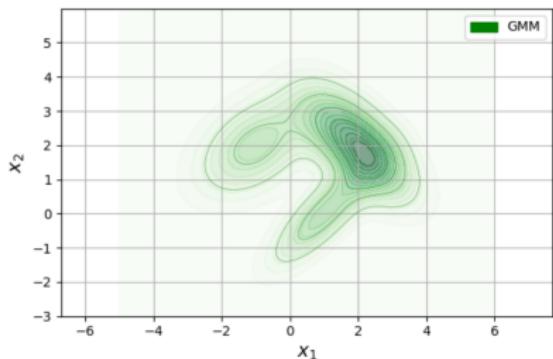
- **ancestral sampling** in BNs to draw a joint (Lecture 7)
  - sampling a Categorical  $Z$
  - sampling  $\mathbf{X}$  from the Gaussian selected by  $Z$
- having a joint sample, a **marginal** sample  $\mathbf{X}$  is obtained by discarding marginalized dimensions  $Z$  (Lecture 4)

# Sampling in GMMs

Example

Sampling a GMM?

- $Z \sim \text{Categorical}(w_1, w_2, w_3)$
- $\mathbf{X} \sim \text{Gaussian}(\mu_z, \Sigma_z)$
- discard  $Z$



$\mathbf{X}$  is an **exact sample** from  $p_{\text{GMM}}$

## Learning GMMs

---

## Learning with Latent Variables?

A GMM can be seen as a joint over  $X_1, X_2, \dots, X_D$  and  $Z$ . However, since  $Z$  is **latent**, we never observe its value, i.e. our dataset is **incomplete** and actually looks like this:

$i$	$x_1$	$x_2$	$\dots$	$x_D$	$z$
1	-0.34	0.04	...	-0.53	?
2	1.33	-0.60	...	0.94	?
:	:	:	⋮	⋮	⋮
$N$	1.64	-1.80	...	1.99	?

## Assuming Complete Data

Let us ignore the problem that  $Z$  is latent (missing) for a moment, and pretend we actually had observed it:

$i$	$x_1$	$x_2$	$\dots$	$x_D$	$z$
1	-0.34	0.04	$\dots$	-0.53	2
2	1.33	-0.60	$\dots$	0.94	1
$\vdots$	$\vdots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$
$N$	1.64	-1.80	$\dots$	1.99	2

How do we learn the GMM parameters  $\theta$  in this case?

Recall from last lecture that in Bayesian Networks, the log-likelihood decomposes over RV families:

$$\begin{aligned}\mathcal{L} &= \sum_{n=1}^N \log \prod_{i=1}^D p(x_i^{(n)} | \mathbf{x}_{\text{pa}_i}^{(n)}, \theta_i) = \sum_{n=1}^N \sum_{i=1}^D \log p(x_i^{(n)} | \mathbf{x}_{\text{pa}_i}^{(n)}, \theta_i) \\ &= \sum_{i=1}^D \underbrace{\sum_{n=1}^N \log p(x_i^{(n)} | \mathbf{x}_{\text{pa}_i}^{(n)}, \theta_i)}_{\mathcal{L}_i: \text{local log-likelihood for } X_i}\end{aligned}$$

Now, since a GMM **is** a Bayesian network over  $Z$  and  $\mathbf{X}$ , we can learn **independently**

- $p(z)$
- $p(\mathbf{x} | z)$



# Learning GMM with Complete Data

Example

Recall the MLE for Categorical:

$$\hat{w}_k = \frac{\sum_{i=1}^N \mathbb{1}[z^{(i)} = k]}{N}$$

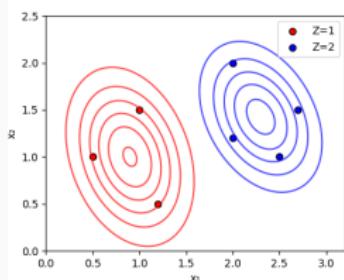
and for multivariate Gaussian

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

We fit Gaussians for each value of  $Z$ .

	$z = 1$	$z = 2$
$w$	$3/7$	$4/7$
$\mu$	$(0.9, 1)^T$	$(2.3, 1.425)^T$
$\Sigma$	$\begin{pmatrix} 0.13 & -0.05 \\ -0.05 & 0.25 \end{pmatrix}$	$\begin{pmatrix} 0.17 & -0.05 \\ -0.05 & 0.19 \end{pmatrix}$

$X_1$	$X_2$	$Z$
1.2	0.5	1
2.7	1.5	2
2	1.2	2
0.5	1	1
1	1.5	1
2	2	2
2.5	1	2



Globally optimal solution for complete data!

## How to deal with missing $Z$ 's?

However,  $Z$ 's are latent and hence missing.

A common heuristic is to “invent” suitable  $Z$ 's, for example by  
**K-means clustering**.

# A Crash-Course in K-Means

- the **K-means** algorithm aims to cluster a dataset  $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$  into  $K$  clusters (groups of similar items)
- it introduces **cluster centers**  $\mu_1, \dots, \mu_K$
- $k^{\text{th}}$  cluster is the set of samples closest to  $\mu_k$ :

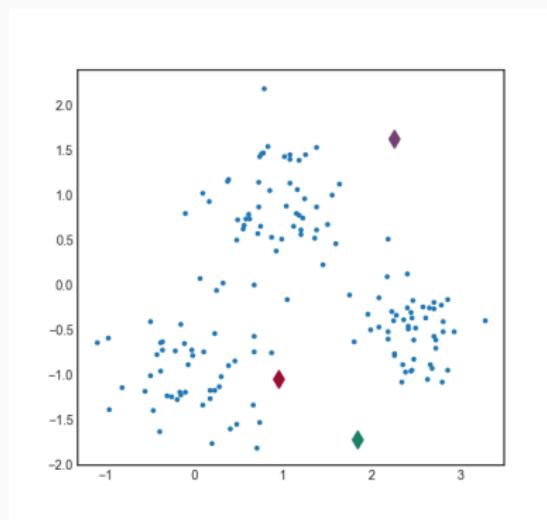
$$\mathcal{D}_k = \left\{ \mathbf{x} \in \mathcal{D} \mid k = \arg \min_I \|\mathbf{x} - \mu^{(I)}\|_2 \right\}$$

- **chicken-egg problem:**
  - if  $\mu_1, \dots, \mu_K$  were known, clusters are determined
  - if cluster were known,  $\mu_1, \dots, \mu_K$  are simply the cluster means
- K-means initializes randomly and iterates between these steps

## K-Means Algorithm (Lloyd's Algorithm)

Example

We pick  $K = 3$  and initialize  $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}$ , e.g. by drawing them uniformly between minimal and maximal range in each dimension.

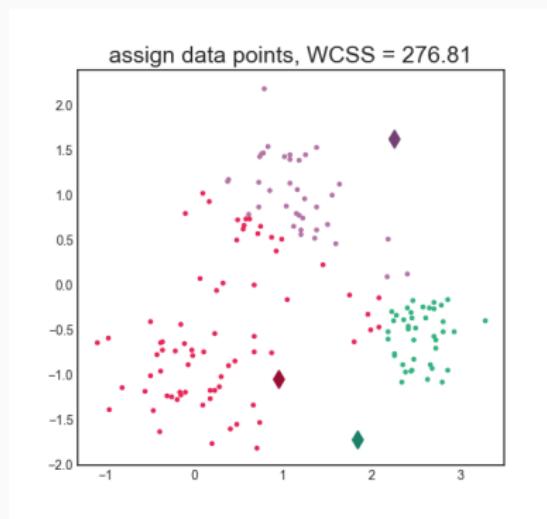


# K-Means Algorithm (Lloyd's Algorithm)

Example

Given  $\mu^{(1)}, \mu^{(2)}, \mu^{(3)}$ , we set assignment indicators

$$z_{i,k} = \begin{cases} 1 & \text{if cluster center } \mu^{(k)} \text{ is closest to } x^{(i)} \\ 0 & \text{otherwise} \end{cases}$$

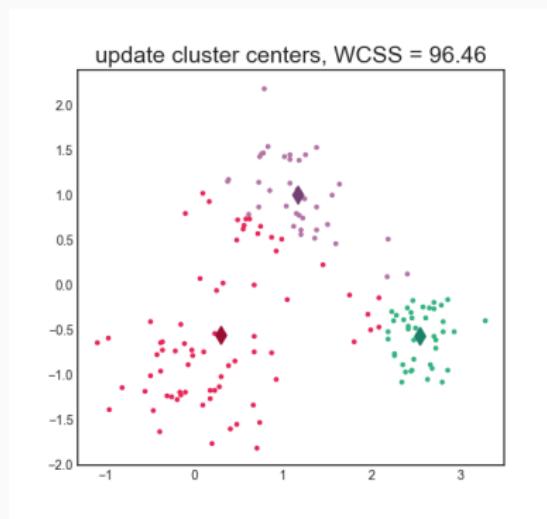


# K-Means Algorithm (Lloyd's Algorithm)

Example

Update cluster centers:

$$\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{x \in \mathcal{D}_k} x$$

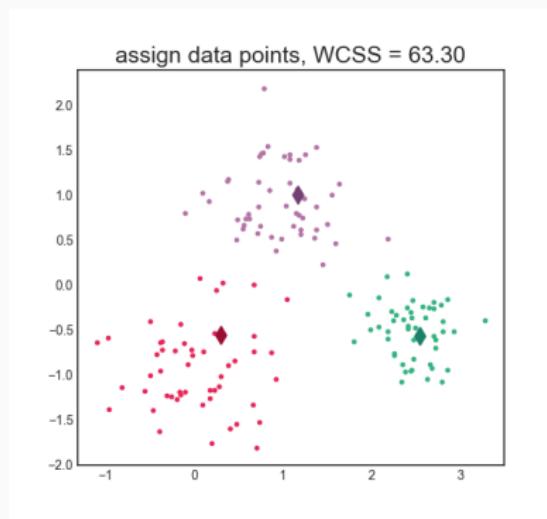


# K-Means Algorithm (Lloyd's Algorithm)

Example

Update assignments:

$$z_{i,k} = \begin{cases} 1 & \text{if cluster center } \mu^{(k)} \text{ is closest to } x^{(i)} \\ 0 & \text{otherwise} \end{cases}$$

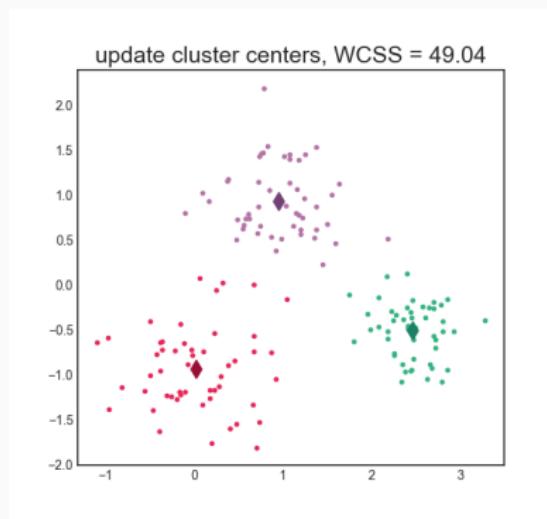


# K-Means Algorithm (Lloyd's Algorithm)

Example

Again, update mean vectors:

$$\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{x \in \mathcal{D}_k} x$$

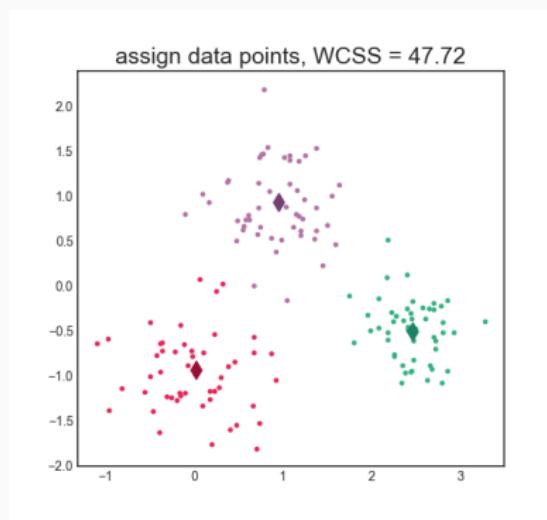


# K-Means Algorithm (Lloyd's Algorithm)

Example

Again, the data points to clusters:

$$z_{i,k} = \begin{cases} 1 & \text{if cluster center } \mu^{(k)} \text{ is closest to } x^{(i)} \\ 0 & \text{otherwise} \end{cases}$$



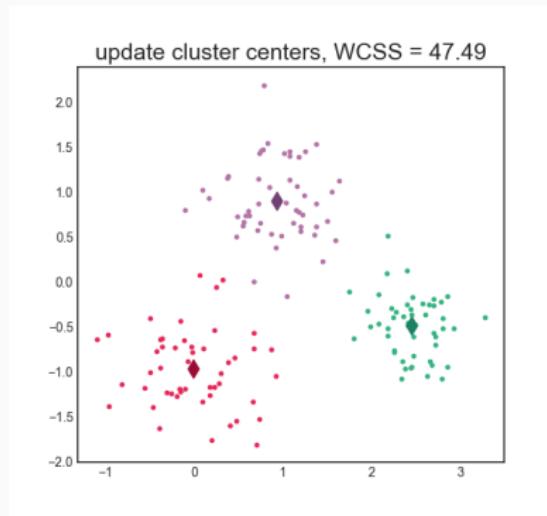
# K-Means Algorithm (Lloyd's Algorithm)

Example

Update means:

$$\mu^{(k)} \leftarrow \frac{1}{|\mathcal{D}_k|} \sum_{x \in \mathcal{D}_k} x$$

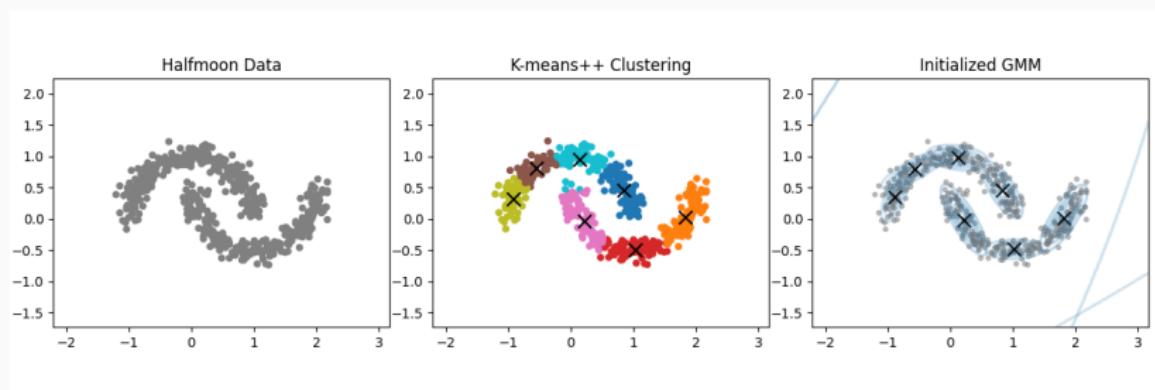
After that, the assignments  $Z$  won't change any more. Hence, also  $\mu^{(k)}$  won't change any more. The algorithm has **converged**.



# K-Means Initialization

## K-means initialization for GMMs:

- apply clustering to data  $\mathcal{D} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(N)}\}$
- use the cluster indices as proxy for  $Z$
- learn GMM on this completed dataset



This is a common and powerful initialization for GMMs.

# Maximum Likelihood under Missing Data

- when  $Z$ 's are observed, we got a simple closed-form solution for the **complete data log-likelihood**  $\sum_i \log p(\mathbf{x}^{(i)}, z^{(i)})$ 
  - MLE for  $Z \sim \text{Categorical}$
  - MLE for conditional  $\mathbf{X}|Z \sim \text{Gaussian}$
- K-means heuristic: impute missing  $Z$ 's and then apply complete data MLE
- however, in fact  $Z$ 's are not observed and the goal would be to optimize the **marginal log-likelihood**:

$$\sum_{i=1}^N \log p(\mathbf{x}^{(i)}) = \sum_{i=1}^N \log \sum_z \underbrace{p(z)}_{w_z} \underbrace{p(\mathbf{x}^{(i)} | z)}_{p(\mathbf{x}^{(i)} | \boldsymbol{\mu}_z, \boldsymbol{\Sigma}_z)},$$

# Optimizing the Marginal Log-Likelihood

- the (marginal) log-likelihood in GMMs is non-convex
- finding a global optimum: **NP-hard** problem
- first idea:
  - initialize parameters randomly
  - apply gradient ascent
- this idea indeed works, but care needs to be taken about constraints

$$w_k \geq 0 \quad \sum_k w_k = 1 \quad (\text{softmax re-parameterization})$$

$$\Sigma_k \text{ positive definite} \quad (\text{low-rank + diagonal re-parameterization})$$

# Expectation-Maximization

- **expectation maximization (EM)** is a general and powerful approach to MLE under missing data
- it applies to **any parametric model** and **any missingness pattern**
- does not introduce hyperparameters such as stepsize
- updates are **monotonic** in the marginal log-likelihood, i.e. they can only increase the log-likelihood
- in GMMs, it yields simple update rules and takes care about constraints
- coincidentally, it is very similar to K-means

## Expectations Maximization Setting (in general)

- assume some parametric joint model (not necessarily GMM):

$$p(\mathbf{x}, \mathbf{z} | \boldsymbol{\theta})$$

over **observed** RVs  $\mathbf{X}$  and **missing/latent** RVs  $\mathbf{Z}$

- in general, the **missingness pattern does not need to be fixed**—the role of observed and missing RVs is allowed to change between samples! (here we assume it fixed though, for simplicity)
- assumption:** the **complete log-likelihood** is easy to optimize, i.e. we can solve

$$\arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \boldsymbol{\theta})$$

# Expectations Maximization (in general)

- **initialize** parameters  $\theta$  (e.g. randomly, heuristically)
- **iterate** until marginal log-likelihood converges:
  - **step 1** — using the current model, **infer** the posterior over missing RVs  $Z^{(i)}$  given observed variables  $X^{(i)}$ :

$$p(z^{(i)} | x^{(i)}, \theta_{\text{cur}})$$

- **step 2** — use posterior to compute **expected complete log-likelihood**:

$$Q(\theta, \theta_{\text{cur}}) = \sum_{i=1}^N \mathbb{E}_{z^{(i)} | x^{(i)}, \theta_{\text{cur}}} [\log p(x^{(i)}, z^{(i)} | \theta)]$$

- **step 3** — maximize:

$$\theta_{\text{new}} = \arg \max_{\theta} Q(\theta, \theta_{\text{cur}}) \text{ and let } \theta_{\text{cur}} \leftarrow \theta_{\text{new}}$$

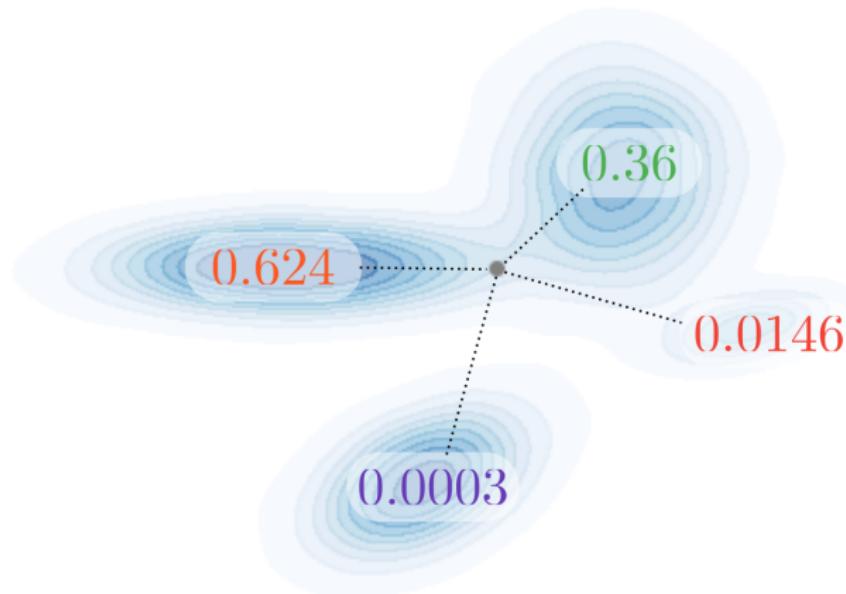
## Expectation Maximization (for GMMs)

**Step 1 — compute posterior**  $p(z^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}_{\text{cur}})$

In GMMs, there is just a single missing  $Z \in \{1, \dots, K\}$ . Bayes rule yields a closed-form:

$$p(z^{(i)} = k | \mathbf{x}^{(i)}) = \frac{\overbrace{p(z^{(i)} = k)}^{w_k} \overbrace{p(\mathbf{x}^{(i)} | z^{(i)} = k)}^{k^{\text{th}} \text{ Gaussian eval. on } \mathbf{x}^{(i)}}}{\underbrace{\sum_{z'} p(z') p(\mathbf{x}^{(i)} | z')}_{\text{GMM eval. on } \mathbf{x}^{(i)}}}$$

Let us use a shorthand  $\gamma_{i,k} = p(z^{(i)} = k | \mathbf{x}^{(i)})$ , often called **components responsibilities**. They might be interpreted as “fractions of a sample”.



## Expectation Maximization (for GMMs) cont'd

### Step 2 — expected complete log-likelihood

$$\begin{aligned} Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{\text{cur}}) &= \sum_{i=1}^N \mathbb{E}_{\mathbf{z}^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}_{\text{cur}}} \left[ \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \boldsymbol{\theta}) \right] \\ &= \sum_{i=1}^N \sum_{k=1}^K \gamma_{i,k} \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \boldsymbol{\theta}) \end{aligned}$$

**Note:** in the case of observed  $Z$  (i.e. when GMM learning is easy) the complete log-likelihood can **be written exactly the same**, where

- $\gamma_{i,k}$  are “hard”, i.e. a one-hot encoding, indicating to which component  $\mathbf{x}^{(i)}$  belongs to
- $\gamma_{i,k}$  are fixed and do not depend on  $\boldsymbol{\theta}_{\text{cur}}$

## Expectation Maximization (for GMMs) cont'd

### Step 3 — maximize

$$Q(\boldsymbol{\theta}, \boldsymbol{\theta}_{\text{cur}}) = \sum_{i=1}^N \sum_{k=1}^K \gamma_{i,k} \log p(\mathbf{x}^{(i)}, \mathbf{z}^{(i)} | \boldsymbol{\theta})$$

This is really just a weighted version of the maximum likelihood problem for complete data. Without proof, the MLE is

$$\hat{w}_k = \frac{\sum_{i=1}^N \gamma_{i,k}}{N}$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^N \gamma_{i,k} \mathbf{x}^{(i)}}{\sum_{i=1}^N \gamma_{i,k}}$$

$$\hat{\Sigma}_k = \frac{\sum_{i=1}^N \gamma_{i,k} (\mathbf{x}^{(i)} - \hat{\mu}_k)(\mathbf{x}^{(i)} - \hat{\mu}_k)^T}{\sum_{i=1}^N \gamma_{i,k}}$$

**Note:** in the case of observed  $Z$ , the solution has actually the same form, just with “hard”  $(0, 1)$  and fixed  $\gamma_{i,k}$ 's

# Expectation Maximization Algorithm for GMMs

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

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

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

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

## E-Step

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

$$\text{let } \gamma_{i,k} \leftarrow \frac{p(\mathbf{x}^{(i)} | \mu_k, \Sigma_k) w_k}{\sum_{k'} p(\mathbf{x}^{(i)} | \mu_{k'}, \Sigma_{k'}) w_{k'}} \quad \triangleright \text{cluster responsibilities}$$

**end**

## M-Step

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

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

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

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

**end**

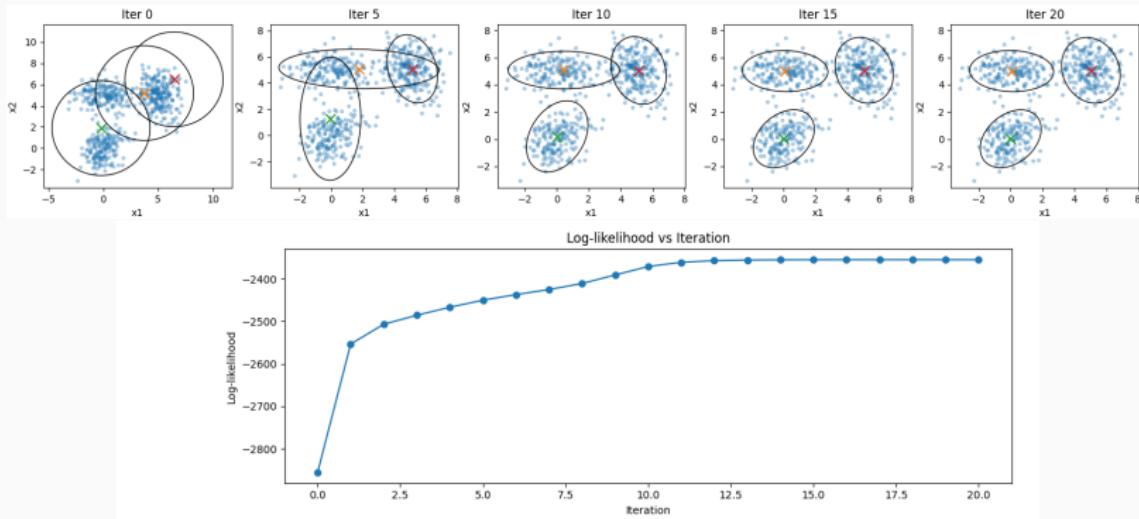
**end**

# Expectations for GMMs

Example

## Initialization

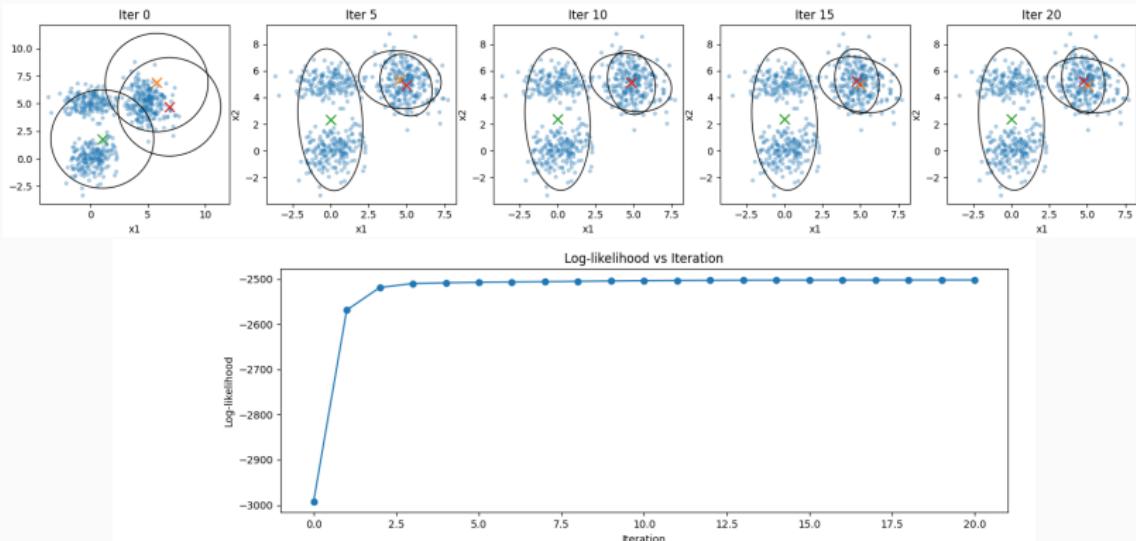
- uniform weights
- random data points as means
- identity matrix as covariances



# Bad Local Maxima

Example

- EM monotonically increases the log-likelihood
- but it can run into pretty bad local maxima
- mitigation:
  - initialization with K-means++
  - random restarts



# Inference

---

# Inference in GMMs

- GMMs are universal approximators of densities
- we discussed how they can be learned from data
- what about **inference**?



- it turns out that many inference tasks in GMMs are **tractable**, i.e. can be computed **exactly and efficiently**
- they basically “inherit” tractability from their Gaussian components

# Gaussian Marginal Distribution

Marginalizing a single Gaussian is easy: one simply needs to discard the parameters mentioning marginalized RVs.

For example, a Gaussian over 5 random variables  $X_1, X_2, X_3, X_4, X_5$ , the marginal over  $X_1, X_4, X_5$  is Gaussian with parameters  $\mu_q$  and  $\Sigma_{qq}$  obtained by deleting rows and columns corresponding to  $X_2, X_3$ :

$$\boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \\ \mu_5 \end{pmatrix} \quad \Sigma = \begin{pmatrix} v_{1,1} & v_{1,2} & v_{1,3} & v_{1,4} & v_{1,5} \\ \cancel{v_{2,1}} & \cancel{v_{2,2}} & \cancel{v_{2,3}} & \cancel{v_{2,4}} & v_{2,5} \\ \cancel{v_{3,1}} & \cancel{v_{3,2}} & \cancel{v_{3,3}} & \cancel{v_{3,4}} & v_{3,5} \\ v_{4,1} & v_{4,2} & v_{4,3} & v_{4,4} & v_{4,5} \\ v_{5,1} & v_{5,2} & v_{5,3} & v_{5,4} & v_{5,5} \end{pmatrix}$$

$$\boldsymbol{\mu}_q = \begin{pmatrix} \mu_1 \\ \mu_4 \\ \mu_5 \end{pmatrix} \quad \Sigma_{qq} = \begin{pmatrix} v_{1,1} & v_{1,4} & v_{1,5} \\ v_{4,1} & v_{4,4} & v_{4,5} \\ v_{5,1} & v_{5,4} & v_{5,5} \end{pmatrix}$$

## GMM Marginal Distribution

The marginals of GMMs are also easy:

$$\begin{aligned} p(\mathbf{x}_q) &= \int \sum_k w_k p(\mathbf{x}_q, \mathbf{x}_m | \boldsymbol{\theta}_k) d\mathbf{x}_m \\ &= \sum_k w_k \underbrace{\int p(\mathbf{x}_q, \mathbf{x}_m | \boldsymbol{\theta}_k) d\mathbf{x}_m}_{\text{Gaussian marginal}} \\ &= \sum_k w_k p(\mathbf{x}_q | \boldsymbol{\theta}_k) \end{aligned}$$

Here we actually didn't use the fact that components are Gaussian—as long as each component has tractable marginals, the mixture has tractable marginals as well.

**marginal of mixture = mixture of marginals!**

# Gaussian Conditional Distributions

In Gaussians with parameters  $\mu$  and  $\Sigma$ , any conditional distribution  $p(x_q | x_e)$  is again Gaussian with **closed form parameters**:

$$\mu_{q|e} = \mu_q + \Sigma_{qe} \Sigma_{ee}^{-1} (x_e - \mu_e)$$

$$\Sigma_{q|e} = \Sigma_{qq} - \Sigma_{qe} \Sigma_{ee}^{-1} \Sigma_{eq}$$

For query variables  $X_1, X_4, X_5$  and evidence variables  $X_2, X_3$ :

$$\mu = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \\ \mu_4 \\ \mu_5 \end{pmatrix} \quad \Sigma = \begin{pmatrix} \Sigma_{ee} & \Sigma_{qq} \\ \Sigma_{qe} & \Sigma_{eq} \end{pmatrix}$$

The diagram illustrates the decomposition of the mean vector  $\mu$  and covariance matrix  $\Sigma$  into components based on evidence variables  $X_2, X_3$  and query variables  $X_1, X_4, X_5$ . The mean vector  $\mu$  is shown as a column of five elements:  $\mu_1, \mu_2, \mu_3, \mu_4, \mu_5$ . The evidence variables  $X_2, X_3$  are represented by the first two columns of the covariance matrix  $\Sigma$ , which are colored green. The query variables  $X_1, X_4, X_5$  are represented by the last three columns of the covariance matrix  $\Sigma$ , which are colored red. The diagonal elements of the covariance matrix  $\Sigma$  represent the variances of the query variables:  $\Sigma_{ee}$  (top-left),  $\Sigma_{qq}$  (top-right),  $\Sigma_{eq}$  (bottom-left), and  $\Sigma_{qe}$  (bottom-right).

$v_{1,1}$	$v_{1,2}$	$v_{1,3}$	$v_{1,4}$	$v_{1,5}$
$v_{2,1}$	$v_{2,2}$	$v_{2,3}$	$v_{2,4}$	$v_{2,5}$
$v_{3,1}$	$v_{3,2}$	$v_{3,3}$	$v_{3,4}$	$v_{3,5}$
$v_{4,1}$	$v_{4,2}$	$v_{4,3}$	$v_{4,4}$	$v_{4,5}$
$v_{5,1}$	$v_{5,2}$	$v_{5,3}$	$v_{5,4}$	$v_{5,5}$

# GMM Conditional Distribution

The conditionals of GMMs are also quite easy. We use the fact that the conditional is proportional to the joint,  $p_{GMM}(\mathbf{x}_q | \mathbf{x}_e) \propto p_{GMM}(\mathbf{x}_q, \mathbf{x}_e)$ :

$$\begin{aligned} p_{GMM}(\mathbf{x}_q | \mathbf{x}_e) &\propto \sum_{k=1}^K w_k p(\mathbf{x}_q, \mathbf{x}_e | \boldsymbol{\theta}_k) \\ &= \sum_{k=1}^K w_k \overbrace{p(\mathbf{x}_q | \mathbf{x}_e, \boldsymbol{\theta}_k) p(\mathbf{x}_e | \boldsymbol{\theta}_k)}^{\text{chain rule}} \\ &= \sum_{k=1}^K \underbrace{w_k p(\mathbf{x}_e | \boldsymbol{\theta}_k)}_{\text{marginal Gauss}} \underbrace{p(\mathbf{x}_q | \mathbf{x}_e, \boldsymbol{\theta}_k)}_{\text{cond. Gauss}} \end{aligned}$$

We can normalize by replacing  $w_k p(\mathbf{x}_e | \boldsymbol{\theta}_k)$  with  $\tilde{w}_k = \frac{w_k p(\mathbf{x}_e | \boldsymbol{\theta}_k)}{\sum_i w_i p(\mathbf{x}_e | \boldsymbol{\theta}_i)}$ .  
Hence

$$p_{GMM}(\mathbf{x}_q | \mathbf{x}_e) = \sum_{k=1}^K \tilde{w}_k p(\mathbf{x}_q | \mathbf{x}_e, \boldsymbol{\theta}_k)$$

## Conditional GMM cont'd

$$p_{GMM}(\mathbf{x}_q | \mathbf{x}_e) = \sum_{k=1}^K \tilde{w}_k p_k(\mathbf{x}_q | \mathbf{x}_e, \boldsymbol{\theta}_k)$$

Again, we did not use the fact that the components are Gaussian—only that we know how to marginalize and condition them.

**Conditional of mixture = a mixture of conditionals!**

... with modified weights

$$\tilde{w}_k = \frac{w_k p(\mathbf{x}_e | \boldsymbol{\theta}_k)}{\sum_i w_i p(\mathbf{x}_e | \boldsymbol{\theta}_i)}$$

# On GANs and GMMs

**Eitan Richardson**

School of Computer Science and Engineering  
The Hebrew University of Jerusalem  
Jerusalem, Israel  
`eitanrich@cs.huji.ac.il`

**Yair Weiss**

School of Computer Science and Engineering  
The Hebrew University of Jerusalem  
Jerusalem, Israel  
`yweiss@cs.huji.ac.il`

- this paper demonstrates that good old GMMs can deliver decent generative models
- uses **low-rank + diagonal** covariances

<https://github.com/eitanrich/gans-n-gmms>

# Implementation

## Covariances

- each covariance matrix is represented as

$$\Sigma = AA^T + \text{diag}(\mathbf{d})$$

where  $A = \mathbb{R}^{D \times L}$  with  $L < D$

- $\Sigma$  constructed in this way is always positive definite
- adding a diagonal matrix  $\text{diag}(\mathbf{d})$  with strictly positive elements in the diagonal ensures that  $\Sigma$  is invertible

## Training

Maximum likelihood with Adam on GPU. Later versions also include an EM implementation.

# Samples

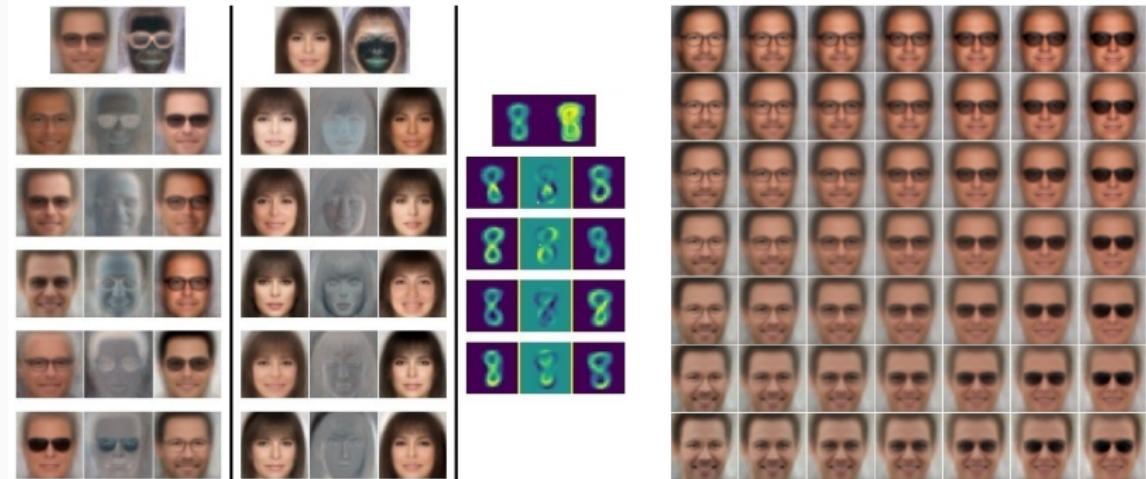
# Example



Samples on *CelebA*, *MNIST*, and *Street View House Numbers* are decent, albeit a bit blurry and oversmoothed

# Factor Analysis

Example



Left: mean and diagonal of a Gaussian components; each row visualizes a column of the  $A$ -matrix (mean minus column, column, mean plus column)

Right: superpositions of two columns

# Outlier Detection

Example



By selecting data samples with low  $\log p$  one can detect outliers and peculiar examples.

## Image Inpainting

Example



A GMM is one of the few models that allow **tractable inference**, such as computing **exact conditionals**. This can be used for **data restoration** and **inpainting**, by sampling from the conditional, or by taking the conditional expectation.