# **Unsupervised Probabilistic Learning with Latent Variables**

Lecture 1: Exact Expectation-Maximisation Algorithms and Variational Autoencoders

Machine Learning Summer School Africa 2023, Cape Town

Xavier Alameda-Pineda

RobotLearn team, Inria at University Grenoble-Alpes, Jean-Kuntzman CNRS Laboratory, Multidisciplinary Institute of Artificial Intelligence









Imagine a system for {person tracking, speech denoising, body pose estimation,  $\ldots$ } in:



[Images under Creative Commons license]

We would like to avoid re-annotating for every new environment  $\rightarrow$  Interest in unsupervised learning (and/or unsupervised domain adaptation).

# And probabilistic learning?

Probabilistic generative models aim to learn a (parametric) distribution  $p_{\theta}(x)$  that approximates the complex data distribution  $p_{data}(x)$ :



- Once learned, we can (ideally) sample new data.
- We can jointly learn them with other probabilistic models using *maximum likelihood*.

## The Kullback-Leibler divergence and the ML formulation

The Kullback-Leilbler (KL) divergence between two distributions writes:

$$D_{\mathsf{KL}}(p(\boldsymbol{x}) \| q(\boldsymbol{x})) = -\mathbb{E}_{p(\boldsymbol{x})} \left[ \log \frac{q(\boldsymbol{x})}{p(\boldsymbol{x})} \right] = -\int_{\mathcal{X}} p(\boldsymbol{x}) \log \frac{q(\boldsymbol{x})}{p(\boldsymbol{x})} \mathrm{d}\boldsymbol{x} \left\{ \begin{array}{l} \geq 0 \\ \neq D_{\mathsf{KL}}(q(\boldsymbol{x}) \| p(\boldsymbol{x})) \\ \mathcal{D}_{\mathsf{KL}}(p(\boldsymbol{x}) \| q(\boldsymbol{x})) = 0 \Leftrightarrow p(\boldsymbol{x}) = q(\boldsymbol{x}). \end{array} \right.$$

#### The Kullback-Leibler divergence and the ML formulation

The Kullback-Leilbler (KL) divergence between two distributions writes:

$$D_{\mathsf{KL}}(p(\boldsymbol{x}) \| q(\boldsymbol{x})) = -\mathbb{E}_{p(\boldsymbol{x})} \left[ \log \frac{q(\boldsymbol{x})}{p(\boldsymbol{x})} \right] = -\int_{\mathcal{X}} p(\boldsymbol{x}) \log \frac{q(\boldsymbol{x})}{p(\boldsymbol{x})} \mathrm{d}\boldsymbol{x} \begin{cases} \geq 0\\ \neq D_{\mathsf{KL}}(q(\boldsymbol{x}) \| p(\boldsymbol{x})) \end{cases}$$
$$\mathcal{D}_{\mathsf{KL}}(p(\boldsymbol{x}) \| q(\boldsymbol{x})) = 0 \Leftrightarrow p(\boldsymbol{x}) = q(\boldsymbol{x}).$$

Given a training set  $\{x_i\}_{i=1}^N, x_i \sim p_{\scriptscriptstyle \mathsf{data}}(x)$ , ML minimizes the KL divergence:

$$\begin{aligned} \boldsymbol{\theta}^* &= \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \quad D_{\mathsf{KL}} \left( p_{\mathsf{data}}(\boldsymbol{x}) \parallel p_{\boldsymbol{\theta}}(\boldsymbol{x}) \right) = \underset{\boldsymbol{\theta}}{\operatorname{argmin}} \quad - \mathbb{E}_{p_{\mathsf{data}}} \left[ \log \frac{p_{\boldsymbol{\theta}}(\boldsymbol{x})}{p_{\mathsf{data}}(\boldsymbol{x})} \right] \\ &= \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \quad \mathbb{E}_{p_{\mathsf{data}}} \left[ \log p_{\boldsymbol{\theta}}(\boldsymbol{x}) \right] \approx \left[ \underset{\boldsymbol{\theta}}{\operatorname{argmax}} \quad \frac{1}{N} \sum_{i=1}^{N} \log p_{\boldsymbol{\theta}}(\boldsymbol{x}_i) \right] \end{aligned}$$

## **Interest of Latent Variables**

- Speech enhancement: noisy speech (observation), clean speech (latent variable)
   → Get rid of noise heard together with speech.
- Person tracking: detections (observation), person positions (latent variable)
   → Track occluded persons, filter over time.
- Representation learning: raw data (observation), representation (latent variable)
  - $\rightarrow$  Compute an optimal compact representation of the raw data.

## **Interest of Latent Variables**

- Speech enhancement: noisy speech (observation), clean speech (latent variable)
   → Get rid of noise heard together with speech.
- Person tracking: detections (observation), person positions (latent variable)
   → Track occluded persons, filter over time.
- Representation learning: raw data (observation), representation (latent variable)  $\rightarrow$  Compute an optimal compact representation of the raw data.

Let z denote the latent variable:

$$\begin{cases} \boldsymbol{z} \sim p_{\boldsymbol{\theta}}(\boldsymbol{z}) \\ \boldsymbol{x} | \boldsymbol{z} \sim p_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{z}) \end{cases} \rightarrow p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \int p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z}) \mathrm{d}\boldsymbol{z} = \int p_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{z}) p_{\boldsymbol{\theta}}(\boldsymbol{z}) \mathrm{d}\boldsymbol{z}$$

**New samples:** Draw  $z_k \sim p_{\theta}(z)$ , then draw a new sample  $x_k \sim p_{\theta}(x|z_k)$ .

## **Interest of Latent Variables**

- Speech enhancement: noisy speech (observation), clean speech (latent variable)
   → Get rid of noise heard together with speech.
- Person tracking: detections (observation), person positions (latent variable)
   → Track occluded persons, filter over time.
- Representation learning: raw data (observation), representation (latent variable)  $\rightarrow$  Compute an optimal compact representation of the raw data.

Let z denote the latent variable:

$$\begin{cases} \boldsymbol{z} \sim p_{\boldsymbol{\theta}}(\boldsymbol{z}) \\ \boldsymbol{x} | \boldsymbol{z} \sim p_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{z}) \end{cases} \rightarrow p_{\boldsymbol{\theta}}(\boldsymbol{x}) = \int p_{\boldsymbol{\theta}}(\boldsymbol{x}, \boldsymbol{z}) \mathrm{d}\boldsymbol{z} = \int p_{\boldsymbol{\theta}}(\boldsymbol{x} | \boldsymbol{z}) p_{\boldsymbol{\theta}}(\boldsymbol{z}) \mathrm{d}\boldsymbol{z}$$

New samples: Draw  $z_k \sim p_{\theta}(z)$ , then draw a new sample  $x_k \sim p_{\theta}(x|z_k)$ .

**Learning:** How to estimate  $\theta$  with the ML formulation?

- Lecture 1: Probabilistic learning with latent variables, GMM and the exact EM algorithm, probabilistic PCA and VAE.
- Lecture 2: The variational EM algoritm, mixtures of VAEs, application to audio-visual speech enhancement.
- Lecture 3: Dynamical variational autoencoders, application to speech enhancement.

Connections to other lectures:

- Matthias Bauer: from PPCA/VAE to diffusion models.
- Steve Kroon: normalising flows and their use with VAEs.

I am fully available for discussing/chatting with y'all!

1 Learning with Latent Variables: the Gaussian Mixture Model

- EM for GMM and beyond
- Probabilistic PCA and VAE
- 4 VAE for audio modeling

# Learning with Latent Variables: the Gaussian Mixture Model

## Simple example: clustering

Definition: find groups of data points without labels.



# Very intuitive algorithm

Very simple algorithm (called the *K*-means algorithm):

- Initialise randomly K centroids.
- Assign each data point to the closes centroid.
- Secompute centroids from the assignments.
- Iterate the past two steps.



<sup>[</sup>Images from https://ai.plainenglish.io/]

## Important points of K-means

The point-to-cluster assignment variable is unknown (latent, z), and must be inferred with the centroids (parameters of the model).

The assignment criterion is the Euclidean distance  $\Rightarrow$  groups are spherical and equally populated a priori.



# Generalising *K*-means: GMM

Defining the Gaussian mixture model (GMM)

• For each  $x_n$  there is a latent variable  $z_n$  taking values from 1 to K:  $z_n \in \{1, \ldots, K\}$ .

# Generalising *K*-means: GMM

Defining the Gaussian mixture model (GMM)

- For each  $x_n$  there is a latent variable  $z_n$  taking values from 1 to K:  $z_n \in \{1, \dots, K\}$ .
- Its prior probability is defined as:  $p(\boldsymbol{z}_n = k) = \pi_k \ge 0$ , with  $\sum_{k=1}^{K} \pi_k = 1$ .

# Generalising *K*-means: GMM

Defining the Gaussian mixture model (GMM)

- For each  $x_n$  there is a latent variable  $z_n$  taking values from 1 to K:  $z_n \in \{1, \dots, K\}$ .
- Its prior probability is defined as:  $p(\boldsymbol{z}_n = k) = \pi_k \ge 0$ , with  $\sum_{k=1}^{K} \pi_k = 1$ .
- Given  $z_n$ , the data point is modeled as a multivariate Gaussian:

$$p(\boldsymbol{x}_n | \boldsymbol{z}_n = k) = \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Advantages:

- Having  $\pi_1, \ldots, \pi_K$  means that groups can be differently populated.
- 2 The shape of the groups is modeled by  $\Sigma_k$ .
- The parameters are:  $\boldsymbol{\theta} = \{\pi_k, \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k\}_{k=1}^K$ .

Let's compute  $p(\boldsymbol{x}_n)$ 

$$p(\boldsymbol{x}_n) = \sum_{k=1}^{K} p(\boldsymbol{x}_n, \boldsymbol{z}_n = k) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k).$$

The log-likelihood:

$$\mathcal{L}(oldsymbol{ heta}|oldsymbol{X}) = \sum_{n=1}^N \log \sum_{k=1}^K \pi_k \mathcal{N}(oldsymbol{x}_n;oldsymbol{\mu}_k,oldsymbol{\Sigma}_k),$$

Computing directly ML by  $\frac{\partial \mathcal{L}}{\partial \pi_k} = 0$ ,  $\frac{\partial \mathcal{L}}{\partial \mu_k} = 0$  or  $\frac{\partial \mathcal{L}}{\partial \Sigma_k} = 0$  is very difficult.

# EM for GMM and beyond

We have seen that  $\log p(\boldsymbol{x})$  does not work well with derivatives. However,  $\log p(\boldsymbol{x}, \boldsymbol{z})$  does!

**Problem**:  $\boldsymbol{z}$  is not observed, thus  $\sum_n \log p(\boldsymbol{x}_n, \boldsymbol{z}_n)$  is a random variable.

We have seen that  $\log p(x)$  does not work well with derivatives. However,  $\log p(x, z)$  does!

**Problem**:  $\boldsymbol{z}$  is not observed, thus  $\sum_n \log p(\boldsymbol{x}_n, \boldsymbol{z}_n)$  is a random variable.

Given an initial value of the parameters,  $\theta^0$ , let's take the expectation w.r.t. the posterior distribution  $p(\boldsymbol{z}|\boldsymbol{x}, \theta^0)$  (we will justify this choice later on):

$$\mathcal{Q}(oldsymbol{ heta},oldsymbol{ heta}^0) = \sum_n \mathbb{E}_{p(oldsymbol{z}_n | oldsymbol{x}_n;oldsymbol{ heta}^0)} \log p(oldsymbol{x}_n,oldsymbol{z}_n;oldsymbol{ heta})$$

This function is called: *expected complete-data log-likelihood*, and is the main mathematical object when working with EM algorithms.

# The EM algorithm for GMM

Given  $\theta^0$ , we use the expected complete-data log-likelihood Q. For iteration  $r = 1, \ldots, R$ :

Expectation [E-step]:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{r-1}) = \mathbb{E}_{p(\boldsymbol{Z}|\boldsymbol{X}; \boldsymbol{\theta}^{r-1})} \log p(\boldsymbol{X}, \boldsymbol{Z}; \boldsymbol{\theta})$$

Maximisation [M-step]:

$$\boldsymbol{\theta}^r = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{r-1})$$

(observations  $X = \{x_n\}_{n=1}^N$ , latent variables  $Z = \{z_n\}_{n=1}^N$ , parameters  $\theta = \{\pi_k, \mu_k, \Sigma_k\}_{k=1}^K$ )

# The EM algorithm for GMM

Given  $\theta^0$ , we use the expected complete-data log-likelihood Q. For iteration  $r = 1, \ldots, R$ :

Expectation [E-step]:

$$\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{r-1}) = \mathbb{E}_{p(\boldsymbol{Z}|\boldsymbol{X}; \boldsymbol{\theta}^{r-1})} \log p(\boldsymbol{X}, \boldsymbol{Z}; \boldsymbol{\theta})$$

Maximisation [M-step]:

$$\boldsymbol{\theta}^r = \arg \max_{\boldsymbol{\theta}} \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^{r-1})$$

(observations  $m{X} = \{m{x}_n\}_{n=1}^N$ , latent variables  $m{Z} = \{m{z}_n\}_{n=1}^N$ , parameters  $m{ heta} = \{\pi_k, m{\mu}_k, m{\Sigma}_k\}_{k=1}^K$ )

We can look back to K-means:

- Infer latent variables (assignment) given the parameters (centroids) [E-step].
- Stimate the parameters (centroids) given the assignments [M-step].

# The EM algorithm for GMM (II)

**E-step** The posterior distribution writes:

$$p(\boldsymbol{z}=k|\boldsymbol{x}_n;\boldsymbol{\theta}^0) = \frac{p(\boldsymbol{z}_n=k;\boldsymbol{\theta}^0)p(\boldsymbol{x}_n|\boldsymbol{z}_n=k;\boldsymbol{\theta}^0)}{\sum_{\ell} p(\boldsymbol{z}_n=\ell;\boldsymbol{\theta}^0)p(\boldsymbol{x}_n|\boldsymbol{z}_n=\ell;\boldsymbol{\theta}^0)} = \frac{\pi_k^0 \mathcal{N}(\boldsymbol{x}_n;\boldsymbol{\mu}_k^0,\boldsymbol{\Sigma}_k^0)}{\sum_{\ell} \pi_\ell^0 \mathcal{N}(\boldsymbol{x}_n;\boldsymbol{\mu}_\ell^0,\boldsymbol{\Sigma}_\ell^0)} = \eta_{nk}$$

# The EM algorithm for GMM (II)

**E-step** The posterior distribution writes:

$$p(\boldsymbol{z}=k|\boldsymbol{x}_n;\boldsymbol{\theta}^0) = \frac{p(\boldsymbol{z}_n=k;\boldsymbol{\theta}^0)p(\boldsymbol{x}_n|\boldsymbol{z}_n=k;\boldsymbol{\theta}^0)}{\sum_{\ell} p(\boldsymbol{z}_n=\ell;\boldsymbol{\theta}^0)p(\boldsymbol{x}_n|\boldsymbol{z}_n=\ell;\boldsymbol{\theta}^0)} = \frac{\pi_k^0 \mathcal{N}(\boldsymbol{x}_n;\boldsymbol{\mu}_k^0,\boldsymbol{\Sigma}_k^0)}{\sum_{\ell} \pi_\ell^0 \mathcal{N}(\boldsymbol{x}_n;\boldsymbol{\mu}_\ell^0,\boldsymbol{\Sigma}_\ell^0)} = \eta_{nk}$$

The expected complete-data log-likelihood writes:

$$\mathcal{Q}(\boldsymbol{ heta}, \boldsymbol{ heta}^0) = \sum_{n=1}^{N} \sum_{k=1}^{K} \eta_{nk} \log \pi_k \mathcal{N}(\boldsymbol{x}_n; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

## The EM algorithm for GMM (II)

**E-step** The posterior distribution writes:

$$p(\boldsymbol{z}=k|\boldsymbol{x}_n;\boldsymbol{\theta}^0) = \frac{p(\boldsymbol{z}_n=k;\boldsymbol{\theta}^0)p(\boldsymbol{x}_n|\boldsymbol{z}_n=k;\boldsymbol{\theta}^0)}{\sum_{\ell} p(\boldsymbol{z}_n=\ell;\boldsymbol{\theta}^0)p(\boldsymbol{x}_n|\boldsymbol{z}_n=\ell;\boldsymbol{\theta}^0)} = \frac{\pi_k^0 \mathcal{N}(\boldsymbol{x}_n;\boldsymbol{\mu}_k^0,\boldsymbol{\Sigma}_k^0)}{\sum_{\ell} \pi_\ell^0 \mathcal{N}(\boldsymbol{x}_n;\boldsymbol{\mu}_\ell^0,\boldsymbol{\Sigma}_\ell^0)} = \eta_{nk}$$

The expected complete-data log-likelihood writes:

$$\mathcal{Q}(oldsymbol{ heta},oldsymbol{ heta}^0) = \sum_{n=1}^N \sum_{k=1}^K \eta_{nk} \log \pi_k \mathcal{N}(oldsymbol{x}_n;oldsymbol{\mu}_k,oldsymbol{\Sigma}_k)$$

**M-step** The optimal value for  $\mu_k$  writes:

$$oldsymbol{\mu}_k^* = rac{1}{\sum_{n=1}^N \eta_{nk}} \sum_{n=1}^N \eta_{nk} oldsymbol{x}_n$$

The main mathematical object in EM is  ${\cal Q}$  (The expected complete-data log-likelihood).

What is the relationship with the log-likelihood? Let's take any distribution of z, q(z):

$$\log p(\boldsymbol{x}) = \mathbb{E}_{q(\boldsymbol{z})} \left[ \log p(\boldsymbol{x}) \right]$$
$$= \mathbb{E}_{q(\boldsymbol{z})} \left[ \log p(\boldsymbol{x}) \frac{p(\boldsymbol{z}|\boldsymbol{x})q(\boldsymbol{z})}{p(\boldsymbol{z}|\boldsymbol{x})q(\boldsymbol{z})} \right]$$
$$= \underbrace{\mathbb{E}_{q(\boldsymbol{z})} \left[ \log \frac{p(\boldsymbol{x})p(\boldsymbol{z}|\boldsymbol{x})}{q(\boldsymbol{z})} \right]}_{\text{M-step - }\mathcal{Q}} + \underbrace{D_{\text{KL}} \left( q(\boldsymbol{z}) \left\| p(\boldsymbol{z}|\boldsymbol{x}) \right)}_{\text{E-step - KL}} \right]$$

## But why does the EM work? (II)

$$\log p(\boldsymbol{x}; \boldsymbol{\theta}) = \underbrace{\mathbb{E}_{q(\boldsymbol{z})} \left[ \log \frac{p(\boldsymbol{x})p(\boldsymbol{z}|\boldsymbol{x})}{q(\boldsymbol{z})} \right]}_{\text{M-step - }\mathcal{Q}} + \underbrace{D_{\text{KL}} \left( q(\boldsymbol{z}) \left\| p(\boldsymbol{z}|\boldsymbol{x}) \right)}_{\text{E-step - }\text{KL}} \right)$$

#### Another interpretation. Given $\theta^0$ :

- Set  $q(z) = p(z|x; \theta^0)$ , minimise KL (thus maximize Q) w.r.t. q(z). The E-step reduces the distance between log-likelihood and Q.
- $\textbf{ aximize } \mathcal{Q} \text{ w.r.t. } \boldsymbol{\theta} \text{:} \qquad \qquad \mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0) = \mathbf{E}_{q(\boldsymbol{z})} \Big[ \log \frac{p(\boldsymbol{x}, \boldsymbol{z}; \boldsymbol{\theta})}{q(\boldsymbol{z})} \Big].$

The M-step pushes  ${\mathcal Q}$  and therefore pushes the log-likelihood.

## The Exact EM

$$\log p(\boldsymbol{x}; \boldsymbol{\theta}) = \underbrace{\mathbb{E}_{q(\boldsymbol{z})} \Big[ \log \frac{p(\boldsymbol{x})p(\boldsymbol{z}|\boldsymbol{x})}{q(\boldsymbol{z})} \Big]}_{\text{M-step - }\mathcal{Q}} + \underbrace{D_{\text{KL}} \Big(q(\boldsymbol{z}) \Big\| p(\boldsymbol{z}|\boldsymbol{x}) \Big)}_{\text{E-step - KL}}$$

Given  $\theta^0$ :

- Set  $q(\boldsymbol{z}) = p(\boldsymbol{z}|\boldsymbol{x}; \boldsymbol{\theta}^0)$ , and compute  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0)$ .
- **2** Maximize  $Q(\theta, \theta^0)$  w.r.t.  $\theta$ .

## The Exact EM

$$\log p(\boldsymbol{x}; \boldsymbol{\theta}) = \underbrace{\mathbb{E}_{q(\boldsymbol{z})} \Big[ \log \frac{p(\boldsymbol{x})p(\boldsymbol{z}|\boldsymbol{x})}{q(\boldsymbol{z})} \Big]}_{\text{M-step - }\mathcal{Q}} + \underbrace{D_{\text{KL}} \Big(q(\boldsymbol{z}) \Big\| p(\boldsymbol{z}|\boldsymbol{x}) \Big)}_{\text{E-step - KL}}$$

Given  $\theta^0$ :

- Set  $q(\boldsymbol{z}) = p(\boldsymbol{z}|\boldsymbol{x}; \boldsymbol{\theta}^0)$ , and compute  $\mathcal{Q}(\boldsymbol{\theta}, \boldsymbol{\theta}^0)$ .
- **2** Maximize  $Q(\theta, \theta^0)$  w.r.t.  $\theta$ .

Several potential issues:

- The posterior distribution does not exist/is not computationally tractable.
- The expectation cannot be taken analytically.
- The maximization w.r.t. heta does not have close-form solutions.

# **Probabilistic PCA and VAE**

Reduce the data dimensionality and extract a compact representation (z) of each datum (x).

Probabilistic version of PCA ( $D \ll F$ ):

$$p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}; \boldsymbol{0}, \boldsymbol{I}), \quad \boldsymbol{z} \in \mathbb{R}^{D}.$$
  
 $p(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}; \boldsymbol{A}\boldsymbol{z} + \boldsymbol{b}, \nu \boldsymbol{I}), \quad \boldsymbol{x} \in \mathbb{R}^{F}.$ 

Reduce the data dimensionality and extract a compact representation (z) of each datum (x).

Probabilistic version of PCA ( $D \ll F$ ):

$$p(oldsymbol{z}) = \mathcal{N}(oldsymbol{z}; oldsymbol{0}, oldsymbol{I}), \quad oldsymbol{z} \in \mathbb{R}^{D}.$$
 $p(oldsymbol{x} | oldsymbol{z}) = \mathcal{N}(oldsymbol{x}; oldsymbol{A} oldsymbol{z} + oldsymbol{b}, 
u oldsymbol{I}), \quad oldsymbol{x} \in \mathbb{R}^{F}.$ 

Posterior distribution  $\leftrightarrow$  data projection:

$$p(\boldsymbol{z}_n | \boldsymbol{x}_n; \boldsymbol{\theta}^0) = \mathcal{N}(\boldsymbol{z}_n; \boldsymbol{W} \boldsymbol{x}_n + \boldsymbol{c}, \boldsymbol{\Omega}),$$

( $oldsymbol{W}$  and  $oldsymbol{\Omega}$  depend on  $oldsymbol{ heta}^0$ )



[Image from Wikimedia Commons]

## **Comments on PPCA - Toward VAE**

- It is possible to derive an exact EM algorithm ightarrow local convergence guaranteed.
- $\bullet\,$  Can only model linear dependencies  $\rightarrow$  limited expressivity.

Non-linear Gaussian model:  $p_{\theta}(x|z) = \mathcal{N}(x; \mu_{\theta}(z), \Sigma_{\theta}(z))$ 

- $\mu_{ heta}(.), \Sigma_{ heta}(.)$ : Non-linear functions implemented as deep nets
- Parameter estimation is challenging, but much more expressive



The non-linear Gaussian model can be optimised via the variational autoencoder (VAE).

# Formalising the generative model (I)

How can we ensure that  $\Sigma_{ heta}(z)$  is a covariance matrix?

• The covariance matrix is assumed to be diagonal:

$$oldsymbol{\Sigma}_{oldsymbol{ heta}}(oldsymbol{z}) = \left(egin{array}{ccccc} 
u_{oldsymbol{ heta},1}(oldsymbol{z}) & 0 & \cdots & 0 \ 0 & 
u_{oldsymbol{ heta},2}(oldsymbol{z}) & \cdots & 0 \ dots & dots & \ddots & dots \ 0 & 0 & \cdots & 
u_{oldsymbol{ heta},F}(oldsymbol{z}) \end{array}
ight)$$

Reduces complexity and memory, but also expressivity.

• We estimate the log-variance:  $\eta_{\theta,f}(z) = \log \nu_{\theta,f}(z)$ :

$$\boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\boldsymbol{z}) = \operatorname{diag}_{f}\left(\exp\left(\eta_{\boldsymbol{\theta},f}(\boldsymbol{z})\right)\right)$$

The values of  $\eta_{\theta,f}(z)$  can be positive or negative.

# Formalising the generative model (II)

In terms of probabilistic dependencies, they are the same as PPCA:



But we can also draw the non-lineariry:

$$egin{array}{c} egin{array}{c} egin{array}$$

The dependency of the parameters w.r.t. z is deterministic.

Denoted by  $f_{\theta}(z) : \mathbb{R}^D \to \mathbb{R}^{2F}$ , this non-linearity is implemented with a deep network, with parameters (weights and biases)  $\theta$ .

Since  $p(\boldsymbol{z}|\boldsymbol{x};\boldsymbol{\theta}^0)$  cannot be computed analytically, it needs to be approximated.

The idea is to find the best candidate within a family of distributions.



The posterior distribution will be approximated with **another** feed-forward network parametrised with  $\phi$ :

$$p(\boldsymbol{z}|\boldsymbol{x}) pprox q(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}(\boldsymbol{z}; \tilde{\boldsymbol{\mu}}_{\boldsymbol{\phi}}(\boldsymbol{x}), \tilde{\boldsymbol{\Sigma}}_{\boldsymbol{\phi}}(\boldsymbol{x}))$$

The approximating family is composed of all the distributions that can be expressed as above, for a certain value of  $\phi$ .

$$\mathcal{G} = \{ \boldsymbol{g}_{\boldsymbol{\phi}} : \mathbb{R}^F o \mathbb{R}^{2D}; \boldsymbol{\phi} \in \boldsymbol{\Phi} \},$$

with  $oldsymbol{g}_{oldsymbol{\phi}}(oldsymbol{x}) = [ ilde{oldsymbol{\mu}}_{oldsymbol{\phi}}(oldsymbol{x}), ilde{m{\Sigma}}_{oldsymbol{\phi}}(oldsymbol{x})].$ 

## **Overall architecture**

If we "chain" the posterior (encoder) and the generative (decoder) model:



This is why we call these architectures variational autoencoders VAE.

But how do we optimise for the parameters  $\theta$  and  $\phi$ ?

# Learning - ELBO

If we recall the formulation for the EM:

$$\log p(\boldsymbol{x}) = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})} \Big[ \log \frac{p(\boldsymbol{x}, \boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})} \Big] + D_{\mathsf{KL}} \Big( q(\boldsymbol{z}|\boldsymbol{x}) \Big\| p(\boldsymbol{z}|\boldsymbol{x}) \Big)$$

# Learning - ELBO

If we recall the formulation for the EM:

$$\log p(\boldsymbol{x}) = \mathbb{E}_{q(\boldsymbol{z}|\boldsymbol{x})} \Big[ \log \frac{p(\boldsymbol{x}, \boldsymbol{z})}{q(\boldsymbol{z}|\boldsymbol{x})} \Big] + D_{\mathsf{KL}} \Big( q(\boldsymbol{z}|\boldsymbol{x}) \Big\| p(\boldsymbol{z}|\boldsymbol{x}) \Big)$$

Problem: the second term cannot be computed! But it's positive:

$$\log p(\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) \geq \mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log \frac{p(\boldsymbol{x}, \boldsymbol{z})}{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \right]$$
$$\log p(\boldsymbol{x}; \boldsymbol{\theta}, \boldsymbol{\phi}) \geq \underbrace{\mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \left[ \log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z}) \right]}_{\text{Reconstruction}} - \underbrace{D_{\mathsf{KL}} \left( q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) \left\| p(\boldsymbol{z}) \right)}_{\text{Regularisation}} \right]$$

This is known as **Evidence Lower-BOund or ELBO**:  $\mathcal{L}_{ELBO}(\theta, \phi)$ .

Be VERY careful with these expressions: They look alike, but they are NOT the same.

# Learning - Sampling

But we still have one problem:

$$\mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \underbrace{\mathbb{E}_{q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x})} \Big\{ \log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\boldsymbol{z}) \Big\}}_{\mathsf{Reconstruction}} - \underbrace{D_{\mathsf{KL}} \Big( q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) \Big\| p(\boldsymbol{z}) \Big)}_{\mathsf{Regularisation}}$$

To compute the "reconstruction" term we need to take the expectation w.r.t.  $q_{\phi}(z|x)$ , but recall that:

$$p_{\theta}(\boldsymbol{x}|\boldsymbol{z}) = \mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}_{\theta}(\boldsymbol{z}), \boldsymbol{\Sigma}_{\theta}(\boldsymbol{z})).$$

Due to the non-linearity, we cannot compute the reconstruction term in closed form  $\rightarrow$  we sample R points  $\hat{z}^{(1)}, \ldots, \hat{z}^{(R)}$  from  $q_{\phi}$ :

$$\mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \underbrace{\frac{1}{R} \sum_{r=1}^{R} \log p_{\boldsymbol{\theta}}(\boldsymbol{x} | \hat{\boldsymbol{z}}^{(r)})}_{\mathsf{Reconstruction}} - \underbrace{D_{\mathsf{KL}} \Big( q_{\boldsymbol{\phi}}(\boldsymbol{z} | \boldsymbol{x}) \Big\| p(\boldsymbol{z}) \Big)}_{\mathsf{Regularisation}}$$

Let's go back to the architecture:

Since there is no closed-form solution for the parameters (due to the non-linearity), we will learn the parameters using stochastic gradient ascent (to maximise the ELBO).

We assume that  $g_{\phi}$  and  $f_{\theta}$  are differentiable (we can compute the gradient).

The sampling operation  $(\hat{z}^{(r)})$  from  $q_{\phi}$  is **NOT** differentiable w.r.t.  $\phi$ .

## Learning - Reparametrisation trick

Instead of sampling directly from the posterior  $(\hat{z}^{(r)} \sim q_{\phi} = \mathcal{N}(\tilde{\mu}_{\phi}, \tilde{\Sigma}_{\phi}))$  we sample as follows:

$$ar{m{z}}^{(r)} = ilde{\Sigma}_{m{\phi}}^{1/2}ar{m{\epsilon}}^{(r)} + ilde{m{\mu}}_{m{\phi}} \hspace{0.5cm} ext{with} \hspace{0.5cm} ar{m{\epsilon}}^{(r)} \sim \mathcal{N}(m{0},m{I})$$

 $\hat{m{z}}^{(r)}$  and  $ar{m{z}}^{(r)}$  follow the same distribution, BUT  $ar{m{z}}^{(r)}$  is differentiable w.r.t.  $\phi!!!$ 

## Learning - Reparametrisation trick

Instead of sampling directly from the posterior  $(\hat{z}^{(r)} \sim q_{\phi} = \mathcal{N}(\tilde{\mu}_{\phi}, \tilde{\Sigma}_{\phi}))$  we sample as follows:

$$ar{m{z}}^{(r)} = ilde{m{\Sigma}}_{m{\phi}}^{1/2}ar{m{\epsilon}}^{(r)} + ilde{m{\mu}}_{m{\phi}} \quad ext{with} \quad ar{m{\epsilon}}^{(r)} \sim \mathcal{N}(m{0},m{I})$$

 $\hat{z}^{(r)}$  and  $ar{z}^{(r)}$  follow the same distribution, BUT  $ar{z}^{(r)}$  is differentiable w.r.t.  $\phi!!!$ 

This is called the reparametrisation trick and can be used with other distributions.



# Learning - Reparametrisation trick (II)



[Image from Wikimedia Commons]

**Implementation note**: Most deep learning libraries implement stochastic gradient DESCENT algorithms, but we would like to maximize the ELBO  $\rightarrow$  we need to change the sign.

- **Implementation note**: Most deep learning libraries implement stochastic gradient DESCENT algorithms, but we would like to maximize the ELBO  $\rightarrow$  we need to change the sign.
- **Posterior collapse**: common phenomenon when the posterior  $q_{\phi}$  gets to close to the standard prior. It can happen for various dimensions of  $z \to$  the VAE stops learning.
  - The KL term dominates the ELBO  $\rightarrow$  weight the KL term with  $\beta < 1$ .
  - The data can be reduced to less dimensions than  $D \rightarrow$  decrease D.

## VAE: Summary

Generative model. Prior:  $p(z) = \mathcal{N}(z; 0, I)$  and decoder:  $p_{\theta}(x|z) = \mathcal{N}(x; \mu_{\theta}(z), \Sigma_{\theta}(z))$ . Inference model (encoder):  $p_{\theta}(z|x) \approx q_{\phi}(z|x) = \mathcal{N}(z; \mu_{\phi}(x), \Sigma_{\phi}(x))$ 



Training criterion (maximise the evidence lower bound):

$$\mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = \underbrace{\log p_{\boldsymbol{\theta}}(\boldsymbol{x} | \hat{\boldsymbol{z}})}_{\mathsf{Reconstruction}} - \underbrace{D_{\mathsf{KL}}\Big(q_{\boldsymbol{\phi}}(\boldsymbol{z} | \boldsymbol{x}) \Big\| p(\boldsymbol{z})\Big)}_{\mathsf{Regularisation}}$$

where  $\hat{\pmb{z}} \sim q_{\pmb{\phi}}$  is sampled using the reparametrisation trick.

# VAE for audio modeling

## Representing audio: time vs. frequency

Fourier domain decomposes the signal in frequencies:



[Image from https://dev.to/trekhleb/]

Problem: we have to choose either time or frequency.

# The short-time Fourier transform (STFT)

STFT: segment the input signal, and apply DFT to each segment.



[Image from Mathworks]

- Sine frequency sweep (pure sine of increasing frequency).
- Leyenda (piano, harmonics).
- O mio babbino caro (opera, vibrato).
- Highway to hell (rock&roll, distortion).

Each observation x will be an F-dimensional **complex** vector:  $x \in \mathbb{C}^{F}$ . The low-dimensional latent variable will be **real** of dimension D:  $z \in \mathbb{R}^{D}$ .

The model:

- Prior:  $p(\boldsymbol{z}) = \mathcal{N}(\boldsymbol{z}; \boldsymbol{0}, \boldsymbol{I}).$
- Decoder:  $p_{\theta}(x|z) = \mathcal{N}_c(x; 0, \Sigma_{\theta}(z))$ , complex Gaussian distribution, see next.
- Posterior:  $p_{\theta}(\boldsymbol{z}|\boldsymbol{x}) \approx q_{\phi}(\boldsymbol{z}|\boldsymbol{x}) = \mathcal{N}\Big(\boldsymbol{z}; \boldsymbol{\mu}_{\phi}(\boldsymbol{x}), \boldsymbol{\Sigma}_{\phi}(\boldsymbol{x})\Big)$

We need to compute the reconstruction and the regularization terms.

Let's recall that the covariance matrices  $\Sigma_{\theta}(z)$  and  $\Sigma_{\phi}(x)$  are diagonal log-variance matrices.

## ELBO for spectrogram VAEs

The centered 1D complex-normal writes:

$$\mathcal{N}_c(x; 0, \nu) = rac{1}{\pi 
u} \exp\left(-rac{|x|^2}{
u}
ight).$$

The 1D complex-normal for the f-th entry of x as a function of the log-var writes:

$$\log \mathcal{N}_c(x_f; 0, \exp \eta_{\boldsymbol{\theta}, f}(\boldsymbol{z})) = -\log(\pi) - \eta_{\boldsymbol{\theta}, f}(\boldsymbol{z}) - \frac{|x_f|^2}{\exp(\eta_{\boldsymbol{\theta}, f}(\boldsymbol{z}))}.$$

For a given sample  $\hat{z}$ , the reconstruction term writes:

$$\log p_{\boldsymbol{\theta}}(\boldsymbol{x}|\hat{\boldsymbol{z}}) = \sum_{f} \log p_{\boldsymbol{\theta}}(x_{f}|\hat{\boldsymbol{z}}) = -F\log(\pi) - \sum_{f} \left( \eta_{\boldsymbol{\theta},f}(\hat{\boldsymbol{z}}) + \frac{|x_{f}|^{2}}{\exp(\eta_{\boldsymbol{\theta},f}(\hat{\boldsymbol{z}}))} \right).$$

## ELBO for spectrogram VAEs (II)

The KL between two *D*-dimensional real Gaussian distributions writes:

$$D_{\mathrm{KL}}(\mathcal{N}_0 \| \mathcal{N}_1) = \frac{1}{2} \left( \mathrm{Tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) - D + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \log \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|} \right)$$

In our case:  $\mu_0 = \tilde{\mu}_{\phi}(x)$ ,  $\Sigma_0 = \operatorname{diag}_d(\exp{(\tilde{\eta}_{\phi,d}(x))})$ ,  $\mu_1 = 0$  and  $\Sigma_1 = I$ . Thus:

$$D_{\mathrm{KL}}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) \| p(\boldsymbol{z})) = \frac{1}{2} \left( \sum_{d} \exp\left(\tilde{\eta}_{\boldsymbol{\phi},d}(\boldsymbol{x})\right) + |\tilde{\mu}_{\boldsymbol{\phi},d}(\boldsymbol{x})|^2 - \tilde{\eta}_{\boldsymbol{\phi},d}(\boldsymbol{x}) \right) - \frac{D}{2}$$

## ELBO for spectrogram VAEs (II)

The KL between two *D*-dimensional real Gaussian distributions writes:

$$D_{\mathrm{KL}}(\mathcal{N}_0 \| \mathcal{N}_1) = \frac{1}{2} \left( \mathrm{Tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) - D + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \log \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|} \right)$$

In our case:  $\mu_0 = \tilde{\mu}_{\phi}(x)$ ,  $\Sigma_0 = \operatorname{diag}_d(\exp{(\tilde{\eta}_{\phi,d}(x))})$ ,  $\mu_1 = 0$  and  $\Sigma_1 = I$ . Thus:

$$D_{\mathrm{KL}}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) \| p(\boldsymbol{z})) = \frac{1}{2} \left( \sum_{d} \exp\left(\tilde{\eta}_{\boldsymbol{\phi},d}(\boldsymbol{x})\right) + |\tilde{\mu}_{\boldsymbol{\phi},d}(\boldsymbol{x})|^2 - \tilde{\eta}_{\boldsymbol{\phi},d}(\boldsymbol{x}) \right) - \frac{D}{2}$$

Therefore the ELBO (w/o constant terms) writes:

$$\mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = -\sum_{f} \left( \eta_{\boldsymbol{\theta}, f}(\hat{\boldsymbol{z}}) + \frac{|\boldsymbol{x}_{f}|^{2}}{\exp(\eta_{\boldsymbol{\theta}, f}(\hat{\boldsymbol{z}}))} \right) - \frac{1}{2} \left( \sum_{d} \exp\left(\tilde{\eta}_{\boldsymbol{\phi}, d}(\boldsymbol{x})\right) + |\tilde{\mu}_{\boldsymbol{\phi}, d}(\boldsymbol{x})|^{2} - \tilde{\eta}_{\boldsymbol{\phi}, d}(\boldsymbol{x}) \right)$$

which is what we give to our optimizer.

## ELBO for spectrogram VAEs (II)

The KL between two *D*-dimensional real Gaussian distributions writes:

$$D_{\mathrm{KL}}(\mathcal{N}_0 \| \mathcal{N}_1) = \frac{1}{2} \left( \mathrm{Tr}(\boldsymbol{\Sigma}_1^{-1} \boldsymbol{\Sigma}_0) - D + (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0)^\top \boldsymbol{\Sigma}_1^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_0) + \log \frac{|\boldsymbol{\Sigma}_1|}{|\boldsymbol{\Sigma}_0|} \right)$$

In our case:  $\mu_0 = \tilde{\mu}_{\phi}(\boldsymbol{x})$ ,  $\Sigma_0 = \operatorname{diag}_d \left( \exp\left( \tilde{\eta}_{\phi,d}(\boldsymbol{x}) \right) \right)$ ,  $\mu_1 = \mathbf{0}$  and  $\Sigma_1 = \boldsymbol{I}$ . Thus:

$$D_{\text{KL}}(q_{\boldsymbol{\phi}}(\boldsymbol{z}|\boldsymbol{x}) \| p(\boldsymbol{z})) = \frac{1}{2} \left( \sum_{d} \exp\left(\tilde{\eta}_{\boldsymbol{\phi},d}(\boldsymbol{x})\right) + |\tilde{\mu}_{\boldsymbol{\phi},d}(\boldsymbol{x})|^2 - \tilde{\eta}_{\boldsymbol{\phi},d}(\boldsymbol{x}) \right) - \frac{D}{2}$$

Therefore the ELBO (w/o constant terms) writes:

$$\mathcal{L}_{\mathsf{ELBO}}(\boldsymbol{\theta}, \boldsymbol{\phi}) = -\sum_{f} \left( \eta_{\boldsymbol{\theta}, f}(\hat{\boldsymbol{z}}) + \frac{|\boldsymbol{x}_{f}|^{2}}{\exp(\eta_{\boldsymbol{\theta}, f}(\hat{\boldsymbol{z}}))} \right) - \frac{1}{2} \left( \sum_{d} \exp\left(\tilde{\eta}_{\boldsymbol{\phi}, d}(\boldsymbol{x})\right) + |\tilde{\mu}_{\boldsymbol{\phi}, d}(\boldsymbol{x})|^{2} - \tilde{\eta}_{\boldsymbol{\phi}, d}(\boldsymbol{x}) \right)$$

which is what we give to our optimizer. NO!!! We should give  $-\mathcal{L}_{\text{ELBO}}(\theta, \phi)$ .

Thank you for your attention

- Lecture 2: The variational EM algoritm, mixtures of VAEs, application to audio-visual speech enhancement.
- Lecture 3: Dynamical variational autoencoders, application to speech enhancement.

#### References

- **1** D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," ICLR, 2014.
- Y. Bando et al., "Statistical speech enhancement based on probabilistic integration of variational autoencoder and non-negative matrix factorization," in Proc. ICASSP, 2018, pp. 716–720
- S. Leglaive et al., "A variance modeling framework based on variational autoencoders for speech enhancement," in Proc. MLSP, 2018.
- C. Bishop, Pattern Recognition and Machine Learning, Springer-Verlag Berlin, Heidelberg, 2006.