

Latent Variable Models

M2 IASD - Y. Chevaleyre

Outline

1. Introduction
2. Gaussian mixtures and the EM algorithm
3. Variational Analysis of the EM algorithm
4. Variational Auto-Encoders
5. Exercise: probabilistic PCA

Note:

latent var appear everywhere, but in this lecture, we will focus on generative problems (\neq supervised learning)

A generative problem is learning a distribution for data $(x_1, \dots, x_N) \in \mathbb{R}^d$ in order to generate new points from this distribution

Outline

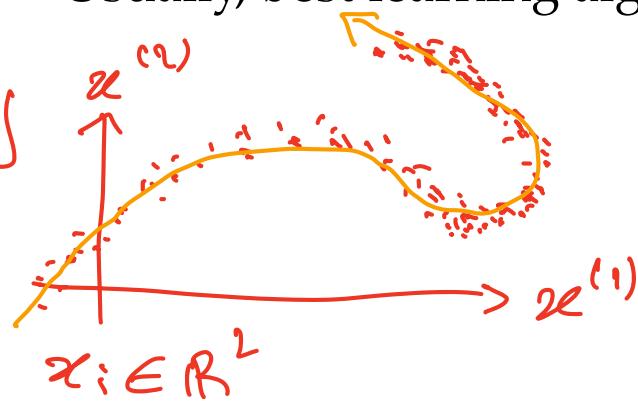
1. **Introduction**
2. Gaussian mixtures and the EM algorithm
3. Variational Analysis of the EM algorithm
4. Variational Auto-Encoders
5. Exercise: probabilistic PCA

Introduction: what are latent variable models ?

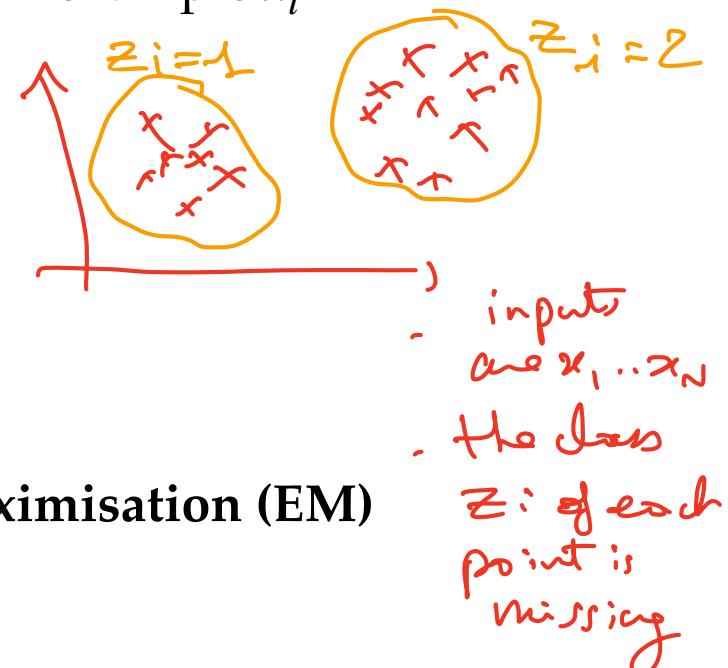
- A latent variable z_i is a missing information about an example x_i

- example:

- In the **clustering** setting
- in **Dimensionality reduction** setting
- in the case of **missing data**
- Usually, best learning algorithm is **Expectation Maximisation (EM)**

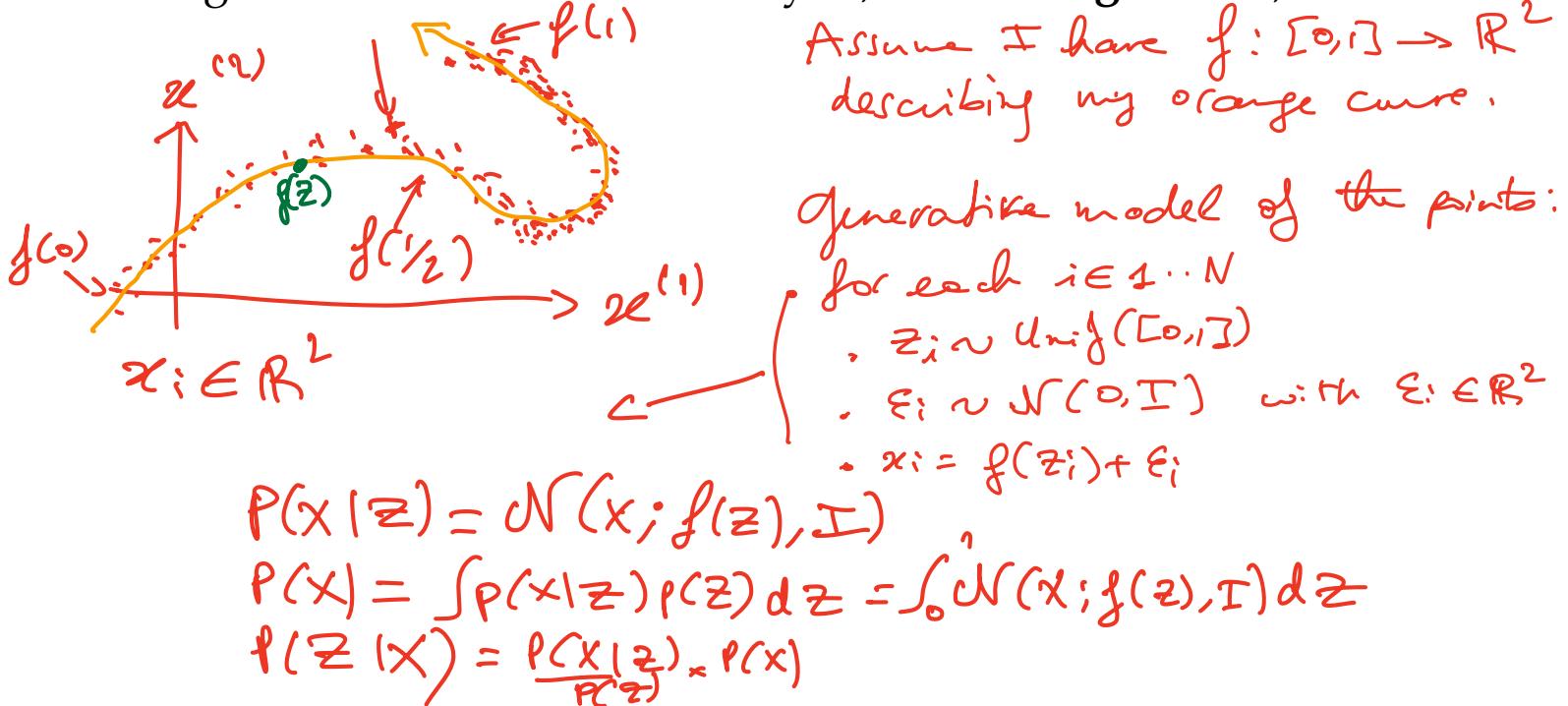


for each x_i , the $z_i \in \mathbb{R}$
is its coordinate on the orange curve



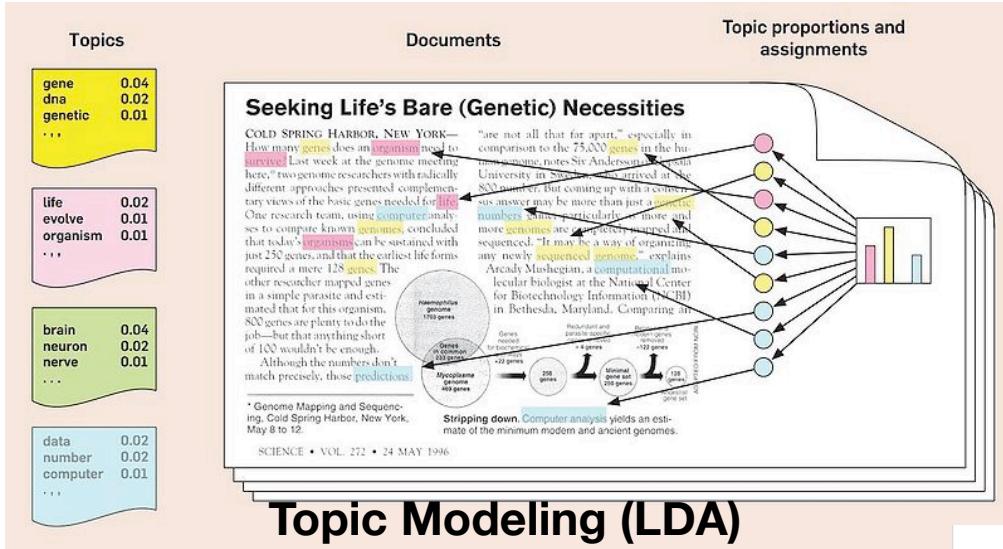
Probabilistic Machine Learning

- The most appropriate framework for handling latent variables is **Probabilistic ML**.
- The Probabilistic Machine Learning is when you assume that your data (x_i or y_i or both) was generated according to some generative model with some unknown parameters
- Learning a model is done by **MLE** or **MAP**
- e.g: Linear Discriminant Analysis, **Linear Regression**, ...

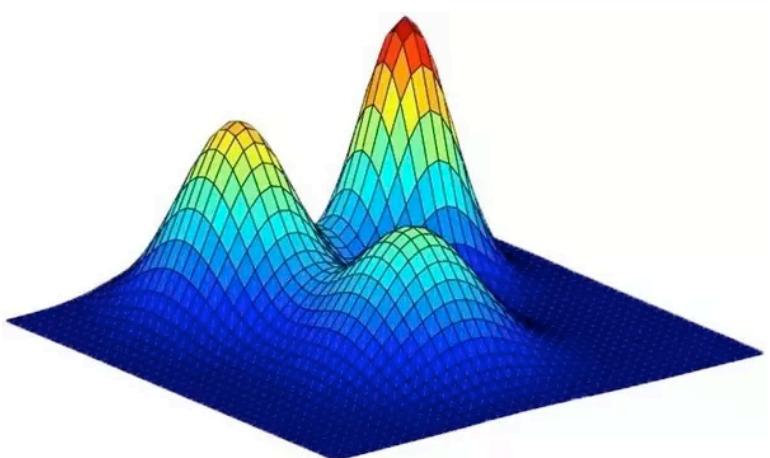


Probabilistic Machine Learning

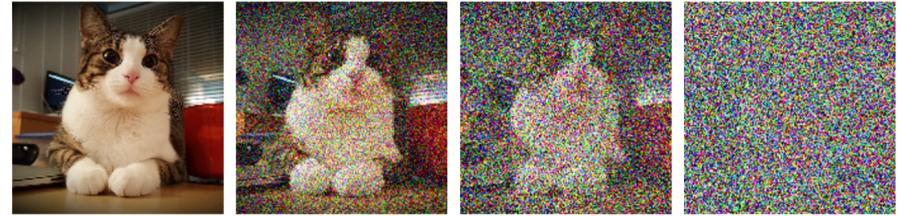
- This generative model may include unobserved variables. We call these *Latent Variables*.



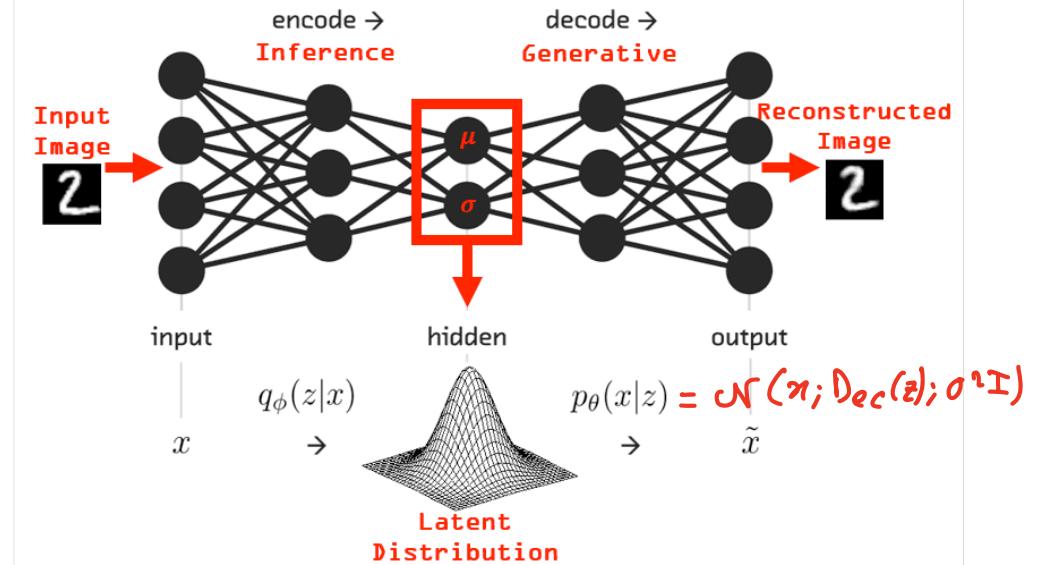
Topic Modeling (LDA)



Gaussian Mixture Models (GMM)



Latent Diffusion Processes (e.g. DALLE-2)



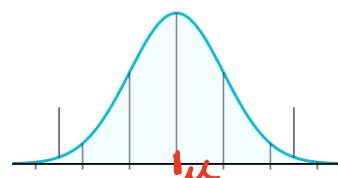
Variational AutoEncoders (VAE)

Outline

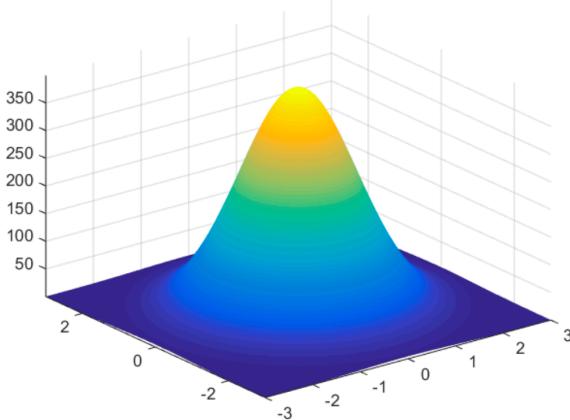
1. Introduction
2. **Gaussian mixtures and the EM algorithm**
3. Variational Analysis of the EM algorithm
4. Variational Auto-Encoders
5. Exercise: probabilistic PCA

Reminder: The multivariate Normal Distribution

- Gaussienne univariée



- Gaussienne multivariée:



A normal (or Gaussian) distribution in 2 variables

- Notation:

- $\mathcal{N}(\mu, \sigma^2)$ and $\mathcal{N}(\mu, \Sigma)$ represent the normal distribution.
- We write $X \sim \mathcal{N}(\mu, \sigma^2)$ or $X \sim \mathcal{N}(\mu, \Sigma)$
- $\mathcal{N}(x | \mu, \Sigma)$ represents the p.d.f. of $\mathcal{N}(\mu, \Sigma)$ at point x

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

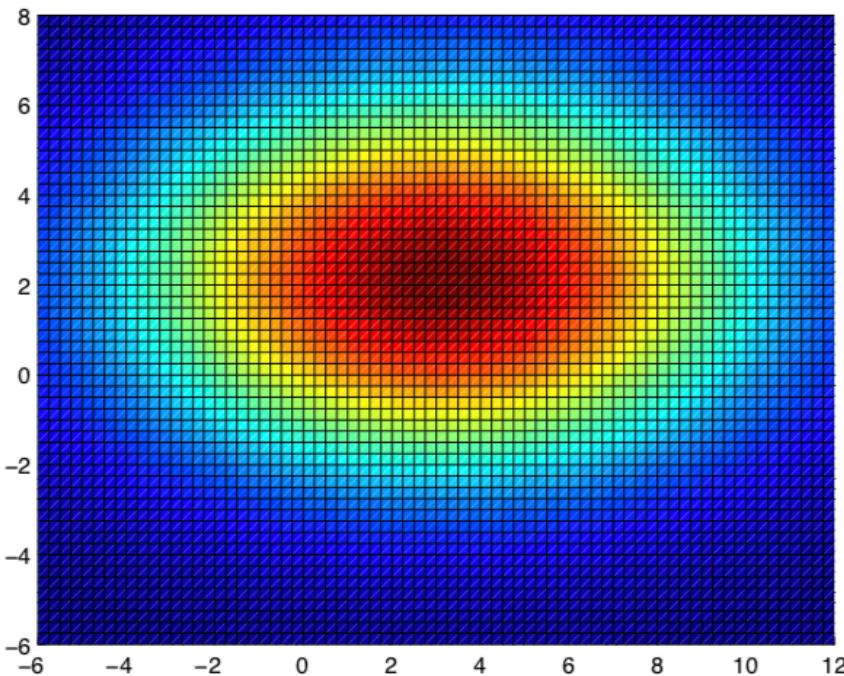
 $\mu \in \mathbb{R}^d, \Sigma \in \mathbb{R}^{d \times d}, x \in \mathbb{R}^d$

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

if $\Sigma = \mathbb{I}_d$
 then $\mathcal{N}(x; \mu, \mathbb{I}) = \frac{1}{(2\pi)^{d/2}} \exp\left\{-\frac{1}{2} \|x-\mu\|^2\right\}$

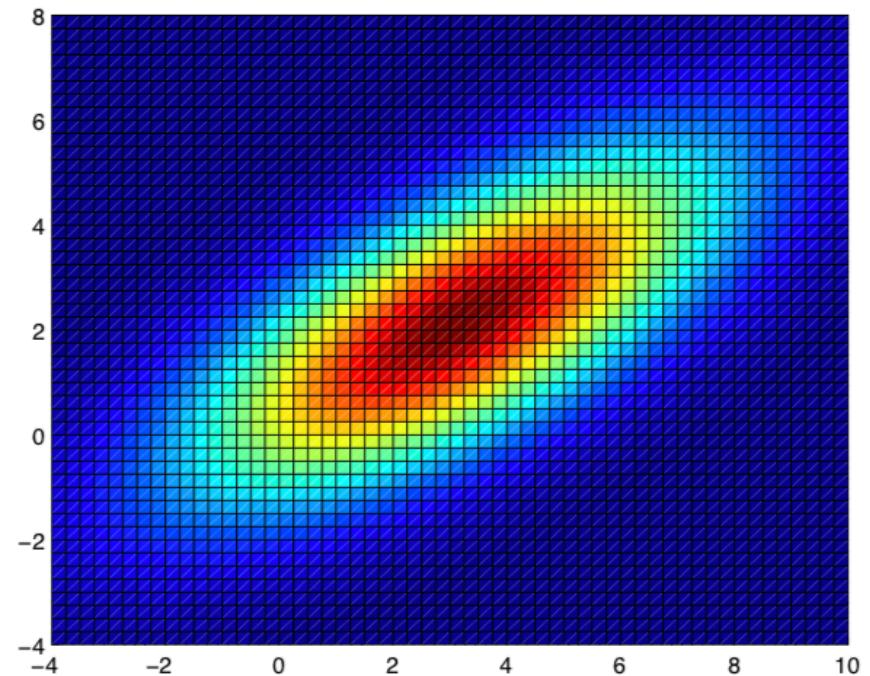
Reminder: The multivariate Normal Distribution

- Gaussienne multivariée:



$$\mu = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 25 & 0 \\ 0 & 9 \end{bmatrix}$$

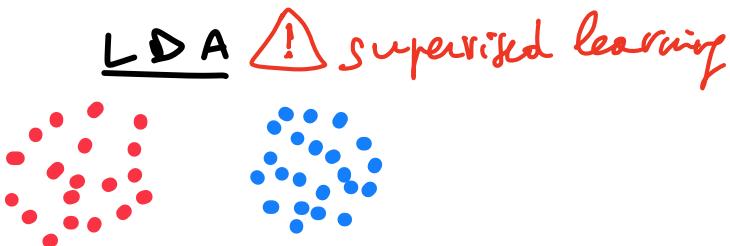


$$\mu = \begin{bmatrix} 3 \\ 2 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 10 & 5 \\ 5 & 5 \end{bmatrix}$$

Gaussian Mixture Models

- One of the simplest latent variable model is *Gaussian Mixture Models (GMM)*
- same as LDA without class information. Recall LDA with 2 classes first:
(linear discriminant analysis)



generative model (with 2 classes)

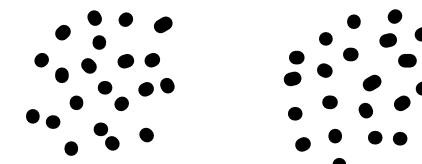
- . two gaussians $\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)$
- . a parameter $\pi \in [0, 1]$

for each $i = 1..N$

- $y_i \sim \text{Ber}(\pi) + 1$
- if $y_i = 1$ then $x_i \sim \mathcal{N}(\mu_1, \Sigma_1)$
if $y_i = 2$, $x_i \sim \mathcal{N}(\mu_2, \Sigma_2)$

Note: this define $P(X, Y)$, $P(X)$, $P(Y)$.

GMM Unsupervised



- . generative model (with 2 classes)
- . two gaussians $\mathcal{N}(\mu_1, \Sigma_1), \mathcal{N}(\mu_2, \Sigma_2)$
- . a parameter $\pi \in [0, 1]$
- . for each $i = 1..N$
 - $z_i \sim \text{Ber}(\pi) + 1$
 - if $z_i = 1$ then $x_i \sim \mathcal{N}(\mu_1, \Sigma_1)$
if $z_i = 2$, $x_i \sim \mathcal{N}(\mu_2, \Sigma_2)$
- . Note: this define $P(X, Z)$, $P(X)$, $P(Z)$..
- . At the end, only $x_1..x_N$ are observed. $z_1..z_N$ are unknown.

Gaussian Mixture Models: The log-likelihood

Exercise:

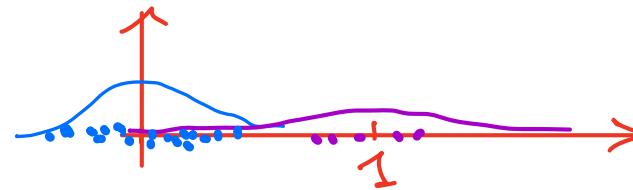
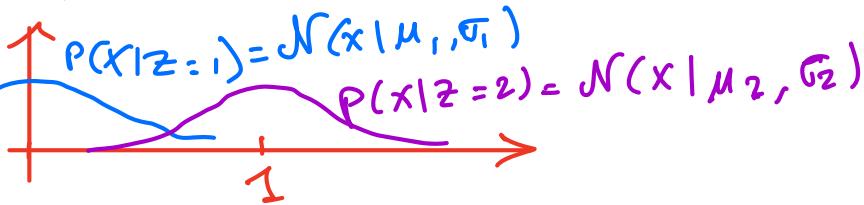
$K=2$ gaussians, $d=1$ $(x_1 \dots x_n) \in \mathbb{R}$

$$\mu_1 = 0 \\ \sigma_1 = 1$$

$$\mu_2 = 1 \\ \sigma_2 = 1$$

$$\pi = P(Z_i=1) = 80\% \\ P(Z_i=2) = 1 - \pi$$

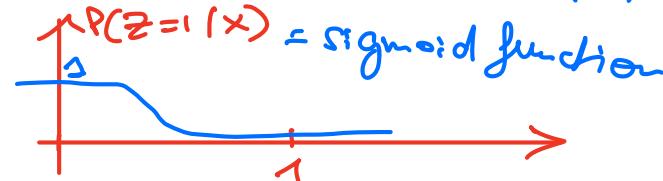
- 1) Draw informally $P(X|Z=1)$, $P(X|Z=2)$ on a graph } these are all functions of X
- 2) Draw $P(X, Z=1)$ and $P(X, Z=2)$ on another graph }



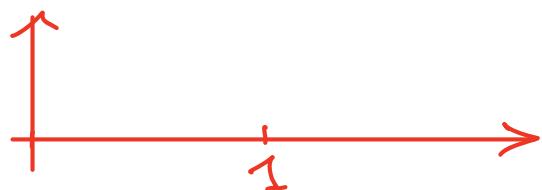
- 3) Draw $P(Z=1|X)$

$$P(Z=1|X) = \frac{P(X, Z=1)}{P(X)} = \frac{P(X|Z=1) P(Z=1)}{\pi P(X|Z=1) + (1-\pi) P(X|Z=2)}$$

$$P(X, Z=1) = P(X|Z=1) \times P(Z=1) = N(x|\mu_1, \sigma_1^2) \times \pi \\ P(X, Z=2) = \dots \dots \dots = N(x|\mu_2, \sigma_2^2) \times (1-\pi)$$



- 4) Draw points "by hand" from this process.



- 5) Exercise:

compute and simplify $P(Z=1|X)$

exercice:

$K=2$ gaussiens, $d=1$ $(x_1 \dots x_n) \in \mathbb{R}$

$$\mu_1 = 0$$

$$\mu_2 = 1$$

$$\sigma_1^2 = 1$$

$$\sigma_2^2 = 1$$

$$\pi = P(Z_i=1) = 80\%$$

$$P(Z_i=2) = 1 - \pi$$

$$P(X|z=c) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}(x-\mu_c)^2}$$

$$P(Z=1|x) = \frac{P(X|Z=1)\pi}{\pi \cdot P(X|Z=1) + (1-\pi) \cdot P(X|Z=2)}$$

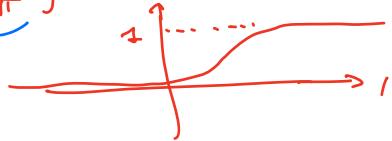
$$= \frac{\pi \exp(-\frac{1}{2}(x-\mu_1)^2)}{\pi \exp(-\frac{1}{2}(x-\mu_1)^2) + (1-\pi) \exp(-\frac{1}{2}(x-\mu_2)^2)}$$

multiply everything
by $\exp(\frac{1}{2}(x-\mu_1)^2) \times \frac{1}{\pi}$

$$= \frac{1}{1 + \frac{1-\pi}{\pi} \exp\left\{-\frac{1}{2}(x-\mu_2)^2 + \frac{1}{2}(x-\mu_1)^2\right\}}$$

$$= \frac{1}{1 + \exp\left\{x \underbrace{x(\mu_1 - \mu_2)}_{a} + \underbrace{\frac{\mu_2^2 - \mu_1^2}{2} + \ln \frac{1-\pi}{\pi}}_{b}\right\}}$$

$$= \frac{1}{1 + \exp(ax+b)}$$



Compute the log likelihood of a GAN model

- K gaussians $\mathcal{N}(\mu_1, \Sigma_1) \dots \mathcal{N}(\mu_K, \Sigma_K)$, with priors $p(z_i=k) = \pi_k$ for $k \in 1 \dots K$

- We have a dataset $x_1 \dots x_n \in \mathbb{R}^N$

- $\Theta = (\mu_1, \Sigma_1, \dots, \mu_K, \Sigma_K, \pi_1, \dots, \pi_K)$

- the log likelihood is

$$\begin{aligned} L(\Theta) &= \log P_\Theta(x_1, \dots, x_N) = \sum_{i=1}^n \log P_\Theta(x_i) \\ &= \sum_{i=1}^n \log \sum_{k=1}^K p_\Theta(x_i | z_i=k) p_\Theta(z_i=k) \quad (\text{total law of proba}) \\ &= \sum_{i=1}^n \log \sum_{k=1}^K \mathcal{N}(x_i | \mu_k, \Sigma_k) \cdot \pi_k \end{aligned}$$

- The max likelihood estimate is $\hat{\Theta} = \underset{\Theta}{\operatorname{argmax}} L(\Theta)$

Note: for $K=1$, computing $\hat{\Theta}$ is easy because we know how to compute it.

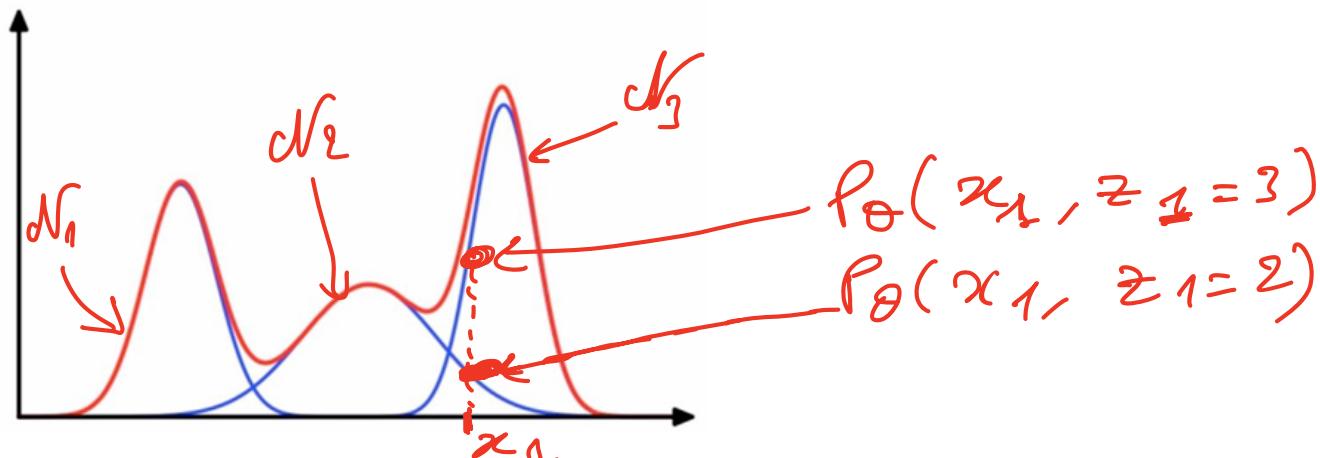
for $K>1$: non convex, often "hard" to compute

We could learn Θ by gradient descent, but much better strategy

GMM with K>1. The Complete Log Likelihood

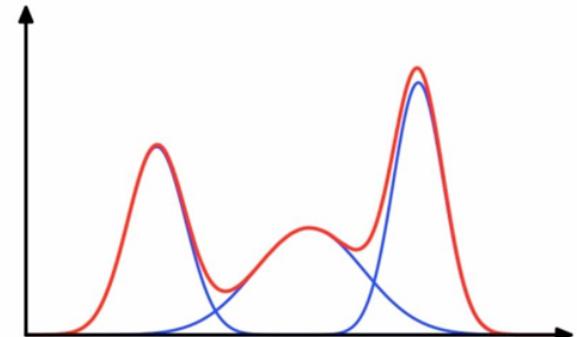
- For $K > 1$, maximising the likelihood is hard
- Idea: optimize a *surrogate* instead: the **Complete Log Likelihood (CLL)**
- The CLL assumes we know $z_1 \dots z_N$ so we still can't compute it either

$$CLL(\theta, z_1 \dots z_N) = \log P_\theta(x_1 \dots x_N, z_1 \dots z_N)$$



GMM with K>1. The Expected Complete Log Likelihood

- Idea: optimize a *surrogate* instead: the **Complete Log Likelihood (CLL)**
- The CLL assumes we know $z_1 \dots z_N$ so we still can't compute it either
- But if for each i we can have an *estimate* the distribution $p(z_i = k | x_i; \hat{\theta})$ then we can compute the **Expected CLL**. Let $q_i(k)$ be this estimation.



$$q_i(k) \approx p_{\theta}(z_i = k | x_i)$$

$$\begin{aligned} \mathcal{L}(q_i, \theta, i) &= \mathbb{E}_{z_i \sim q_i} \log P_{\theta}(x_i, z_i) = \mathbb{E}_{z_i \sim q_i} \log (P_{\theta}(x_i | z_i) \times \pi_i) \\ &= \sum_{k=1}^K q_i(k) \log (P_{\theta}(x_i | z_i = k) \cdot \pi_k) \end{aligned}$$

$$\begin{aligned} ECLL = \mathcal{L}(q, \theta) &= \sum_{i=1}^N \mathcal{L}(q_i, \theta, i) = \mathbb{E}_{z_1 \sim q_1, \dots, z_N \sim q_N} [\text{CLL}(\theta, z_1, z_N)] \\ &= \sum_{i=1}^N \sum_{k=1}^K q_i(k) \log (P_{\theta}(x_i | z_i = k) \cdot \pi_k) \end{aligned}$$

⚠ $\mathcal{L}(q, \theta)$ is convex in $\mu_1 \dots \mu_K$!

The EM Algorithm

- EM Algorithm:

1. Initialize $\hat{\theta}$ arbitrarily

2. E-Step: for each $i \in \{1 \dots N\}$, $k \in \{1 \dots K\}$

$$\text{compute } q_i(k) = p(z_i = k | x_i; \hat{\theta}) = p_{\hat{\theta}}(z_i = k | x_i) \neq p(z_i = k | x_i)$$

3. M-Step: compute $\hat{\theta} \leftarrow \arg \max_{\theta} \mathcal{L}(q, \theta)$

↑ (expected complete d.).
the opt. is convex.

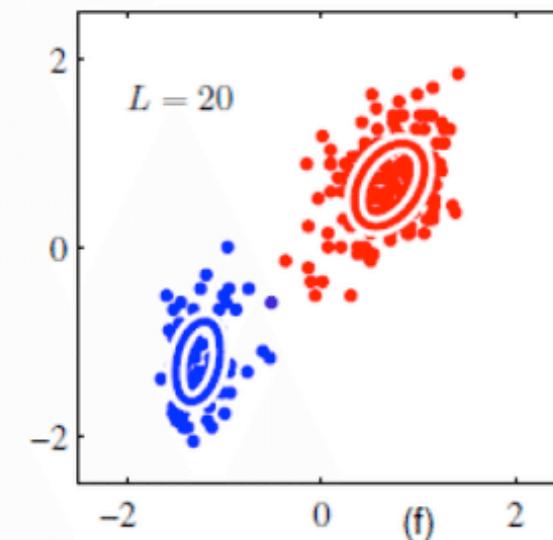
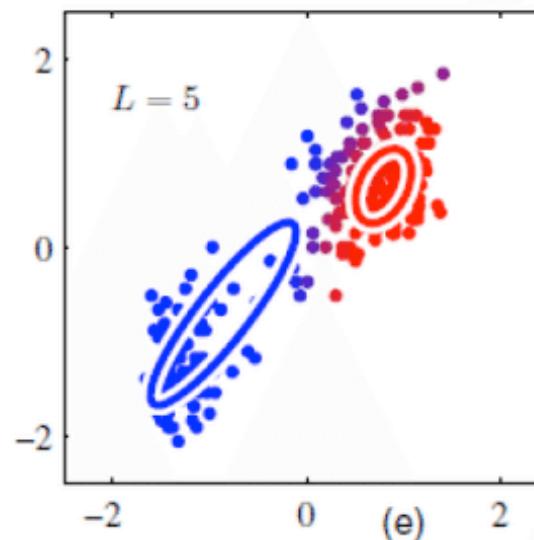
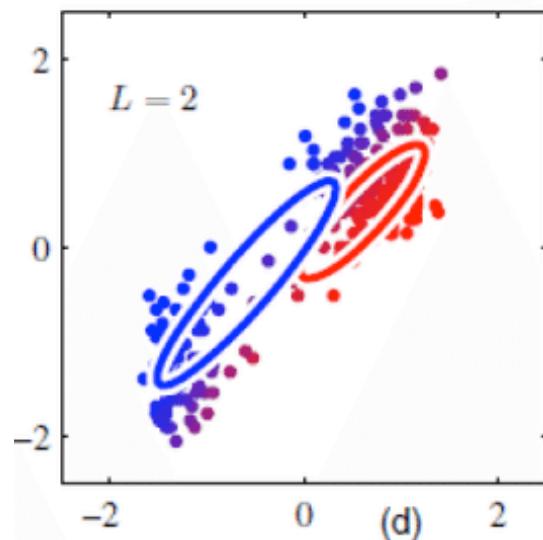
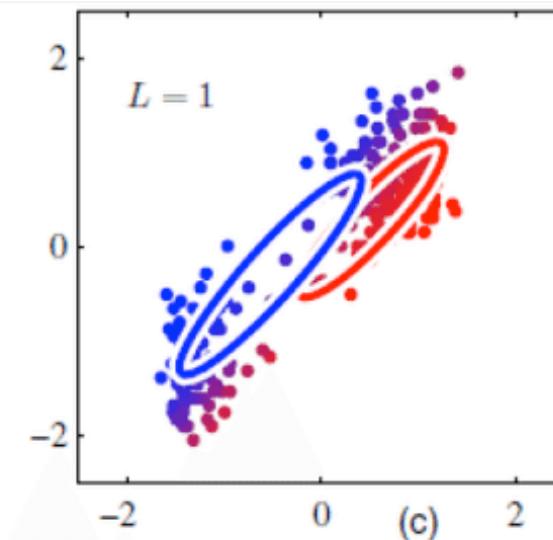
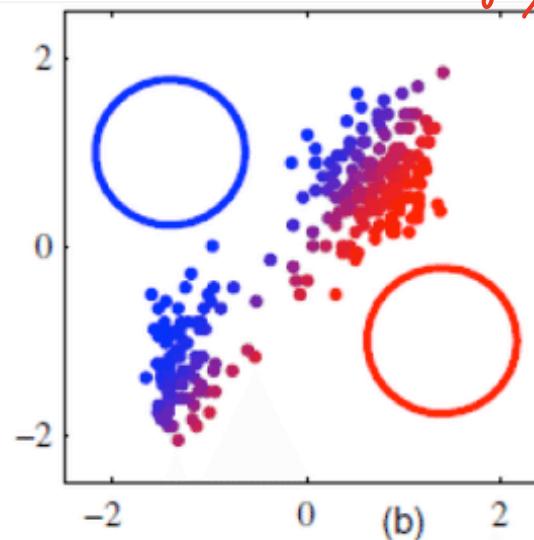
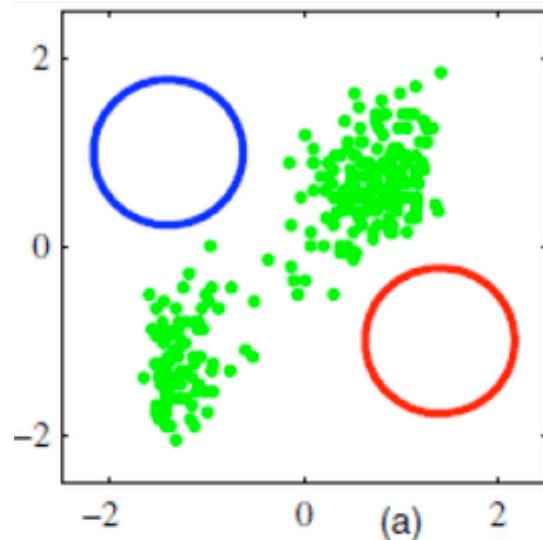
- Details: (in 1D):

$$\arg \max_{\theta} \mathcal{L}(q, \theta) = \arg \max_{\theta} \sum_{i=1}^N \sum_{k=1}^K q_i(k) \left(\text{cste} - \ln \sigma_k - \frac{(x_i - \mu_k)^2}{2\sigma_k^2} \right)$$

For fixed $\sigma_k = r$, I get $\hat{\mu}_k = \left(\frac{1}{\sum_{i=1}^N q_i(k)} \right) \times \sum_{i=1}^N x_i \times q_i(k)$

The EM Algorithm

E-step: estimating



Outline

1. Introduction
2. Gaussian mixtures and the EM algorithm
3. **Variational Analysis of the EM algorithm**
4. Variational Auto-Encoders
5. Exercise: probabilistic PCA

Variational Analysis of EM: The ELBO (1)

- We really want to optimise the likelihood $p(x_1 \dots x_N | \theta)$, but instead we optimise the a surrogate: the expected CLL $\mathcal{L}(q, \theta)$. What is the link ?

- Recall that $\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathcal{L}(q, \theta, i)$ where $\mathcal{L}(q, \theta, i) = \mathbb{E}_{z_i \sim q_i} [p_{\theta}(X = x_i, z_i = k | \theta)]$

- Let us pick a single point x and compute its ~~log~~ likelihood

$$\log p_{\theta}(x) = \ln \sum_{k=1}^K p_{\theta}(x, z=k) = \ln \sum_{k=1}^K q_x(k) \times \frac{p_{\theta}(x, z=k)}{q_x(k)}$$

$$q_x(k) \approx p_{\theta}(z=k | x=x)$$

Jensen's inequality: if f is concave then $\mathbb{E}(f(x)) \leq f(\mathbb{E}(x))$

$$\log p_{\theta}(x) = \ln \mathbb{E}_{k \sim q_x} \left[\frac{p_{\theta}(x, z=k)}{q_x(k)} \right] \geq \mathbb{E}_{k \sim q_x} \ln \left[\frac{p_{\theta}(x, z=k)}{q_x(k)} \right]$$

$$\log p_{\theta}(x) \geq \sum_{k=1}^K q_x(k) \ln p_{\theta}(x, z=k) - \sum_{k=1}^K q_x(k) \ln q_x(k)$$

$\mathcal{L}(q_x, \theta, x) = \text{ECLL}$
for a single example

Entropy

ELBO(x, θ, q_x) Evidence Lower Bound

Variational Analysis of EM: The ELBO (2)

- We really want to optimise the likelihood $p(x_1 \dots x_N | \theta)$, but instead we optimise the a surrogate: the expected CLL $\mathcal{L}(q, \theta)$. What is the link?

- Recall that $\mathcal{L}(q, \theta) = \sum_{i=1}^N \mathcal{L}(q, \theta, i)$ where $\mathcal{L}(q, \theta, i) = \mathbb{E}_{z_i \sim q_i} [p(X = x_i, z_i = k | \theta)]$

How close the log likelihood and the ELBO are?

$$\begin{aligned} \ln p_\theta(x) - \text{ELBO}(x, \theta, q_x) &= \ln p_\theta(x) - \mathbb{E}_{k \sim q_x} \left[\ln \left(\frac{p_\theta(x, z=k)}{q_x(k)} \right) \right] \\ &= \mathbb{E}_{k \sim q_x} \left[\ln \frac{p_\theta(x) q_x(k)}{p_\theta(x, z=k)} \right] = \mathbb{E} \left[\ln \frac{q_x(k)}{p_\theta(z=k|x)} \right] \\ &= \text{KL}\left(q_x(z) \parallel p_\theta(z|x)\right) \end{aligned}$$

Properties of the Kullback Leibler divergence:

- $\text{KL}(q || q') = \mathbb{E}_{z \sim q} \left[\ln \left(\frac{q(z)}{q'(z)} \right) \right]$
- $\text{KL}(q || q') \geq 0 \quad \forall q, q'$
- $\text{KL}(q || q') = 0 \text{ iff } q = q'$
- $\text{KL}(q || q') \neq \text{KL}(q' || q)$

- $\ln p_\theta(x) = \text{ELBO}(x, \theta, q_x) \text{ iff } q_x(z) = p_\theta(z|x)$
- $\underset{q_x(\cdot)}{\text{argmax}} \text{ELBO}(x, \theta, q_x) = p_\theta(z|x)$
- $\underset{\theta}{\text{argmax}} \text{ELBO}(x, \theta, p_\theta(z|x)) = \underset{\theta}{\text{argmax}} \ln p_\theta(x)$
- $\underset{\theta, q_x}{\text{argmax}} \text{ELBO}(x, \theta, q) = \underset{\theta}{\text{argmax}} \ln p_\theta(x)$

Variational Analysis of EM: reinterpreting the algorithm

- EM Algorithm:

1. Initialize $\hat{\theta}$ arbitrarily

2. E-Step: for each $i \in \{1 \dots N\}, k \in \{1 \dots K\}$

$$\text{compute } q_i(k) = p(z_i = k | x_i; \hat{\theta})$$

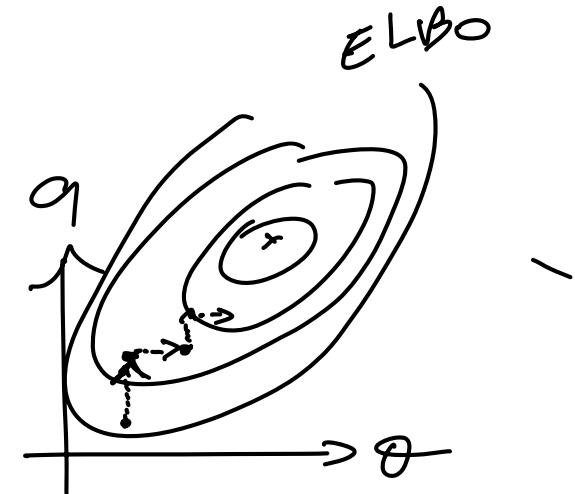
$$q_1 \dots q_N = \arg \max_{q_1 \dots q_N} \sum_i \text{ELBO}(x_i, \hat{\theta}, q_i)$$

$$\leq p_\theta(z_1 | x_1) \dots p_\theta(z_N | x_N)$$

3. M-Step: compute $\hat{\theta} \leftarrow \arg \max_{\theta} \mathcal{L}(q, \theta)$

$$\hat{\theta} = \arg \max_{\theta} \sum_i \text{ELBO}(x_i, \theta, q_i)$$

- Because $\text{ELBO}(x_i, \theta, q_i) = \ln p_\theta(x_i) - \text{KL}(q_i(z) \| p_\theta(z | x_i))$, the E-step minimizes the KL and the M-step maximizes the log likelihood while possibly increasing the KL.
- Each step (E and M) increases the ELBO. Unless we get stuck in a local minimum, this will reach the optimal ELBO which coincides with the optimal log likelihood.

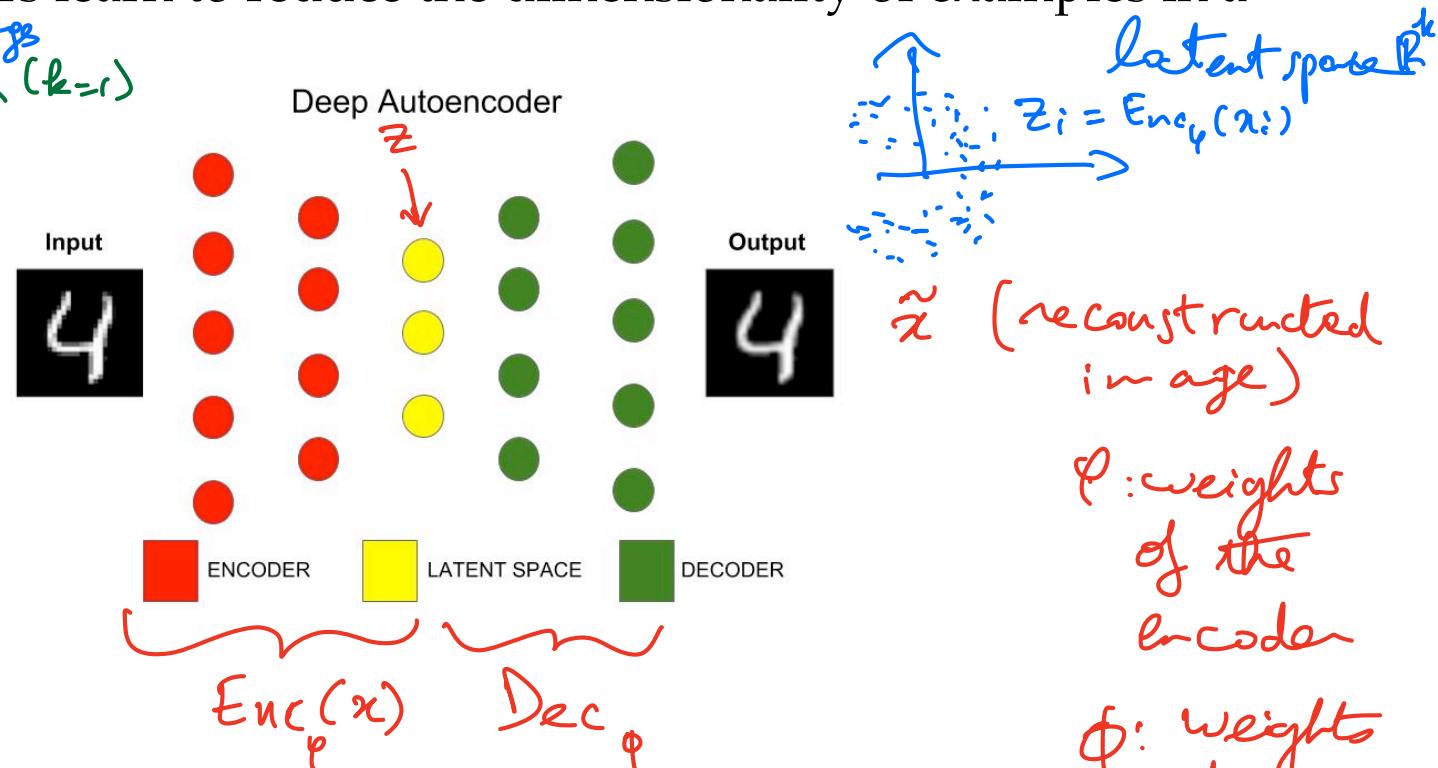
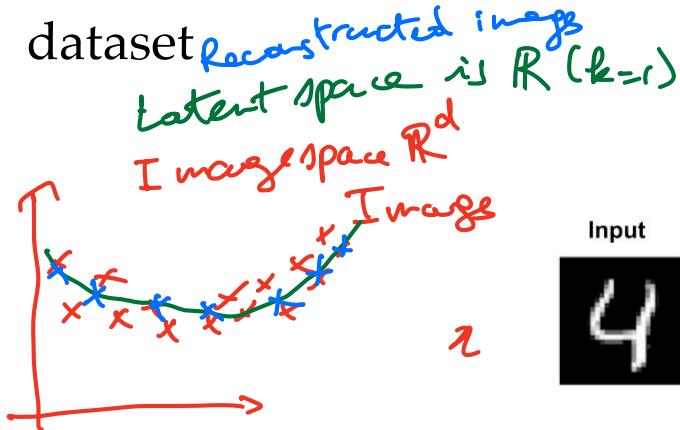


Outline

1. Introduction
2. Gaussian mixtures and the EM algorithm
3. Variational Analysis of the EM algorithm
- 4. Variational Auto-Encoders**
5. Exercise: probabilistic PCA

A word about auto-encoders

- Standard auto-encoders learn to reduce the dimensionality of examples in a dataset



To train this:

$$\underset{\varphi, \phi}{\operatorname{argmin}} \sum_{i=1}^N \|x_i - \tilde{x}_i\|^2$$

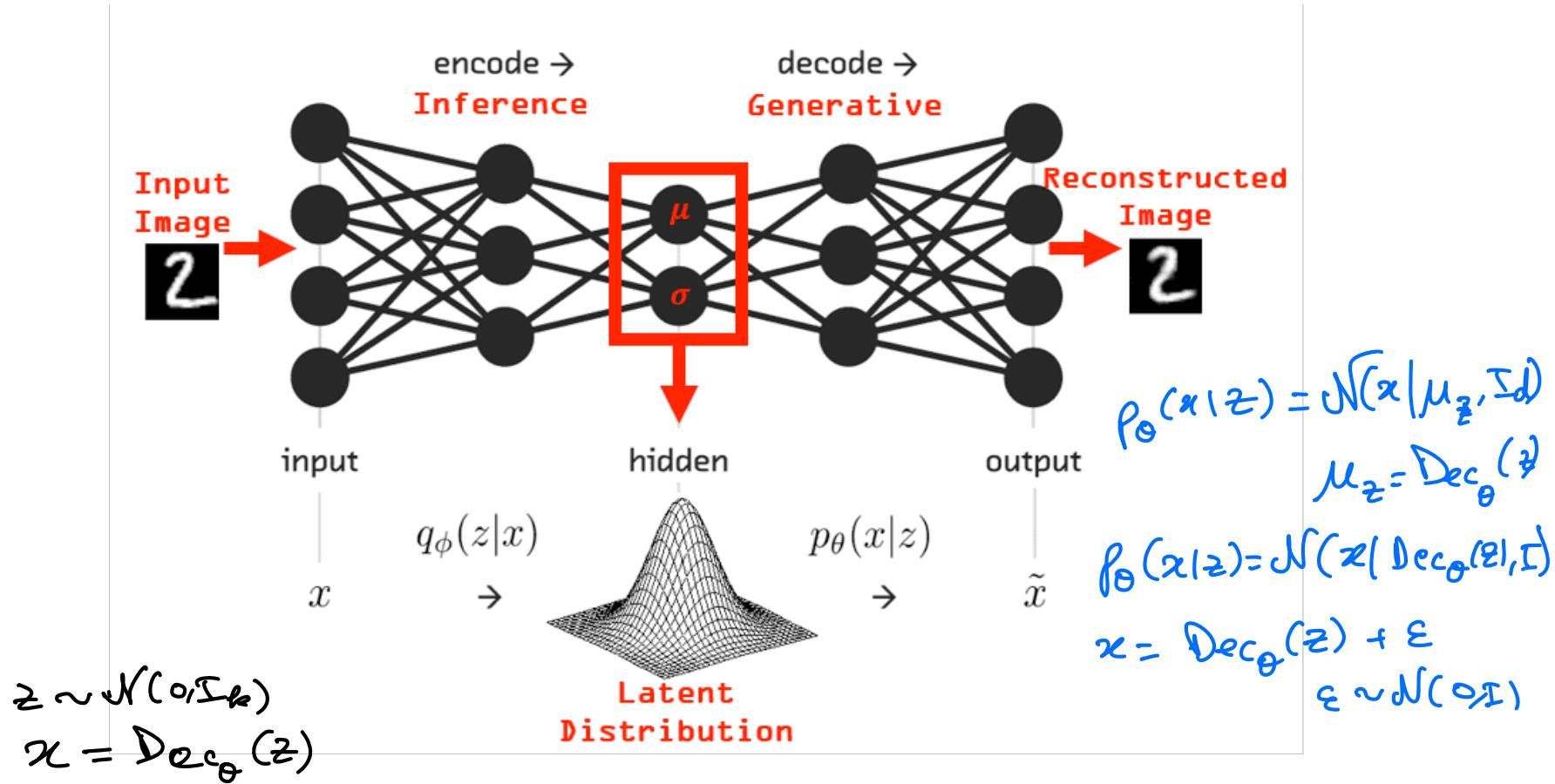
variant: denoising auto encoder
add noise ϵ_i to the input

$$\underset{i}{\operatorname{argmin}} \sum_i \|x_i - \text{Dec}_\phi(\text{Enc}_\varphi(x_i + \epsilon_i))\|^2$$

where $\tilde{x}_i = \text{Dec}_\phi(\text{Enc}_\varphi(x_i))$

A word about *variational* auto-encoders

- Variational auto-encoders (VAE) learn to reduce the dimensionality of examples in a dataset *and* produces for each example x_i a distribution of latent vectors z_i such that overall, latent vectors $z_1 \dots z_N$ are normally distributed.



Training the decoder alone

Assume we have a dataset $\{(x_i, z_i)\}_{i=1..N}$
 $x_i \in \mathbb{R}^d$, $z_i \in \mathbb{R}^k$ $k < d$.

$$\begin{aligned}\hat{\theta} = \arg\max_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i | z_i) &= \arg\max_{\theta} \sum_{i=1}^N \log N(x_i | \text{Dec}_{\theta}(z_i), I) \\ &= \arg\max_{\theta} \sum_{i=1}^N -\frac{1}{2} \|x_i - \text{Dec}_{\theta}(z_i)\|^2 \\ &= \arg\min_{\theta} \sum_{i=1}^N \|x_i - \text{Dec}_{\theta}(z_i)\|^2\end{aligned}$$

Of course, in practice we only have access to $\{x_i\}_{i=1..N}$. Thus, we will train the encoder to learn the distribution of z for each x_i .

Understanding $p(x)$ and deriving the ELBO

$z_i \sim N(0, I_k)$ I_k is $R^{k \times k}$ identity matrix.

$P_\theta(x|z)$.

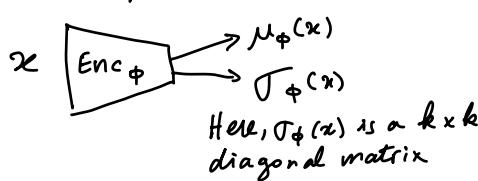
$$P_\theta(x) = \int_{R^k} P_\theta(x|z) p(z) dz \xrightarrow[\text{is to train}]{\text{the goal}} \hat{\theta} = \arg \max \sum_{i=1}^n \log P_\theta(x_i)$$

Unfortunately, this integral is hard to estimate because for almost all $z \sim p(z)$, we have $P_\theta(x|z)$ close to zero

Idea: apply the ELBO which requires

$$q_x(z) = q(z|x)$$

To represent $q(z|x)$, we will use an encoder



$$q_\phi(z|x) = N(z | \underbrace{\mu_\phi(x)}_{R^k}, \underbrace{\Sigma_\phi(x)}_{R^{k \times k}})$$

ELBO objective

$$\log P_\theta(x) \geq \text{ELBO}(n, \theta, q_\phi)$$

$$= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\ln \frac{p_\theta(x, z)}{q_\phi(z|x)} \right]$$

$$= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\ln p_\theta(x|z) - \ln \frac{q_\phi(z|x)}{p(z)} \right]$$

$$= \mathbb{E}_z \left[\ln p_\theta(x|z) \right] - KL(q_\phi(z|x), p(z))$$

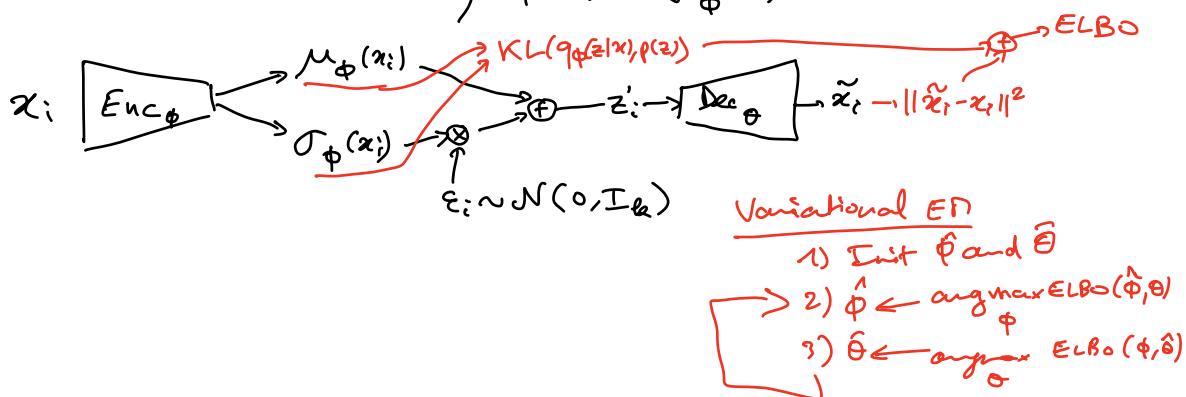
$$= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[-\frac{1}{2} \|x - \text{Dec}_\theta(z)\|^2 \right] - \underbrace{KL(q_\phi(z|x), p(z))}_{\frac{1}{2} \|\mu_\phi(x)\|^2 + \text{tr}(\Sigma_\phi(x))} + \text{const}$$

$$- \ln \log |\Sigma_\phi(x)|$$

Reparametrization trick

- Pb: ELBO will not backpropagate on ϕ (encoder) because of $z \sim q_\phi$
 \Rightarrow can't compute $\nabla_\phi \text{ELBO} \Rightarrow$ no training possible.

- trick: $z \sim q_\phi(\cdot|x) = \mathcal{N}(\mu_\phi(x), \sigma_\phi^2(x))$
is distributed identically to
 $z' = \mu_\phi(x) + \sigma_\phi(x) \cdot \varepsilon$ where $\varepsilon \sim \mathcal{N}(0, I_k)$



Variational EP

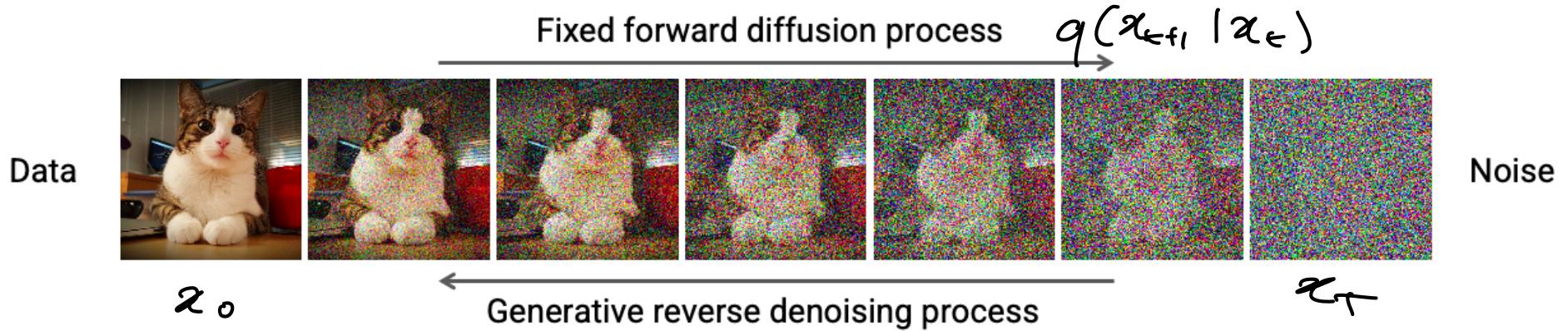
- 1) Init $\hat{\phi}$ and $\hat{\theta}$
- 2) $\hat{\phi} \leftarrow \underset{\phi}{\operatorname{argmax}} \text{ELBO}(\hat{\phi}, \theta)$
- 3) $\hat{\theta} \leftarrow \underset{\theta}{\operatorname{argmax}} \text{ELBO}(\phi, \hat{\theta})$

Outline

1. Introduction
2. Gaussian mixtures and the EM algorithm
3. Variational Analysis of the EM algorithm
4. Variational Auto-Encoders
5. **Exercise: probabilistic PCA**

Probabilistic PCA

Denoising Diffusion Models



$q(x_0)$: distribution of image in the dataset

$$q(x_t | x_{t-1}) = \mathcal{N}(x_t | \sqrt{1 - \beta_t} \cdot x_{t-1} + \beta_t \mathbb{I}_d)$$

$$x_t = x_{t-1} \times \sqrt{1 - \beta_t} + \beta_t \cdot \varepsilon_t \quad \varepsilon_t \sim \mathcal{N}(0, \mathbb{I})$$

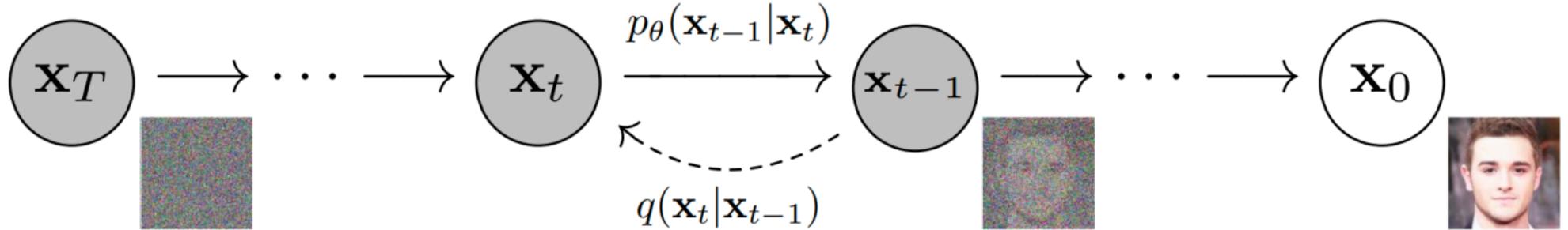
$$\beta_t \in]0, 1[$$

$$\beta_1 < \beta_2 < \beta_3 \dots$$

$$q(x_0 \dots x_T | x_0) = \prod_{t=0}^T q(x_t | x_{t-1})$$

(chain rule of probability)

Denoising Diffusion Models



Algorithm 1 Training

```

1: repeat
2:    $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ 
3:    $t \sim \text{Uniform}(\{1, \dots, T\})$ 
4:    $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
5:   Take gradient descent step on
       $\nabla_\theta \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$ 
6: until converged
  
```

Algorithm 2 Sampling

```

1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ 
2: for  $t = T, \dots, 1$  do
3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$ 
4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$ 
5: end for
6: return  $\mathbf{x}_0$ 
  
```

observation

- for simplicity, $\alpha_t = 1 - \beta_t$ $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$

$$x_t = \sqrt{\alpha_t} \cdot x_{t-1} + \sqrt{1-\alpha_t} \cdot \varepsilon_{t-1}$$

$$x_{t-1} = \sqrt{\alpha_{t-1}} \cdot x_{t-2} + \sqrt{1-\alpha_{t-1}} \cdot \varepsilon_{t-2}$$

:

$$x_t = \sqrt{\bar{\alpha}_t} \cdot x_0 + \sqrt{1-\bar{\alpha}_t} \cdot \varepsilon$$

$$q(x_t | x_0) = \mathcal{N}(x_t | \sqrt{\bar{\alpha}_t} x_0; (1-\bar{\alpha}_t) I_d)$$

- $p_\theta(x_0 \dots x_T) = p(x_T) \cdot \prod_{t=1}^T p_\theta(x_{t-1} | x_t)$

Assume $p_\theta(x_{t-1} | x_t) = \mathcal{N}(x_{t-1} | \mu_\theta(x_t, t), \Sigma_\theta(x_t, t))$

- Ideally, to learn $p_\theta(x_{t-1} | x_t)$,
we would like to compute $q(x_{t-1} | x_t) \rightarrow$ hard.
- Instead, compute $q(x_{t-1} | x_t, x_0)$

$$q(x_{t-1} | x_t, x_0)$$

$$= q(x_t | x_{t-1}, x_0) q(x_{t-1} | x_0)$$

known normally distributed

$\xrightarrow{\quad}$

$\xrightarrow{\quad}$

$q(x_t | x_0)$

conditional Bayes rule

$$P(A|BC) = \frac{P(B|AC) P(A|C)}{P(B|C)}$$

$$= \mathcal{N}(x_{t-1} | \tilde{\mu}_t(x_t, x_0), \tilde{\Sigma}_t I)$$

$$\tilde{\mu}_t(x_t, x_0) = \frac{\sqrt{\alpha_t} \cdot (1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \cdot x_t + \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \cdot x_0$$

$$x_0 = \frac{1}{\sqrt{\bar{\alpha}_t}} (x_t - \sqrt{1-\bar{\alpha}_t} \varepsilon_t)$$

$$\Rightarrow \tilde{\mu}_t = \frac{1}{\sqrt{\bar{\alpha}_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \varepsilon_t \right)$$

At the end, we want $q(x_{t-1}, x_t)$ to be close
 $\Rightarrow p_\theta(x_{t-1}, x_t)$

We want to max the ELBO, or minimize $-ELBO$

$$\begin{aligned}
 -ELBO &= \mathbb{E}_q \left[-\log \frac{p_\theta(x_0 \dots x_T)}{q(x_1 \dots x_T | x_0)} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(x_T) - \sum_{t \geq 1} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1})} \right] \\
 &= \mathbb{E}_q \left[-\log p_\theta(x_T) - \sum_{t \geq 2} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_t | x_{t-1}, x_0)} - \log \frac{p_\theta(x_0 | x_1)}{q(x_1 | x_0)} \right] \\
 &\quad \downarrow \text{Conditional Bayes rule} \quad \text{does not change the value of } q(\cdot | \cdot) \\
 &= \mathbb{E}_q \left[-\log p_\theta(x_T) - \sum_{t \geq 2} \log \left\{ \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} \cdot \frac{q(x_{t-1} | x_0)}{q(x_t | x_0)} \right\} - \log \frac{p_\theta(x_0 | x_1)}{q(x_1 | x_0)} \right] \\
 &= \mathbb{E}_q \left[-\log \frac{p_\theta(x_T)}{q(x_T | x_0)} - \sum_{t \geq 2} \log \frac{p_\theta(x_{t-1} | x_t)}{q(x_{t-1} | x_t, x_0)} - \log p_\theta(x_0 | x_1) \right] \\
 &= \mathbb{E}_q \left[\underbrace{\text{KL}(q(x_T | x_0), p_\theta(x_T))}_{L_T = \text{cote}} + \sum_{t \geq 2} \underbrace{\text{KL}(q(x_{t-1} | x_t, x_0), p_\theta(x_{t-1} | x_t))}_{L_{t-1}} - \log p_\theta(x_0 | x_1) \right]
 \end{aligned}$$

Let us focus on L_{t-1}

$$\begin{aligned}
 L_{t-1} &= \text{KL} \left(\underbrace{q(x_{t-1} | x_t, x_0)}_{\mathcal{N}(x_{t-1} | \tilde{\mu}_t, \tilde{\beta}_t I)}, \underbrace{p_\theta(x_{t-1} | x_t)}_{\mathcal{N}(x_{t-1} | M_\theta(x_t), \tilde{\beta}_t I)} \right) \\
 &= \| \tilde{\mu}_t(x_t, x_0) - M_\theta(x_t) \|^2 \times \text{cote}
 \end{aligned}$$

recall that

$$\hat{x}_t = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \cdot \varepsilon_t \right)$$

let us define

$$u_0(x_t) = \frac{1}{\sqrt{\alpha_t}} \left(x_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \cdot \varepsilon_0(x_t, t) \right)$$

$$\begin{aligned} L_{t-1} &= \| \varepsilon_t - \varepsilon_0(x_t, t) \|^2 \times \text{cste}' \\ &= \| \varepsilon_t - \varepsilon_0(\sqrt{\bar{\alpha}_t} \cdot x_0 + \sqrt{1-\bar{\alpha}_t} \cdot \varepsilon_t, t) \|^2 \times \text{cste}' \end{aligned}$$