

Generative modelling challenge

Second session

Alain Durmus
École polytechnique
November 15, 2022

Variational Autoencoders

Introduction to generative models

Let $X \subset \mathbb{R}^d$.

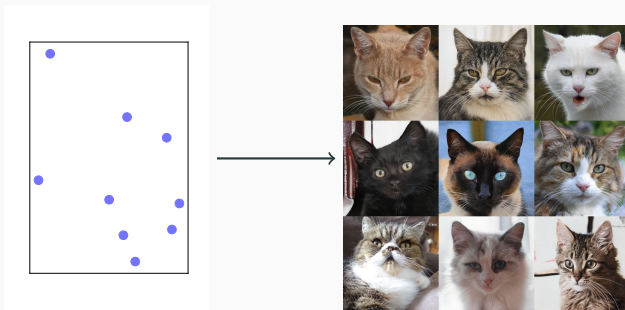
- **Input**

Data $\{x_i\}_{i=1}^N$: N i.i.d. observations from an unknown distribution $\mu^* \in \mathcal{P}(X)$.

Notation: $\mathcal{P}(X)$ space of probability measures on (X, \mathcal{X}) .

- **Output**

New samples from μ^*



thiscatdoesnotexist.com

- **Input** For $X \subset \mathbb{R}^d$.
Data $\{x_i\}_{i=1}^N : N$ i.i.d. observations from $\mu^* \in \mathcal{P}(X)$ unknown
- **Output**
New samples from μ^*
- Consider a parametric family of distributions $\{\mu_\theta : \theta \in \Theta\}$.
- **Minimum distance estimation :**
 - minimize $\theta \mapsto \mathbf{D}(\mu_\theta | \mu^*)$ where \mathbf{D} is a divergence over the space of probability measure on X .
 - Sample a new observation from μ_{θ^*} .

Divergence over $\mathcal{P}(X)$

- A divergence on $\mathcal{P}(X)$, is a function $\mathbf{D} : \mathcal{P}(X)^2 \rightarrow \mathbb{R}_+$ which satisfies *the most important* axiom of a distance: for μ, ν two probability measures

$$\mathbf{D}(\mu|\nu) = 0 \text{ if and only if } \mu = \nu .$$

- Any distance on $\mathcal{P}(X)$ is a divergence.
- Do not satisfy the triangle inequality in general except if it is a distance...
- It is not symmetric in general.
- Important example:

$$\text{KL}(\mu \parallel \nu) = \begin{cases} \int \log(d\mu/d\nu) d\nu & \text{if } \mu \ll \nu \\ +\infty & \text{otherwise} . \end{cases}$$

Minimum Distance Estimation: an ideal?

- **Input** For $X \subset \mathbb{R}^d$.
Data $\{x_i\}_{i=1}^N$: N i.i.d. observations from $\mu^* \in \mathcal{P}(X)$ unknown
- **Output**
New samples from μ^*
- Consider a parametric family of distributions $\{\mu_\theta : \theta \in \Theta\}$.
- **Minimum distance estimation** :
 - minimize $\theta \mapsto \mathbf{D}(\mu_\theta | \mu^*)$ where \mathbf{D} is a divergence over the space of probability measure on X .
 - Sample a new observation from μ_{θ^*} .
- Problem: μ^* is not known and in general $\mathbf{D}(\mu_\theta | \mu^*)$ is not tractable.
- One solution presented here: likelihood estimation.

- Consider the case $X = \mathbb{R}^p$ and $\Theta \subset \mathbb{R}^d$.
- Choice for the family $\{\mu_\theta : \theta \in \Theta\}$?

$$\{\mu_\theta : \mu_\theta \ll \text{Leb}, \quad p_\theta = d\mu_\theta/d\text{Leb}\}.$$

Example:

p_θ : density w.r.t. Leb of $\mathcal{N}(m, \sigma^2)$.

- We should be able to sample from p_θ for any θ ...
- Choice for the divergence **D**?

Maximum likelihood estimation

- Consider the case $X = \mathbb{R}^p$ and $\Theta \subset \mathbb{R}^d$.
- Choice for the family $\{\mu_\theta : \theta \in \Theta\}$?

$$\{\mu_\theta : \mu_\theta \ll \text{Leb}, \quad p_\theta = d\mu_\theta/d\text{Leb}\} .$$

- Choice for the divergence **D**?
- **D** = KL: problem $\text{KL}(\mu_\theta \parallel \hat{\mu}_N) = \infty \dots$
- Recall that ideally if **D** = KL, we would like to minimize

$$\text{KL}(\mu^\star \parallel \mu_\theta) = - \int d\mu^\star \log \left(\frac{d\mu_\theta}{d\mu^\star} \right) .$$

- This is equivalent to maximize if $\mu^\star \ll \text{Leb}$,

$$\theta \mapsto \int d\mu^\star \log \left(\frac{d\mu_\theta}{d\text{Leb}} \right) .$$

- Solution: replace the integral by an empirical version

$$\theta \mapsto N^{-1} \sum_{i=1}^N \log p_\theta(x_i) .$$

Density estimation

- Assume now

$$\mu^* \ll \text{Leb} , \text{ with density } p^* .$$

- Choice for the family $\{\mu_\theta : \theta \in \Theta\}$?

$$\{\mu_\theta : \mu_\theta \ll \text{Leb} , \quad p_\theta = d\mu_\theta/d\text{Leb}\} .$$

- We should be able to sample from μ_θ for any θ and in the same time the family has to be sufficiently rich/large.
- First solution:

$$\mu_\theta = (T_\theta)_\# \nu_0 ,$$

where

$$\nu_0 \in \mathcal{P}(\mathbb{R}^p) \text{ with density } q_0 , \quad T_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^p \text{ is a } C^1 \text{ diffeomorphism} .$$

- In that case, we have

$$p_\theta(x) = q_0(T_\theta^{-1}(x)) \mathbf{Jac}_x[T_\theta^{-1}](x) .$$

- Example: $T_\theta(x) = m + \Sigma x$, $\theta = (m, \Sigma)$.

- Approximating the function $p^*(x)$ is called the density estimation problem.
- Suppose we (somehow) optimize a generator p_θ

$$\theta \mapsto \int d\mu^* \log \left(\frac{d\mu_\theta}{d\text{Leb}} \right) \text{ or } \theta \mapsto \sum_{i=1}^N \log p_\theta(x_i) ,$$

such that $p_\theta \approx p^*$.

- If we can compute $T_\theta^{\leftarrow}(x)$ and $\text{Jac}_x[T_\theta^{\leftarrow}](x)$ then we have a density estimator.
- In many cases, generation (sampling) is easier than density estimation.

- How to optimize?

$$\theta \mapsto \int d\mu^* \log \left(\frac{d\mu_\theta}{d\text{Leb}} \right) \text{ or } \theta \mapsto \sum_{i=1}^N \log p_\theta(x_i) . \quad (1)$$

- Stochastic gradient descent:

$$\theta_{k+1} = \theta_k + \gamma_{k+1} \sum_{x_i \in B_{k+1}} \nabla_\theta [\log p_{(\cdot)}(x_i)](\theta_k) ,$$

where

$(B_k)_k$ is a sequence of random batch of data points ,

$(\gamma_k)_k$ is a sequence of stepsizes/learning rates .

- Under some assumptions, it can be shown that almost surely $(\theta_k)_{k \in \mathbb{N}}$ converges to some minimizers of (1).

Density estimation

- Choice for the family $\{\mu_\theta : \theta \in \Theta\}$? $\{\mu_\theta : \mu_\theta \ll \text{Leb}, \quad p_\theta = d\mu_\theta/d\text{Leb}\}$.

- First solution:

$$\mu_\theta = (T_\theta)_\# \nu_0 ,$$

where

$\nu_0 \in \mathcal{P}(\mathbb{R}^p)$ with density q_0 , $T_\theta : \mathbb{R}^p \rightarrow \mathbb{R}^p$ is a C^1 diffeomorphism .

- In that case, we have

$$p_\theta(x) = q_0(T_\theta^{-1}(x)) \text{Jac}_{T_\theta^{-1}}(x) .$$

- What people thought it was hard:

find $\{T_\theta : \theta \in \Theta\}$ such that $\theta \mapsto \sum_{i=1}^N \log p_\theta(x_i)$ easy to optimize...

- It turns out that such constructions are now possible using neural networks; see normalizing flows [Rezende and Mohamed \(2015\)](#)!
- Here we present a first alternative using MLE: latent variable models.

- Here we aim to construct a family $\{p_\theta : \theta \in \Theta\}$ which is expressive enough.
- Given some samples $\{x_i\}_{i=1}^N$ and a samplable prior r .
- Generative latent variable model:
 1. $z \sim r$;
 2. $x \sim p_\theta(x|z)$
- Typically, $\log p_\theta$ is the error function of a Neural Network:

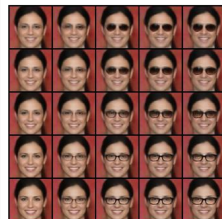
$$\log p_\theta(x, z) = -\ell(T_\theta(z), x) .$$

- This corresponds to the marginal likelihood:

$$\begin{aligned} p_\theta(x) &= \int p_\theta(x|z)r(z)dz \\ &= \int p_\theta(x, z)dz . \end{aligned}$$

Motivation for latent variables

- Conditioning on the latent variable z may give to the samples x more global coherence.
- Learned latent distribution might reveal structure in the data distribution.
- The latent variables could be useful for downstream tasks or interpretability.



Klys, Snell, and Zemel,
Neurips 2018

MLE for latent variable models?

- Marginal likelihood:

$$p_{\theta}(x) = \int p_{\theta}(x|z)r(z)dz = \int p_{\theta}(x, z)dz .$$

- Based on samples $\{x_i\}_{i=1}^N$:

$$\text{fit the MLE: } \theta^* \in \arg \max \left[N^{-1} \sum_{i=1}^N \log p_{\theta}(x_i) \right] .$$

- This doesn't look promising since $p_{\theta}(x_i)$ are intractable...

The evidence lower bound (ELBO)

- A first option to approximate $p_\theta(x_i)$ is to use importance sampling.
- Using Jensen inequality, we can show that for any condition distribution q :

$$\log p_\theta(x) = \log \int \frac{p_\theta(x, z)}{q(z|x)} q(z|x) dz \geq \int \log \left[\frac{p_\theta(x, z)}{q(z|x)} \right] q(z|x) dz .$$

- The RHS is a lower-bound on the marginal log-likelihood, referred to as ELBO [MacKay \(1992\)](#).
- This is this quantity that we maximize through an empirical version:

$$\int \log \left[\frac{p_\theta(x, z)}{q(z|x)} \right] q(z|x) dz \approx M^{-1} \sum_{i=1}^M \log \left[\frac{p_\theta(x_i, z)}{q(z_i|x)} \right] ,$$

with $z_i \stackrel{\text{iid}}{\sim} q(\cdot|x)$.

Do we lose something compared to usual MLE?

- Define the ELBO:

$$\text{ELBO}(x_{1:N}; \theta, q) = N^{-1} \sum_{i=1}^N \int \log \left[\frac{p_{\theta}(x_i, z)}{q(z|x_i)} \right] q(z|x_i) dz .$$

- This satisfies:

$$N^{-1} \sum_{i=1}^N \log p_{\theta}(x_i) \geq \text{ELBO}(x_{1:N}; \theta, q) ,$$

with equality iif:

$$q(z|x) = p_{\theta}(z|x) .$$

- Therefore, the MLE can be re-written:

$$\begin{aligned} \theta^* &\in \arg \max \left[N^{-1} \sum_{i=1}^N \log p_{\theta}(x_i) \right] \\ &\in \arg \max_{\theta, q} \text{ELBO}(x_{1:N}; \theta, q) . \end{aligned}$$

- How to estimate/learn $q(z|x) \approx p_{\theta}(z|x)$?

- How to estimate/learn $q(z|x) \approx p_\theta(z|x)$?
- Let $q_\phi(z|x)$ be a family of density estimators with parameters ϕ .
- People sometimes call this amortized inference.
- The approximate problem, we consider then

maximize over θ, ϕ the function $\text{ELBO}(x_{1:N}; \theta, \phi)$,

writing $\text{ELBO}(x_{1:N}; \theta, \phi) = \text{ELBO}(x_{1:N}; \theta, q_\phi)$.

- This option has been popularized in [Kingma and Welling \(2013\)](#).

- Recall that

$$\text{ELBO}(x_{1:N}; \theta, \phi) = N^{-1} \sum_{i=1}^N f_i(\theta, \phi), \quad f_i = \int \log \left[\frac{p_\theta(x_i, z)}{q_\phi(z|x_i)} \right] q_\phi(z|x_i) dz.$$

- How to optimize this function over θ, ϕ ?
- Solution: SGD
- For any i , we have the Monte Carlo approximation:

$$\nabla_\theta f_i(\theta, \phi) \approx M_i \sum_{k=1}^{M_i} \nabla_\theta \log \left[\frac{p_\theta(x_i, z_k)}{q_\phi(z_k|x_i)} \right]$$

with $z_k \stackrel{\text{iid}}{\sim} q_\phi(\cdot|x_i)$.

- However, the computation/approximation of $\nabla_\phi f_i(\theta, \phi)$ is not that easy.

The reparametrization Trick

- Recall that

$$\text{ELBO}(x_{1:N}; \theta, \phi) = N^{-1} \sum_{i=1}^N f_i(\theta, \phi), \quad f_i = \int \log \left[\frac{p_\theta(x_i, z)}{q_\phi(z|x_i)} \right] q_\phi(z|x_i) dz.$$

- To estimate $\nabla_\phi f_i(\theta, \phi)$, we always assume that $q_\phi(\cdot|x)$ is the distribution of some random variable

$$T_\phi(x, \varepsilon), \quad \varepsilon \sim r_0,$$

for some samplable distribution r_0 .

- Example: $\mu_\phi(x) + \sigma_\phi(x)\varepsilon, \varepsilon \sim \mathcal{N}(0, \text{Id})$
- This is called the reparametrization Trick.
- In that case:

$$\nabla_\phi f_i(\theta, \phi) \approx K_i^{-1} \sum_{k=1}^{K_i} \nabla_\phi \log \left[\frac{p_\theta(x_i, T_\phi(x, \varepsilon_k))}{q_\phi(T_\phi(x, \varepsilon_k)|x_i)} \right],$$

with $\varepsilon_k \stackrel{\text{iid}}{\sim} r_0$.

References

- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- David JC MacKay. Bayesian interpolation. *Neural computation*, 4(3):415–447, 1992.
- Danilo Rezende and Shakir Mohamed. Variational inference with normalizing flows. In *International conference on machine learning*, pages 1530–1538. PMLR, 2015.
- J. Wolfowitz. The minimum distance method. *Ann. Math. Statist.*, 28(1):75–88, 03 1957. doi: 10.1214/aoms/1177707038.