

Neural nets and latent variable models

Statistiques au sommet de Rochebrune

Pierre Gloaguen

March 25, 2024

Université Bretagne Sud

Neural networks

- A parameterized function:

$$\begin{aligned} f_{\theta} : \mathbb{R}^{d_x} &\rightarrow \mathbb{R}^{d_y} \\ x &\mapsto f(x|\theta) \end{aligned} ,$$

which is a composition of non linear functions.

- The function depends both on θ and some intermediate non linear functions.

- A parameterized function:

$$\begin{aligned} f_{\theta} : \mathbb{R}^{d_x} &\rightarrow \mathbb{R}^{d_y} \\ x &\mapsto f(x|\theta) , \end{aligned}$$

which is a composition of non linear functions.

- The function depends both on θ and some intermediate non linear functions.
- For instance, for $d_x = 2, d_y = 1$, consider a set of *biases* and **Weights**
 $\theta = \{b_1 \in \mathbb{R}^4, b_2 \in \mathbb{R}, \mathbf{W}_1 \in \mathbb{R}^{4 \times 2}, \mathbf{W}_2 \in \mathbb{R}^{1 \times 4};$

$$f(x|\theta) = \tanh(\mathbf{W}_2 \cdot \text{sigmoid}(\mathbf{W}_1 x + b_1) + b_2) .$$

We can introduce some intermediary notations and write

$$\sigma_1(z) := \text{sigmoid}(z)$$

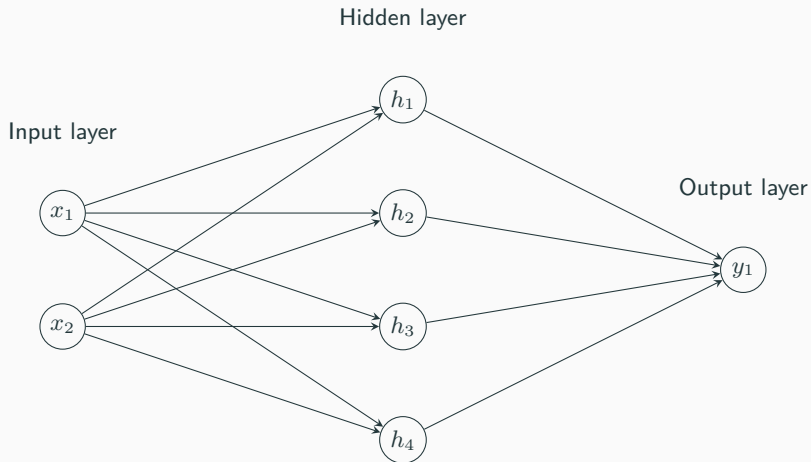
$$\sigma_2(z) := \tanh(z)$$

$$h^{(1)} := \sigma_1(\mathbf{W}_1 x + b_1) ,$$

We then have

$$f(x|\theta) = \sigma_2(\mathbf{W}_2 h^{(1)} + b_2) .$$

Graphical representation of a neural network



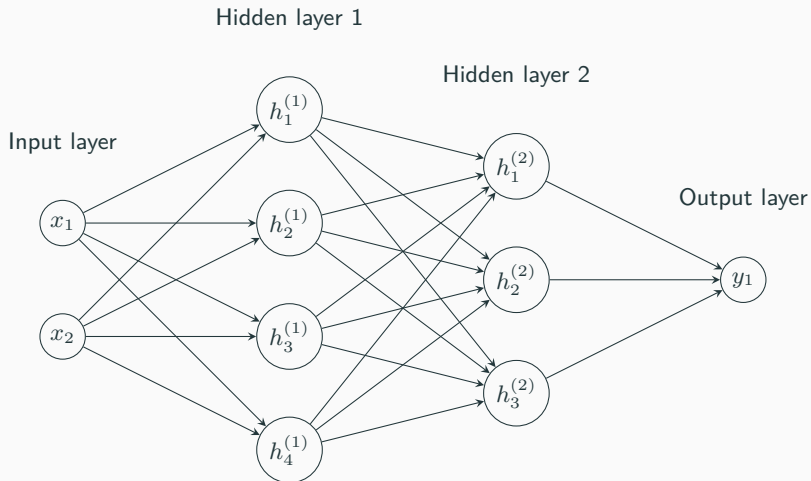
- Build a sequence of vectors $(h^{(k)})_{0 \leq k \leq D}$ where D is the depth of a network:

$$h^{(0)} = x$$

$$h^{(k)} = \sigma_k(\mathbf{W}_k h^{(k-1)} + b_k), \text{ for } 1 \leq k \leq D$$

$$f(x|\theta) = h^{(D)}$$

Graphical representation of a deep neural network



- There exists a dictionary of classical non linear functions
 - $\text{ReLU}(x) = \max(0, x)$;
 - $\text{sigmoid}(x) = \frac{1}{1+e^{-x}}$;
 - $\text{tanh}(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$;
 - etc...

- In the fully connected context, each neuron is a linear combination of all neurons of previous layer;
- No relation is imposed between each linear combination (over the rows of \mathbf{W}).
- For computational efficiency, and to impose some invariance, one can impose that linear combinations only apply locally, and are the same on the different regions of each layer.

Simple convolution filter

- For instance, to create a layer which gives the local differences in the vector:

$$x = \begin{pmatrix} 2 \\ 1 \\ -6 \\ 4 \\ 5 \end{pmatrix}$$

- Define a kernel

$$\begin{pmatrix} -1 \\ 0 \\ 1 \end{pmatrix}$$

The resulting *convolutional layer* is a vector h of size $5 - 3 + 1$ where, for $1 \leq i \leq 3$:

$$h = \begin{pmatrix} -1 & 0 & 1 & 0 & 0 \\ 0 & -1 & 0 & 1 & 0 \\ 0 & 0 & -1 & 0 & 1 \end{pmatrix} \times x = \begin{pmatrix} -8 \\ 3 \\ 11 \end{pmatrix}$$

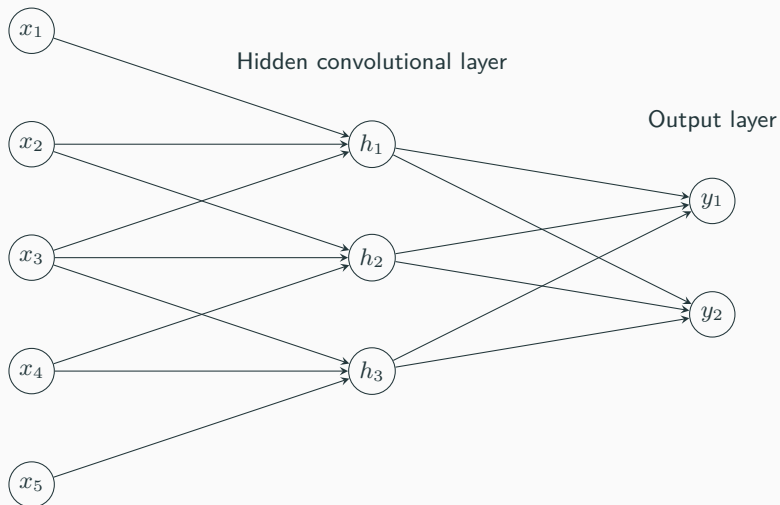
- *Remark* Of course, technical extensions exists to deal with bounds of x (stride and padding).

- A convolutional layer imposes a sparse weight matrix whose parameters are shared among rows:
- For instance, if the convolutional filter implies 3 elements at once:

$$\mathbf{W} = \begin{pmatrix} w_1 & w_2 & w_3 & 0 & \dots & 0 \\ 0 & w_1 & w_2 & w_3 & 0 & \dots \\ \vdots & \vdots & \ddots & \ddots & \ddots & \vdots \\ 0 & \dots & 0 & w_1 & w_2 & w_3 \end{pmatrix}$$

Convolutional neural network

Input layer



- Neural nets were first used for classical learning tasks:
 - Classification;
 - Regression.
- In this context, we have a data set $(x_{1:n}, y_{1:n})$:
 - x_k are the features (explanatory variables):
 - y_k is the response, a variable to predict (a label or a value).
- $f(x_k|\theta)$ is a prediction for y_k .

Neural nets as generators

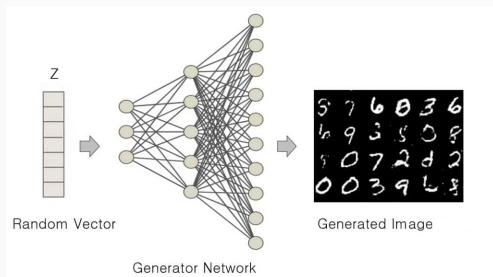
- Another use of neural nets is for generative modelling;
- Probabilistic framework where the high dimensional data \mathbf{Y} are assumed to arise from a lower dimension random variable \mathbf{Z} with a known distribution.

$$\mathbf{Z} \sim \text{Distribution}(\cdot)$$

$$\mathbf{Y}|\mathbf{Z} = f(\mathbf{Z}|\theta) + \mathbf{E},$$

where:

- f is a neural network;
- \mathbf{E} is some random variable (reconstruction error term).



Classical losses functions

- The parameter θ is learnt by minimizing a loss function $\mathcal{L}(\theta)$:
 - In the regression case (for $y_i \in \mathbb{R}$):

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x_i|\theta))^2$$

- In the K classes classification (encoding $y_i = (0, \dots, 0, \overset{k\text{-th value}}{1}, 0, \dots, 0)^T$), and with $f(x_i|\theta)$ being a vector of probabilities:

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n y_i^T \log f(x_i|\theta).$$

- In the generative model case, **ideally**, the loglikelihood:

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \log p(y_i|\theta) = \frac{1}{n} \sum_{i=1}^n \log \int p(y_i|z\theta)p(z)dz$$

- Note that, in general, we have:

$$\mathcal{L}(\theta) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}_i(\theta),$$

Learning parameters of a neural net: stochastic gradient descent (SGD)

Gradient descent to find $\hat{\theta} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta)$.

Starting from $\theta^{(0)}$ and building a sequence $(\theta^{(\ell)})_{\ell \geq 0}$ such that, for $\ell \geq 1$:

$$\theta^{(\ell)} = \theta^{(\ell-1)} + \gamma_{\ell} \nabla_{\theta} \mathcal{L}(\theta^{(\ell-1)}),$$

where γ_{ℓ} is a sequence of step sizes decreasing (at an appropriate rate) to 0.

Learning parameters of a neural net: stochastic gradient descent (SGD)

Gradient descent to find $\hat{\theta} = \operatorname{argmin}_{\theta} \mathcal{L}(\theta)$.

Starting from $\theta^{(0)}$ and building a sequence $(\theta^{(\ell)})_{\ell \geq 0}$ such that, for $\ell \geq 1$:

$$\theta^{(\ell)} = \theta^{(\ell-1)} + \gamma_{\ell} \nabla_{\theta} \mathcal{L}(\theta^{(\ell-1)}),$$

where γ_{ℓ} is a sequence of step sizes decreasing (at an appropriate rate) to 0.

Stochastic gradient descent

In general, with n observations, we have:

$$\begin{aligned} \nabla_{\theta} \mathcal{L}(\theta) &= \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} \mathcal{L}_i(\theta) \\ &= \mathbb{E}_J[\nabla_{\theta} \mathcal{L}_J(\theta)], \end{aligned} \quad \text{where } J \sim \mathcal{U}[[1; n]].$$

Thus, for $1 \leq m \ll n$,

$$\widehat{\nabla}_{\theta} \mathcal{L}(\theta) := \frac{1}{m} \sum_{k=1}^m \nabla_{\theta} \mathcal{L}_{J_k}(\theta) \quad \text{where } J_k \stackrel{\text{i.i.d.}}{\sim} \mathcal{U}[[1; n]],$$

is an unbiased Monte Carlo estimator of the wanted gradient.

- SGD still requires to compute the gradient of $\nabla \mathcal{L}_i(\theta)$;
- Recall that $f(x|\theta)$ is built by setting:

$$h^{(0)} = \begin{pmatrix} 1 \\ x \end{pmatrix} \in \mathbb{R}^{d_0} \quad \text{with } d_0 = d_x + 1$$

$$h^{(1)} = \sigma_1(\theta_1 h^{(k-1)}), \quad \text{with } \theta_1 = (b_1, \mathbf{W}_1), \quad b_1 \in \mathbb{R}^{d_1}, \quad \mathbf{W}_1 \in \mathbb{R}^{d_1 \times d_0}$$

\vdots

$$f(x|\theta) = \sigma_D(\theta_D h^{(D-1)}), \quad \text{with } \theta_d \in \mathbb{R}^{d_y \times d_{D-1}}$$

- Then¹:

$$\nabla_{\theta} \mathcal{L}_i(\theta) = \begin{pmatrix} \frac{\partial \mathcal{L}_i(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial \mathcal{L}_i(\theta)}{\partial \theta_D} \end{pmatrix}$$

¹As the gradient is a vector, all partial derivatives are implicitly vectorized here

Notations

For $1 \leq k \leq D$

- $z^{(k)} := \theta_k h^{((k-1))}$, and recall that $h^{(k)} = \sigma_k(z^{(k)})$;
- $\frac{\partial h^{(k)}}{\partial z^{(k)}}$ is a \mathbb{R}^{d_k} -vector of elementwise derivatives;
- \odot : elementwise product for a vector, \times : matrix product;

Algorithm

- For all $1 \leq k \leq D$:

$$\frac{\partial \mathcal{L}_i(\theta)}{\partial \theta_k} \stackrel{d_k \times d_{k-1} \text{ matrix}}{=} \left(\delta_k \odot \frac{\partial h^{(k)}}{\partial z^{(k)}} \right) \times h^{(k-1)T},$$

where $\delta_k := \frac{\partial \mathcal{L}_i(\theta)}{\partial h^{(k)}} \in \mathbb{R}^{d_k}$

- For any $1 \leq k < D$, δ_k satisfies the **backward recursion**:

$$\delta_k = \theta_{k+1}^T \times \left(\delta_{k+1} \odot \frac{\partial h^{(k+1)}}{\partial z^{(k+1)}} \right) .$$

Learning parameters of a neural net:

- Note that most times are used multiple times, and easily computed;
- Computation of simple derivatives is easily done using automatic differentiation during the forward pass;
- A backward pass computes the wanted gradients.

Latent variable models

Neural network generative model

For $1 \leq k \leq n$:

$$Z_k \sim \mathcal{N}_{d_z}(0, I_{d_z})$$

Prior distribution

$$Y_k | Z_k \sim \mathcal{N}_{d_y}(\text{NeuralNet}(Z|\theta), \sigma^2 I_{d_y}),$$

Observation model

- Used to generate complex structures

- On n sites, observers count abundances of p species resulting in a $n \times p$ matrix \mathbf{Y}
- On each site k , m covariates are measure, resulting in a matrix $n \times m$ matrix \mathbf{X}
- The intensity of presence is supposed to depend on covariates through a linear combination:

Poisson log-normal model

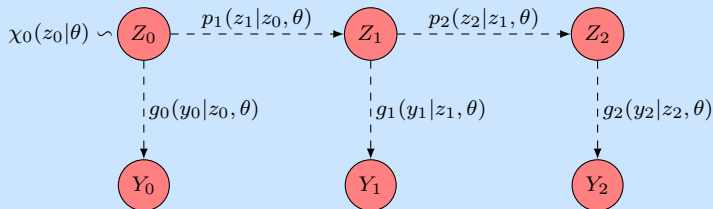
$\mathbf{Z} = (Z_{i,j})_{1 \leq i \leq n, 1 \leq j \leq p}, Z_{i,j} \stackrel{\text{ind.}}{\sim} \mathcal{N}(0, 1)$ Prior distribution

$\mathbf{Y}|\mathbf{Z} \sim \text{Poisson}(\exp(\mathbf{X}\theta + \mathbf{Z}))$, Observation model

where θ is a $m \times p$ unknown matrix to estimate.

- Used in joint species modelling

Hidden Markov models



$$\begin{aligned} Z_0 &\sim \chi_0(z_0, \theta) \\ Z_t | \{Z_{t-1} = z_{t-1}\} &\sim p_t(z_t | z_{t-1}, \theta) \\ Y_t | \{Z_t = z_t\} &\sim g_t(y_t | z_t, \theta) \end{aligned}$$

Prior distribution

Observation model

- Use to model incomplete time series

- Classical learning objectives:
 - Computing the log-likelihood:

$$\log p(y_{1:n}) = \log \int_{\Omega_{\mathbf{z}}} p(z)p(y_{1:n}|z)dz$$

- Learning the posterior distribution:

$$p(z|y_{1:n}) = \frac{p(z)p(y_{1:n}|z)}{p(y_{1:n})}$$

- In a bayesian context, learn the full posterior

$$p(z, \theta|y_{1:n}) = \frac{p(z, \theta)p(y_{1:n}|z, \theta)}{p(y_{1:n})}$$

Variational autoencoder

Poisson log-normal model

$\mathbf{Z} = (Z_{i,j})_{1 \leq i \leq n, 1 \leq j \leq p}, Z_{i,j} \stackrel{\text{ind.}}{\sim} \mathcal{N}(0, 1)$ Prior distribution

$\mathbf{Y}|\mathbf{Z} \sim \text{Poisson}(\exp(\mathbf{X}\theta + \mathbf{Z}))$, Observation model

where θ is a $m \times p$ unknown matrix to estimate.

- In this context, we aim at learning θ .
- Ideally, we aim at maximizing the log-likelihood:

$$\log p(\mathbf{Y}|\theta) = \log \int_{\Omega_{\mathbf{Z}}} p(\mathbf{Y}|z, \theta) p(z) dz$$

- However the integral is not tractable in general;
- Approximating it by Monte Carlo is not likely to work;

Evidence lower Bound

- Let $q(z)$ be a probability density function of Ω_Z :

$$\begin{aligned}\log p(\mathbf{Y}|\theta) &= \log \int_{\Omega_Z} p(\mathbf{Y}|z, \theta)p(z)dz \\ &= \log \int_{\Omega_Z} \frac{p(\mathbf{Y}|z, \theta)p(z)}{q(z)}q(z)dz \\ &= \log \mathbb{E}_q \left[\frac{p(\mathbf{Y}|Z, \theta)p(Z)}{q(Z)} \right] \\ &\geq \mathbb{E}_q \left[\log \frac{p(\mathbf{Y}|Z, \theta)p(Z)}{q(Z)} \right] && \text{Jensen's inequality} \\ &= \mathbb{E}_q[\log p(\mathbf{Y}|Z, \theta)] - \mathbb{E}_q \left[\log \frac{q(Z)}{p(Z)} \right] \\ &= \mathbb{E}_q[\log p(\mathbf{Y}|Z, \theta)] - \text{KL}(q(z) \parallel p(z)) \\ &:= \text{ELBO}(q, \theta)\end{aligned}$$

- Note that there is equality iff $q(z) = p(z|\mathbf{Y}, \theta)$.

- Instead of maximizing the likelihood, variational inference maximizes the ELBO with respect to q and θ :

$$(\hat{q}, \hat{\theta}) = \operatorname{argmax}_{q, \theta} \text{ELBO}(q, \theta).$$

- As q is a p.d.f., this maximization is unfeasible (unless $p(z|\mathbf{Y}, \theta)$ is known).
- In variational inference, the family of possible distributions for q is restricted to a “simple” parametric family;
- For instance, when there is one latent variable per observation, one can set

$$q(z_{1:n}) = \prod_{i=1}^n q(z_i),$$

where $q(z_i)$ is the p.d.f. of a $\mathcal{N}_{d_z}(\mu_i, \sigma_i^2 I_{d_z})$.

- Note that choosing a family for q imposes to choose both a dependence structure and parametric distributions.

Example on PLN

Pose a mean field gaussian family for the approximation of $\mathbf{Z}|Y$:

$$q(\mathbf{Z}) = \prod_{i=1}^n \prod_{j=1}^p q_{i,j}(z_{i,j}), \text{ where } q_{i,j}(z_{i,j}) : \text{ p.d.f. of } \mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2)$$

$$\begin{aligned} \text{ELBO}(q, \theta) &= \mathbb{E}_q[\log p(\mathbf{Y}|Z, \theta)] - \text{KL}(q(z) \parallel p(z)) \\ &= \sum_{i=1}^n \sum_{j=1}^p (\mathbb{E}_{q_{i,j}}[\log p(Y_{i,j}|Z_{i,j}, \theta)] - \text{KL}(q_{i,j}(z) \parallel p(z))). \end{aligned}$$

Example on PLN

Pose a mean field gaussian family for the approximation of $\mathbf{Z}|Y$:

$$q(\mathbf{Z}) = \prod_{i=1}^n \prod_{j=1}^p q_{i,j}(z_{i,j}), \text{ where } q_{i,j}(z_{i,j}) : \text{ p.d.f. of } \mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2)$$

$$\begin{aligned} \text{ELBO}(q, \theta) &= \mathbb{E}_q[\log p(\mathbf{Y}|Z, \theta)] - \text{KL}(q(z) \parallel p(z)) \\ &= \sum_{i=1}^n \sum_{j=1}^p (\mathbb{E}_{q_{i,j}}[\log p(Y_{i,j}|Z_{i,j}, \theta)] - \text{KL}(q_{i,j}(z) \parallel p(z))) . \end{aligned}$$

- For a current estimate $\hat{\theta}^{(0)}$, maximizing w.r.t. q results in $n \times p$ maximizations of the functions

$$f(\mu_{i,j}, \sigma_{i,j}^2) = Y_{i,j} \mu_{i,j} - \exp\left(\left(\mathbf{X}\hat{\theta}^{(0)}\right)_{i,j} + \mu_{i,j}^2 + \frac{\sigma_{i,j}^2}{2}\right) - \frac{1}{2}(\mu_{i,j}^2 + \sigma_{i,j}^2 - \log \sigma_{i,j}^2)$$

- This results in $2 \times n \times p$ parameters for the distribution of $\mathbf{Z}|Y$.
- Note that no explicit link is made between each estimated $(\mu_{i,j}, \sigma_{i,j}^2)$ especially, the dependence on $Y_{i,j}$ is not explicit.

- To potentially reduce the number of parameters, and create functional link between $Y_{i,j}$ and $\mu_{i,j}$, one can add an *amortization* constraint:

$$q(\mathbf{Z}) = \prod_{i=1}^n \prod_{j=1}^p q_{i,j}(z_{i,j}),$$
$$q_{i,j}(z_{i,j}) \sim \mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2)$$
$$(\mu_{i,j}, \sigma_{i,j}^2) = \text{NeuralNet}(Y_{i,j} | \lambda),$$

where λ is a parameter common to all $Y_{i,j} \mapsto (\mu_{i,j}, \sigma_{i,j}^2)$.

- To potentially reduce the number of parameters, and create functional link between $Y_{i,j}$ and $\mu_{i,j}$, one can add an *amortization* constraint:

$$q(\mathbf{Z}) = \prod_{i=1}^n \prod_{j=1}^p q_{i,j}(z_{i,j}),$$
$$q_{i,j}(z_{i,j}) \sim \mathcal{N}(\mu_{i,j}, \sigma_{i,j}^2)$$
$$(\mu_{i,j}, \sigma_{i,j}^2) = \text{NeuralNet}(Y_{i,j}|\lambda),$$

where λ is a parameter common to all $Y_{i,j} \mapsto (\mu_{i,j}, \sigma_{i,j}^2)$.

- The ELBO becomes a function of (θ, λ) . The part to maximize in λ is:

$$f(\lambda) = \sum_{i,j=1}^{n,p} \left\{ Y_{i,j} \mu^\lambda(Y_{i,j}) - \exp \left((\mathbf{X}\hat{\theta}^{(0)})_{i,j} + (\mu^\lambda(Y_{i,j}))^2 + \frac{\sigma^\lambda(Y_{i,j})}{2} \right) \right. \\ \left. - \frac{1}{2} \left((\mu^\lambda(Y_{i,j}))^2 + \sigma^\lambda(Y_{i,j}) - \log \sigma^\lambda(Y_{i,j}) \right) \right\}$$

Variational autoencoder

- This idea comes from neural nets generative models;
- This framework enables representation learning (*encoding*) within a generating (*decoding*) model;

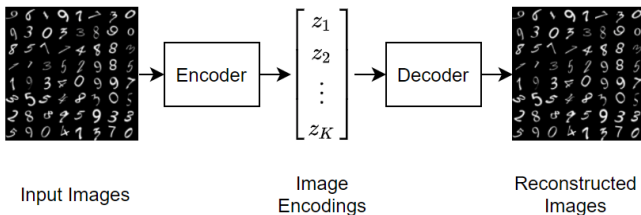


Figure 2: Source: fr.mathworks.com/

Gradient of the ELBO

- In general, the ELBO²

$$\text{ELBO}(\lambda, \theta) = \mathbb{E}_{q^\lambda} [\log p(\mathbf{Y}|Z, \theta)] - \text{KL} (q^\lambda(z) \parallel p(z)) ,$$

does not have an explicit form, because of the expectation.

- The KL term is not a problem in general;
- Stochastic gradient approach requires to compute

$$\nabla_{\lambda, \theta} \mathbb{E}_{q^\lambda} [\log p(\mathbf{Y}|Z, \theta)] .$$

- For a fixed λ , we have that:

$$\nabla_{\theta} \mathbb{E}_{q^\lambda} [\log p(\mathbf{Y}|Z, \theta)] = \mathbb{E}_{q^\lambda} [\nabla_{\theta} \log p(\mathbf{Y}|Z, \theta)] ,$$

which can be estimated by Monte Carlo, using samples from q^λ .

- For a fixed θ , however:

$$\nabla_{\lambda} \mathbb{E}_{q^\lambda} [\log p(\mathbf{Y}|Z, \theta)] \neq \mathbb{E}_{q^\lambda} [\nabla_{\lambda} \log p(\mathbf{Y}|Z, \theta)] ,$$

~~as the integration is w.r.t. q^λ .~~

²We now denote q^λ to highlight the dependence of the distribution in this parameter

$$\begin{aligned}\nabla_{\lambda} \mathbb{E}_{q^{\lambda}} [\log p(\mathbf{Y}|Z, \theta)] &= \nabla_{\lambda} \int_{\Omega_{\mathbf{Z}}} \log p(\mathbf{Y}|Z, \theta) q^{\lambda}(z) dz \\ &\neq \mathbb{E}_{q^{\lambda}} [\nabla_{\lambda} \log p(\mathbf{Y}|Z, \theta)] \quad ,\end{aligned}$$

- Nonetheless, using the fact $\nabla_{\lambda} q^{\lambda} = q^{\lambda} \times \nabla_{\lambda} \log q^{\lambda}$, we have that:

$$\begin{aligned}\nabla_{\lambda} \mathbb{E}_{q^{\lambda}} [\log p(\mathbf{Y}|Z, \theta)] &= \mathbb{E}_{q^{\lambda}} [\nabla_{\lambda} \log p(\mathbf{Y}|Z, \theta)] \\ &\quad + \mathbb{E}_{q^{\lambda}} [\log p(\mathbf{Y}|Z, \theta) \times \nabla_{\lambda} \log q^{\lambda}(\mathbf{Z})] \quad ,\end{aligned}$$

which can be estimated by Monte Carlo, using samples from q^{λ} .

- This direct Monte Carlo estimator can have high variance;
- Different variance reduction techniques can be used;
- A popular one is the reparametrization trick:

- Suppose that you can write

$$\mathbf{Z} = r(\varepsilon, \lambda) \quad ,$$

where ε is a random variable whose distribution does not depend on λ , and $r(\cdot, \lambda)$ is a known function

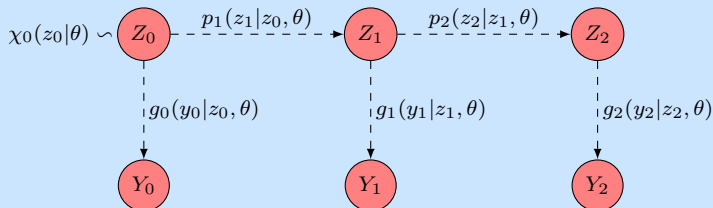
- For instance, if $\mathbf{Z} \sim \mathcal{N}(\lambda_1, \lambda_2^2)$, take $\varepsilon \sim \mathcal{N}(0, 1)$, and $r(\varepsilon, \lambda) = \lambda_1 + \lambda_2 \varepsilon$.
- In this case:

$$\nabla_{\lambda} \mathbb{E}_{q_{\lambda}} [\log p(\mathbf{Y}|Z, \theta)] = \mathbb{E}_{\varepsilon} [\nabla_{\lambda} \log p(\mathbf{Y}|r(\varepsilon, \lambda), \theta)] \quad ,$$

which can be estimated without bias, by Monte Carlo sampling, using samples from ε .

Variational inference and Hidden Markov models

Hidden Markov models



$$\begin{aligned}
 Z_0 &\sim \chi_0(z_0, \theta) \\
 Z_t | \{Z_{t-1} = z_{t-1}\} &\sim p_t(z_t | z_{t-1}, \theta) && \text{Prior distribution} \\
 Y_t | \{Z_t = z_t\} &\sim g_t(y_t | z_t, \theta) && \text{Observation model}
 \end{aligned}$$

In this context, the posterior distribution $Z_{0:n}|Y_{0:n}$ satisfies the following factorization:

$$p_{0:n}(z_{0:n}|Y_{0:n}, \theta) := p(z_n|Y_{0:n}, \theta) \prod_{t=1}^n p(z_{t-1}|z_t, y_{0:(t-1)}, \theta)$$

and has the structure of Markov chain.

Variational HMM, or VHMM

Variational HMM, or VHMM



- Idea (Campbell et al. (2021)), parameterize the variational with the same backward decomposition as the target distribution:

$$q_{0:n}^\lambda(z_{0:n}) = q_n^\lambda(z_n) \prod_{k=1}^n q_{t-1|t}^\lambda(z_{k-1}|z_k).$$

- How can we parameterize the distribution of each brick for efficient estimation?

- For each $t \geq 0$, $q_t^\lambda(z_t)$ aims at mimicking the filtering distribution (i.e. $Z_t|Y_{0:t}$);
- We impose that $q_t^\lambda(z_t)$ is the p.d.f. of a gaussian distribution with natural parameter η_t^λ defined in the following way:
 - $\eta_0 = \text{NeuralNet}(h_0|\lambda)$, where $\text{NeuralNet}(y_0|\lambda)$;
 - For $t > 0$, $\eta_t = \text{NeuralNet}(h_t|\lambda)$, where $\text{NeuralNet}(h_{t-1}, y_t|\lambda)$.
- This recursion (called a recurrent neural network), mimicks the predict-update steps of filtering recursions in HMM.

³Recall that $p(z_{t-1}|z_t, y_{0:t-1}, \theta) \propto p(z_{t-1}|y_{0:t-1}, \theta)p(z_t|z_t, \theta)$

- For each $t \geq 0$, $q_t^\lambda(z_t)$ aims at mimicking the filtering distribution (i.e. $Z_t|Y_{0:t}$);
- We impose that $q_t^\lambda(z_t)$ is the p.d.f. of a gaussian distribution with natural parameter η_t^λ defined in the following way:
 - $\eta_0 = \text{NeuralNet}(h_0|\lambda)$, where $\text{NeuralNet}(y_0|\lambda)$;
 - For $t > 0$, $\eta_t = \text{NeuralNet}(h_t|\lambda)$, where $\text{NeuralNet}(h_{t-1}, y_t|\lambda)$.
- This recursion (called a recurrent neural network), mimicks the predict-update steps of filtering recursions in HMM.
- Impose $q_{t-1|t}^\lambda(z_{k-1}|z_k)$ to be the p.d.f. of a gaussian distribution with parameter $\eta_{t-1|t}^\lambda$;
- To link³ $q_{t-1|t}^\lambda(z_{k-1}|z_k)$ to q_{t-1}^λ and z_t , we impose that:

$$\eta_{t-1|t}^\lambda = \eta_{t-1}^\lambda + \vec{\eta}_t^\lambda,$$

where $\vec{\eta}_t^\lambda = \text{NeuralNet}(z_t|\lambda)$ is a gaussian natural parameter.

³Recall that $p(z_{t-1}|z_t, y_{0:t-1}, \theta) \propto p(z_{t-1}|y_{0:t-1}, \theta)p(z_t|z_t, \theta)$

- This allows to learn representation of times series;
- This factorization is prone to online learning (i.e. learning of (λ, θ) seeing each y_t once);
- **Example model:** Chaotic recurrent neural network with dimension 20 for both X and Y

$$X_0 \sim \mathcal{N}(0, Q), X_t = X_{t-1} + \frac{\Delta}{\tau} (\gamma W \tanh(X_{t-1}) - X_{t-1}) + \eta, t \geq 1, \eta \sim \mathcal{N}(0, Q)$$
$$Y_t = X_t + \epsilon, \epsilon \sim St(2), t \geq 0,$$

Variational learning for HMM

- Online learning of the law of $X_{0:n}|Y_{0:n}$;
- Learnt on a train time series, used to represent a test time series.

Sequence	Smoothing RMSE	Filtering RMSE
Training	0.281	0.311
Eval	0.278 (± 0.01)	0.305 (± 0.014)

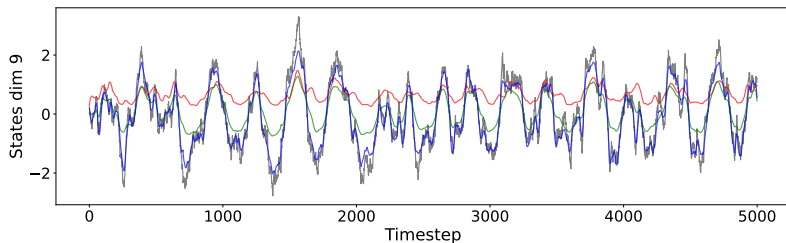


Figure 3: *Black:* True hidden signal, *Red, Green, Blue:* Estimated signal after processing of 1, 10000, 100000 observations

- Neural networks are useful tools to design non linear functions;
- Variational autoencoder is classical variational inference with an additional constraints;
- These ideas are useful bricks that can be implemented for many classical latent variable models.

- Blei, David M, Alp Kucukelbir, and Jon D McAuliffe. 2017. "Variational Inference: A Review for Statisticians." *Journal of the American Statistical Association*, no. just-accepted.
- Campbell, Andrew, Yuyang Shi, Thomas Rainforth, and Arnaud Doucet. 2021. "Online Variational Filtering and Parameter Learning." In *Advances in Neural Information Processing Systems*, 34:18633–45.
- Kingma, Diederik P, and Max Welling. 2013. "Auto-Encoding Variational Bayes." *arXiv Preprint arXiv:1312.6114*.
- Mallat, Stéphane. 2019. "L'apprentissage Par Réseaux de Neurones Profonds." Cours du collège de France, <https://www.college-de-france.fr/fr/agenda/cours/apprentissage-par-reseaux-de-neurones-profonds>.