



KUNGLIGA TEKNISKA HÖGSKOLAN

MACHINE LEARNING ADVANCED COURSE

Assignment 1B

Student :

Yohan PELLERIN

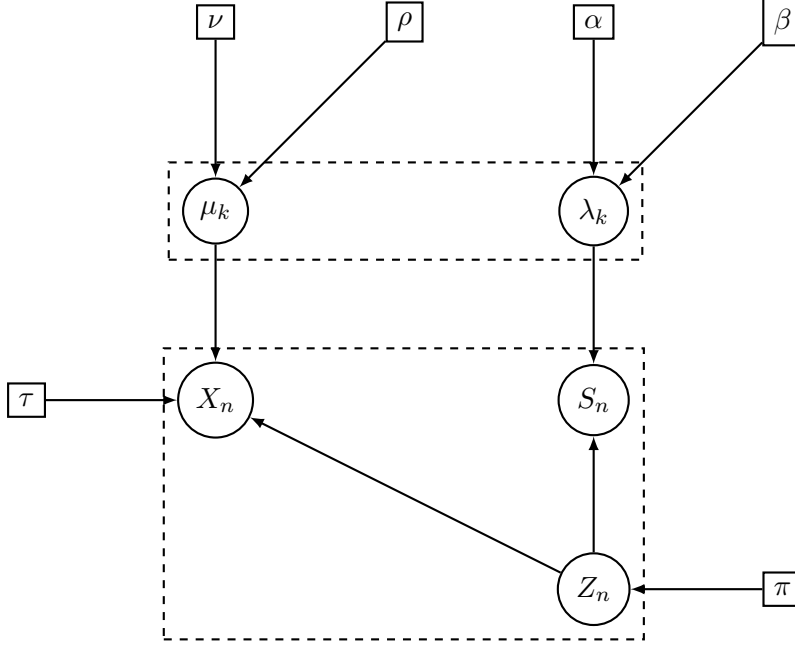
Course coordinator and teachers :

Jens LAGERGREN

December 10, 2023

1 CAVI for Earth quakes

Question 1.1.1



Question 1.1.2

Using conditional Independence, we can write :

$$\begin{aligned}
 p(X, S, Z, \lambda, \mu | \pi, \tau, \alpha, \beta, \nu, \rho) &= \prod_n [p(X_n | Z_n, \mu, \tau) p(S_n | \lambda, Z_n) p(Z_n | \pi)] \prod_k [p(\mu_k | \beta, \alpha) p(\lambda | \nu, \rho)] \\
 p(X, S, Z, \lambda, \mu | \pi, \tau, \alpha, \beta, \nu, \rho) &= \prod_n \left[p(Z_n | \pi) \prod_k [p(X_n | Z_n = k, \mu_k, \tau)^{\mathbb{I}_{Z_n=k}} p(S_n | \lambda, Z_n = k)^{\mathbb{I}_{Z_n=k}}] \right] * \\
 &\quad \prod_k [p(\mu_k | \beta, \alpha) p(\lambda | \nu, \rho)] \\
 \log p(X, S, Z, \lambda, \mu | \pi, \tau, \alpha, \beta, \nu, \rho) &= \sum_n \left[\log p(Z_n | \pi) + \sum_k \mathbb{I}_{Z_n=k} [\log p(X_n | Z_n = k, \mu_k, \tau) \right. \\
 &\quad \left. + \log p(S_n | \lambda, Z_n = k)] \right] + \sum_k \log p(\mu_k | \beta, \alpha) + \log p(\lambda | \nu, \rho)
 \end{aligned}$$

Question 1.1.3

The CAVI update equations are the followings.

$$\begin{aligned}
 \log q^*(\mu_i) &\stackrel{+}{=} E_{-\mu_i} [P(X, S, Z, \lambda, \mu)] \\
 \log q^*(Z_i) &\stackrel{+}{=} E_{-Z_i} [P(X, S, Z, \lambda, \mu)] \\
 \log q^*(\lambda_i) &\stackrel{+}{=} E_{-\lambda_i} [P(X, S, Z, \lambda, \mu)]
 \end{aligned}$$

The first equation

Let's focus on the first equation, using the expression from the previous question, and only keeping the terms that depend on μ_i , we get :

$$\log q^*(\mu_i) \stackrel{+}{=} \sum_{n,k} E_{-\mu_i} [\mathbb{I}_{Z_n=k} \log p(X_n|Z_n=k, \mu_k, \tau)] + \sum_k E_{-\mu_i} [\mathbb{I}_{Z_n=k} \log p(\mu_k|\nu, \rho)]$$

Now, let's turn our attention to the first expectation. We know $X_n|Z_n=k, \mu, \tau$ follows a 2D Normal($\mu_i, \tau.I$)

If $k=i$, $\log p(X_n|Z_n=k, \mu_i, \tau) \stackrel{+}{=} -\frac{1}{2\tau}(\mu_i - X_n)^T(\mu_i - X_n)$

But if $k \neq i$ this term doesn't depend on μ_i

$$\begin{aligned} \sum_{k,n} E_{-\mu_i} [\mathbb{I}_{Z_n=k} \log p(X_n|Z_n=k, \mu, \tau)] &\stackrel{+}{=} -\frac{1}{2\tau} \sum_n (\mu_i - X_n)^T(\mu_i - X_n) q(Z_n=i) \\ &\stackrel{+}{=} -\frac{1}{2\tau} \left[\mu_i^T \mu_i - \left(\sum_n q(Z_n=i) X_n^T \right) \mu_i - \mu_i^T \sum_n q(Z_n=i) X_n \right] \end{aligned}$$

For the second expectation, in a similar way, as we know $\lambda_k|\nu, \rho$ follows a 2D Normal($\nu, \rho.I$) we get :

$$\begin{aligned} \sum_k E_{-\mu_i} [\log p(\mu_k|\nu, \rho)] &\stackrel{+}{=} -\frac{1}{2\rho}(\mu_i - \nu)^T(\mu_i - \nu) \\ &\stackrel{+}{=} -\frac{1}{2\rho} [\mu_i^T \mu_i - \nu^T \mu_i - \mu_i^T \nu] \end{aligned}$$

Then, we rebuild the square and we get :

$$\log q^*(\mu_i) \stackrel{+}{=} -\frac{1}{2\tau_i^*}(\mu_i - \mu_i^*)^T(\mu_i - \mu_i^*)$$

with, $\frac{1}{\tau^*} = \frac{1}{\tau} + \frac{1}{\rho}$ and $\mu_i^* = \frac{\sum_n q(Z_n=i) X_n + \frac{\nu}{\rho}}{\frac{1}{\tau} + \frac{1}{\rho}}$

Then, we can say that the distribution (over q) of μ_i is a 2D Normal($\mu_i^*, \tau^*.I$)

The second equation

Let's focus on the second equation, using the expression from the previous question, and only keeping the terms that depend on Z_i , we get :

$$\log q^*(Z_i) \stackrel{+}{=} \log p(Z_i|\pi) + \sum_k \mathbb{I}_{Z_i=k} E_{-Z_i} [\log p(X_i|Z_i=k, \mu_k, \tau) + \log p(S_i|\lambda_k, Z_i=k)]$$

As $Z_i|\pi$ is a Categorical, $\log p(Z_i|\pi) = \sum_k \mathbb{I}_{Z_i=k} \log(\pi_k)$

As $X_i|Z_i=k, \mu, \tau$ follows a 2D Normal($\mu_k, \tau.I$)

$$\begin{aligned} \log p(X_i|Z_i=k, \mu_k, \tau) &= -\log(\tau^2 2\pi) - \frac{1}{2\tau}(\mu_k - X_i)^T(\mu_k - X_i) \\ E_{-Z_i} [\log p(X_i|Z_i=k, \mu, \tau)] &= -\log(\tau^2 2\pi) - \frac{1}{2\tau} (E_{-Z_i} [\mu_k^T \mu_k] - X_i^T E_{-Z_i} [\mu_k] - E_{-Z_i} [\mu_k^T] X_i + X_i^T X_i) \\ E_{-Z_i} [\log p(X_i|Z_i=k, \mu, \tau)] &= -\log(\tau^2 2\pi) - \frac{1}{2\tau} (E_{\mu_k} [\mu_k^T \mu_k] - X_i^T E_{\mu_k} [\mu_k] - E_{\mu_k} [\mu_k^T] X_i + X_i^T X_i) \end{aligned}$$

Moreover, as we know that μ_k follows a $2DNormal(\mu_k^*, \tau^*.I)$,

$$E_{\mu_k}[\mu_k] = \mu_k^* \text{ and } E_{\mu_k}[\mu_k^T] = \mu_k^{*T}$$

As $\tau^*.I$ is diagonal matrix, the first dimension and the second aren't correlated.

$$\begin{aligned} \mu_k &= \begin{bmatrix} \mu_{k1} \\ \mu_{k2} \end{bmatrix} \\ E_{\mu_k} \begin{bmatrix} \mu_k^T \mu_k \end{bmatrix} &= E_{\mu_{k1}}[\mu_{k1}^2] + E_{\mu_{k2}}[\mu_{k2}^2] = Var(\mu_{k2}) + E_{\mu_{k2}}[\mu_{k2}^2] + Var(\mu_{k1}) + E_{\mu_{k1}}[\mu_{k1}^2] \\ E_{\mu_k} \begin{bmatrix} \mu_k^T \mu_k \end{bmatrix} &= 2 * \tau^* + \mu_k^{*T} \mu_k^* \end{aligned}$$

As $S_i|\lambda, Z_i = k$ follows a $Poisson(\lambda_k)$,

$$E_{-Z_i}[\log p(S_i|\lambda, Z_i = k)] = S_i E_{-Z_i}[\log(\lambda_k)] - \log(S_i!) - E_{-Z_i}[\lambda_k]$$

$$E_{-Z_i}[\log p(S_i|\lambda, Z_i = k)] = S_i E_{\lambda_k}[\log(\lambda_k)] - \log(S_i!) - E_{\lambda_k}[\lambda_k]$$

Moreover, we can also calculate the two expectations because the last equation will show that λ_k follows a Gamma distribution with a_k^* and b_k^* as parameter.

$$E_{\lambda_k}[\log(\lambda_k)] = \psi(a_k^*) - \log(b_k^*) \text{ and } E_{\lambda_k}[\lambda_k] = \frac{a_k^*}{b_k^*}$$

Then, after regrouping all the terms we get :

$$\begin{aligned} \log q^*(Z_i) &\stackrel{+}{=} \sum_k \mathbb{I}_{Z_i=k} \left[\log(\pi_k) + \log(\lambda_k) - \log(S_i!) - E_{\lambda_k}[\lambda_k] - \log(\tau^2 2\pi) \right. \\ &\quad \left. - \frac{1}{2\tau} (E_{\mu_k}[\mu_k^T \mu_k] - X_i^T E_{\mu_k}[\mu_k] - E_{\mu_k}[\mu_k^T] X_i + X_i^T X_i) \right] \end{aligned}$$

Then, we can say that the distribution (over q) of Z_i is Categorical parameterized by π_i^* , where

$$\begin{aligned} \log(\pi_{ki}^*) &= \log(\pi_k) + \log(\lambda_k) - \log(S_i!) - E_{\lambda_k}[\lambda_k] - \log(\tau^2 2\pi) \\ &\quad - \frac{1}{2\tau} (E_{\mu_k}[\mu_k^T \mu_k] - X_i^T E_{\mu_k}[\mu_k] - E_{\mu_k}[\mu_k^T] X_i + X_i^T X_i) \end{aligned}$$

The last equation

Let's focus on the second equation, using the expression from the previous question, and only keeping the terms that depend on λ_i , we get :

$$\log q^*(\lambda_i) \stackrel{+}{=} \log p(\lambda_i|\alpha, \beta) + \sum_n E_{-\lambda_i}[\mathbb{I}_{Z_n=i} \log p(S_i|\lambda_i, Z_n = i)]$$

As $\lambda_i|\alpha, \beta$ follows a gamma distribution with parameters α and β

$$\log p(\lambda_i|\alpha, \beta) \stackrel{+}{=} (\alpha - 1) \log(\lambda_i) - \beta \lambda_i$$

For the second term, as $S_i|\lambda_i, Z_n = i$ follows a $Poisson(\lambda_i)$:

$$E_{-\lambda_i}[\mathbb{I}_{Z_n=i} \log p(S_i|\lambda_i, Z_n = i)] \stackrel{+}{=} q(Z_n = i)(S_n \log(\lambda_i) - \lambda_i)$$

Then regrouping the two terms, we get:

$$\log q^*(\lambda_i) \stackrel{+}{=} (a_i^* - 1) \log(\lambda_i) - \lambda_i (b_i^*)$$

With , $a_i^* = \sum_n [q(Z_n = i) S_n] + \alpha$

$b_i^* = \beta + \sum_n q(Z_n = i)$ Then, we can say that the distribution (over q) of λ_i is a gamma distribution parameterized by a_i^* and b_i^*

2 Reparametrization and the score function

Question 1.3.4

The expression of the ELBO is the following

$$L(\phi) = \mathbb{E}_{z \sim q} [\log p(x|z) + \log p(z) - \log q_\phi(z|x)]$$

Using the reparameterization trick, we can express a sample z from a parametric distribution $q_\phi(z)$ as a deterministic function of a random variable ϵ with some fixed distribution and the parameters ϕ of q_ϕ , i.e., $z = t(\epsilon, \phi)$. Under such a parameterization of z , we can decompose the total derivative (TD) of the integrand of the ELBO.

$$\begin{aligned}\hat{\nabla}_{TD}(\epsilon, \phi) &= \nabla_\phi [\log p(x|z) + \log p(z) - \log q_\phi(z|x)] \\ &= \nabla_\phi [\log p(z|x) + \log p(x) - \log q_\phi(z|x)] \\ &= \nabla_z [\log p(z|x) - \log q_\phi(z|x)] \nabla_\phi t(\epsilon, \phi) - \nabla_\phi \log q_\phi(z|x),\end{aligned}$$

Where, $\nabla_\phi \log q_\phi(z|x)$ is the score function, and $\nabla_z [\log p(z|x) - \log q_\phi(z|x)] \nabla_\phi t(\epsilon, \phi)$ is the path derivative.

Question 1.3.5

$$\begin{aligned}\mathbb{E}_{z \sim q} [\nabla_\phi \log q_\phi(z|x)] &= \int_z (\nabla_\phi \log q_\phi(z|x)) q_\phi(z|x) dz = \int_z \nabla_\phi q_\phi(z|x) dz \\ &= \nabla_\phi \left(\int_z q_\phi(z|x) dz \right) = \nabla_\phi(1) = 0\end{aligned}$$

The expectation of the score function is zero.

Question 1.3.6

The score function term has a high variance, then by removing this term, (because the expectation of the score function is zero) we can keep an unbiased estimator of the true gradient which is needed for stochastic gradient descent to converge.

Question 1.3.7

The score function can act in some situations as a control variate.

3 Reparameterization of common distributions

Question 1.4.8

An exponential distribution as tractable inverse CDF. In this case, let ϵ follows a $U(0,1)$, and let $g_\phi(\epsilon, x)$ be the inverse CDF of $q_\phi(z|x)$. The inverse of an exponential distribution is $f(y) = -\frac{\log(1-y)}{\lambda}$

Then $g_\phi(\epsilon, x) = -\frac{\log(1-\epsilon)}{\lambda}$, Where ϵ follows a $U(0,1)$ and λ is the parameter. See the implementation in Q1 of the notebook.

Question 1.4.8

To reparameterize the categorical distribution, we approximate it by the Gumbel-Softmax distribution. First, we sample the Gumbel(0,1) distribution using inverse transform sampling by drawing u follows a Uniform(0, 1) and computing $g = -\log(-\log(u))$. Then, we add $\log a$ the parameter of the categorical. Then, we apply the softmax function which then allows for differentiation during training, and using the argmax function for evaluation. See the implementation in Q2 of the notebook.

Appendix

Reparameterization of common distributions

We will work with Torch throughout this notebook.

```
In [1]: import torch
from torch.distributions import Beta, Exponential, Categorical #, ... import the distributions you need here
from torch.nn import functional as F
```

A helper function to visualize the generated samples:

```
In [2]: import matplotlib.pyplot as plt
def compare_samples(samples_1, samples_2, bins=100, range=None):
    fig = plt.figure()
    if range is not None:
        plt.hist(samples_1, bins=bins, range=range)
        plt.hist(samples_2, bins=bins, range=range)
    else:
        plt.hist(samples_1, bins=bins)
        plt.hist(samples_2, bins=bins)
    plt.xlabel('value')
    plt.ylabel('number of samples')
    plt.legend(['direct', 'via reparameterization'])
    plt.show()
```

Q1. Exponential Distribution

Below write a function that generates N samples from $Exp(\lambda)$.

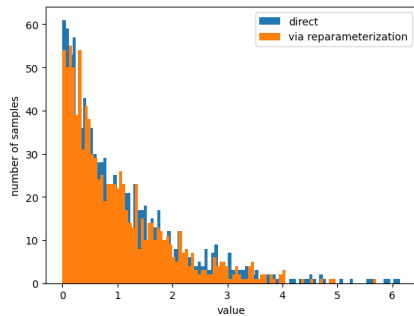
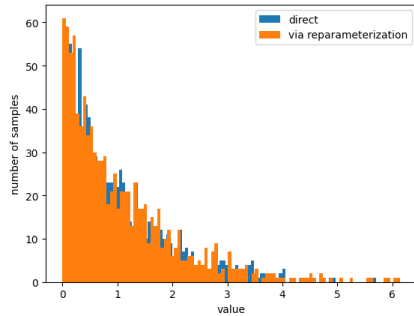
```
In [3]: def exp_sampler(1, N):
# insert your code
rates = torch.full((N,), float(1))
samples = Exponential(rates).sample()
return samples # should be N-by-1
```

Now, implement the reparameterization trick:

```
In [4]: def exp_reparametrize(1,N):
# this function should return N samples via reparameterization,
# insert your code
e = torch.rand(N)
samples = -torch.log(1-e)/1
return samples
```

Generate samples for $\lambda = 1$ and compare:

```
In [5]: l = 1 #lambda
N = 1000
direct_samples = exp_sampler(1, N)
reparametrized_samples = exp_reparametrize(1, N)
compare_samples(direct_samples, reparametrized_samples)
compare_samples(reparametrized_samples, direct_samples)
```



Q2. Categorical Distribution

Below write a function that generates N samples from Categorical(a), where $a = [a_0, a_1, a_2, a_3]$.

```
In [6]: def categorical_sampler(a, N):
# insert your code
rates = a.repeat(N,1)
samples = Categorical(rates).sample()
return samples # should be N-by-1
```

Now write a function that generates samples from Categorical(a) via reparameterization:

```
In [7]: # Hint: approximate the Categorical distribution with the Gumbel-Softmax distribution
def categorical_reparametrize(a, N, temp=0.1, eps=1e-20): # temp and eps are hyperparameters for Gumbel-Softmax
# insert your code
e = torch.rand(N, a.size(0))
g = - torch.log(-torch.log(e+eps)*eps)

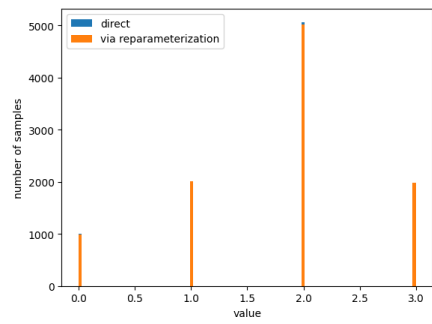
logits = (torch.log(a) + g) / temp
softmax_output = torch.nn.functional.softmax(logits, dim=1)

# Apply the reparameterization trick
samples = torch.argmax(softmax_output, dim=1)

return samples # make sure that your implementation allows the gradient to backpropagate
```

Generate samples when $a = [0.1, 0.2, 0.5, 0.2]$ and visualize them:

```
In [8]: a = torch.tensor([0.1, 0.2, 0.5, 0.2])
N = 1000
direct_samples = categorical_sampler(a, N)
reparametrized_samples = categorical_reparametrize(a, N, temp=0.1, eps=1e-20)
compare_samples(direct_samples, reparametrized_samples)
```



In []:

