# SVAE objective

July 26, 2017

## 1  Definition of the model

Define the following model:

$$z_n|\theta \sim \text{discrete}(\pi^0)$$
$$x_n|z_n, \theta \sim \mathcal{N}(\mu^0_{z_n}, \Sigma^0_{z_n})$$
$$\text{where} \quad \theta = (\pi^0, \{\mu^0_k, \Sigma^0_k\}_{k\in[0,K]})$$
$$\text{and} \quad y_n|x_n, \gamma \sim \text{Bernoulli}(f_\gamma(x_n), \text{M})$$

where $f_\gamma(x_n)$ is the output of a neural network parametrized by $\gamma$ to the input $x_n$, K is the number of mixtures (10 for MNIST) and M is the dimension of the observable variables $y_n$ ($28 \times 28$ for MNIST). Now using the notation under the exponential family form, we can write:

$$p_{z|\theta}(z_n) = \exp\left(\langle \eta^0_z(\theta), t_z(z_n)\rangle - \log\left(Z_z(\eta^0_z(\theta))\right)\right)$$

where:

$$\eta^0_z(\theta)^\top = \left[\ln(\pi^0_1), \ldots, \ln(\pi^0_K)\right]$$
$$t_z(z_n)^\top = \left[\delta(z=1), \ldots, \delta(z=K)\right]$$
$$\log\left(Z_z(\eta)\right) = 0 \quad \forall \eta$$

For the mixture of gaussians part $p_{x|z,\theta}$, we can write:

$$p_{x|z,\theta}(x) = \exp\left(\langle \eta^0_x(\theta, z), t_x(x)\rangle - \log\left(Z_x(\eta^0_x(\theta, z))\right)\right)$$

where: $t_x(x)^\top = \left[x, vec(xx^\top)^\top\right]$ and $\forall k \in [1, K], \eta^0_x(\theta, z=k)^\top = \left[\Sigma^{0-1}_k \mu^0_k, vec(-\frac{1}{2}\Sigma^{0-1}_k)^\top\right]$.

Using the conjugacy of the model, we can rewrite $p_{x|z,\theta}$ as follow:

$$
\begin{aligned}
p_{x|z,\theta}(x) &= \frac{1}{\sqrt{|2\pi\Sigma_z^0|}} \exp\Big( -\frac{1}{2}(x-\mu_z^0)^\top \Sigma_z^{0-1}(x-\mu_z^0) \Big) \\
&= \prod_{k=1}^{K} \Big[ \frac{1}{\sqrt{|2\pi\Sigma_k^0|}} \exp\Big( -\frac{1}{2}(x-\mu_k^0)^\top \Sigma_k^{0-1}(x-\mu_k^0) \Big) \Big]^{\delta(z=k)} \\
&= \exp\Big( \sum_{k=1}^{K} \delta(z=k) \Big[ -\frac{1}{2}\ln\left(|2\pi\Sigma_k^0|\right) - \frac{1}{2}(x-\mu_k^0)^\top \Sigma_k^{0-1}(x-\mu_k^0) \Big] \Big) \\
&= \exp\Big( \sum_{k=1}^{K} \delta(z=k) \Big[ -\frac{1}{2}\ln\left(|2\pi\Sigma_k^0|\right) - \frac{1}{2}\mu_k^{0\top}\Sigma_k^{0-1}\mu_k^0 + x^\top\Sigma_k^{0-1}\mu_k^0 - \frac{1}{2}x^\top\Sigma_k^{0-1}x \Big] \Big) \\
&= \exp\Big( \sum_{k=1}^{K} \delta(z=k) \Big[ -\frac{1}{2}\ln\left(|2\pi\Sigma_k^0|\right) - \frac{1}{2}\mu_k^{0\top}\Sigma_k^{0-1}\mu_k^0 + \eta_x^0(\theta, z=k)^\top t_x(x) \Big] \Big) \\
&= \exp\Big( \sum_{k=1}^{K} \delta(z=k) \Big[ \eta_x^0(\theta)[k]^\top \begin{bmatrix} t_x(x) \\ 1 \end{bmatrix} \Big] \Big)
\end{aligned}
$$

and thus, we have:

$$
p_{x|z,\theta}(x) = \exp\Big( \Big\langle t_z(z), \eta_x^0(\theta)^\top \begin{bmatrix} t_x(x) \\ 1 \end{bmatrix} \Big\rangle \Big) \tag{1}
$$

where $\eta_x^0(\theta) = \big[\eta_x^0(\theta)[1], \ldots, \eta_x^0(\theta)[K]\big]$ is a $\big(N(N+1)+1\big) \times K$ matrix with its columns $\eta_x^0(\theta)[k] \in \mathbb{R}^{N+N^2+1}$, $\forall k \in [1,K]$, given by:

$$
\eta_x^0(\theta, z=k)^\top = \Big[ \Sigma_k^{0-1}\mu_k^0, vec(-\frac{1}{2}\Sigma_k^{0-1})^\top, -\frac{1}{2}\ln\left(|2\pi\Sigma_k^0|\right) - \frac{1}{2}\mu_k^{0\top}\Sigma_k^{0-1}\mu_k^0 \Big]
$$

# 2 SVAE objective

## 2.1 SVAE objective Definition

Taking the respective variational factors in the corresponding exponential families, we define:

$$
q_z(z_n) = \exp\Big( \langle \eta_z, t_z(z_n) \rangle - \log\big(Z_z(\eta_z)\big) \Big)
$$

$$
q_x(x_n) = \exp\Big( \langle \eta_x, t_x(x_n) \rangle - \log\big(Z_x(\eta_x)\big) \Big)
$$

We can now define the mean field objective of the problem $\mathcal{L}$:

$$
\mathcal{L}(\eta_z, \eta_x) = \mathbb{E}_{q_z q_x}\Big[ \log\Big( \frac{p_{z|\theta}(z_n)p_{x|z,\theta}(x_n)p_{y|x,\gamma}(y_n)}{q_z(z_n)q_x(x_n)} \Big) \Big]
$$

In the $SVAE$ algorithm, we express $\eta_x$ and $\eta_z$ as function of the remaining parameters using a surrogate objective $\widehat{\mathcal{L}}$ and introducing a recognition network $r_\phi$ defined as follow:

$$
\widehat{\mathcal{L}}(\eta_z, \eta_x, \phi) = \mathbb{E}_{q_z q_x}\Big[ \log\Big( \frac{p_{z|\theta}(z_n)p_{x|z,\theta}(x_n)\exp(\psi(x_n;\phi))}{q_z(z_n)q_x(x_n)} \Big) \Big]
$$

where $\psi(x_n;\phi) = \langle r_\phi(y_n), t_x(x_n)\rangle$ and $r_\phi(y_n)$ a the output of a neural network parametrized by $\phi$, the recognition network, to the input $y_n$.

We now partially optimize the surrogate objective $\widehat{\mathcal{L}}$ w.r.t $\eta_x$ and $\eta_z$ defining $\eta_x^*$ and $\eta_z^*$:

$$\eta_x^*(\phi), \eta_z^*(\phi) = argmax_{\eta_x,\eta_z}\widehat{\mathcal{L}}(\eta_z, \eta_x, \phi)$$

And finally, the SVAE objective, $\mathcal{L}_{\text{SVAE}}$ as:

$$\mathcal{L}_{\text{SVAE}}(\gamma, \theta, \phi) = \mathcal{L}(\eta_z^*(\phi), \eta_x^*(\phi))$$

Expanding the expression of $\mathcal{L}_{\text{SVAE}}$, we can write:

$$\mathcal{L}_{\text{SVAE}}(\gamma, \theta, \phi) = \mathbb{E}_{q_x^*}\left[\log\left(p_{y|x,\gamma}(y_n)\right)\right] - KL(q_x^* q_z^* || p_{z|\theta}p_{x|z,\theta}) \tag{2}$$

## 2.2 SVAE objective derivation

The first term of the *R.H.S* of the equation can be computed using the reparametrization trick, sampling $\hat{x}$ from respectively $q_x^*$ and using the following approximation:

$$\mathbb{E}_{q_x^*}\left[\log\left(p_{y|x,\gamma}(y_n)\right)\right] \approx \log\left(p_{y|x,\gamma}(y_n|\hat{x},\gamma)\right) \tag{3}$$

where $\hat{x} \sim q_x^*$

We will denote the last term of the *R.H.S* as the *local meanfield* term of the objective. Using the definition of the *KL*, we can expand the *local meanfield* as:

$$KL(q_z^* q_x^* || p_{z|\theta}p_{x|z,\theta}) = KL(q_z^* || p_{z|\theta}) + \mathbb{E}_{q_z^*}\left[KL(q_x^* || p_{x|z,\theta})\right]$$

Using the result for the *KL* in the exponential family case, we have:

$$
\begin{aligned}
KL(q_z^* || p_{z|\theta}) = {} & \left\langle \eta_z^*(\phi) - \eta_z^0(\theta), \mathbb{E}_{q_z^*}[t_z(z)]\right\rangle \\
& - \left(\log\left(Z_z(\eta_z^*(\phi))\right) - \log\left(Z_z(\eta_z^0(\theta))\right)\right)
\end{aligned}
\tag{4}
$$

and

$$
\begin{aligned}
\mathbb{E}_{q_z^*}\left[KL(q_x^* || p_{x|z,\theta})\right] = {} & \left\langle \eta_x^*(\phi) - \mathbb{E}_{q_z^*}[\eta_x^0(\theta, z)], \mathbb{E}_{q_x^*}[t_x(x)]\right\rangle \\
& - \left(\log\left(Z_x(\eta_x^*(\phi))\right) - \mathbb{E}_{q_z^*}[\log\left(Z_x(\eta_x^0(\theta, z))\right)]\right)
\end{aligned}
\tag{5}
$$

Thus, we can rewrite our *local meanfiled* as:

$$
\begin{aligned}
KL(q_z^* q_x^* || p_{z|\theta}p_{x|z,\theta}) = {} & \left[\left\langle \eta_x^*(\phi) - \mathbb{E}_{q_z^*}[\eta_x^0(\theta, z)], \mathbb{E}_{q_x^*}[t_x(x)]\right\rangle - \log\left(Z_x(\eta_x^*(\phi))\right)\right] \\
& + \left[\left\langle \eta_z^*(\phi) - \eta_z^0(\theta), \mathbb{E}_{q_z^*}[t_z(z)]\right\rangle + \mathbb{E}_{q_z^*}[\log\left(Z_x(\eta_x^0(\theta, z))\right)] \right. \\
& \left. - \left(\log\left(Z_z(\eta_z^*(\phi))\right) - \log\left(Z_z(\eta_z^0(\theta))\right)\right)\right]
\end{aligned}
\tag{6}
$$

# 3   Block coordinate ascent algorithm

So it remains to compute the partial optimum of our surrogate objective, $\eta_x^*(\phi)$ and $\eta_z^*(\phi)$ as functions of the other parameters. Using the classical result for meanfield, we have, omitting everything constant *w.r.t* $x$ in the $\propto$ :

$$q_x^*(x; \theta, \phi) \propto \exp\left(\mathbb{E}_{q_z^*}[\log\left(p_{z|\theta}(z)p_{x|z,\theta}(x)\exp(\psi(x;\phi))\right)]\right)$$

$$\propto \exp\left(\mathbb{E}_{q_z^*}[\langle\eta_x^0(\theta, z)], t_x(x)\rangle + \langle r_\phi(y), t_x(x)\rangle]\right)$$

$$\propto \exp\left(\langle\mathbb{E}_{q_z^*}[\eta_x^0(\theta, z)] + r_\phi(y), t_x(x)\rangle\right)$$

and thus:

$$\eta_x^*(\phi) = \mathbb{E}_{q_z^*}[\eta_x^0(\theta, z)] + r_\phi(y) \tag{7}$$

In the same way, we can express $\eta_z^*(\phi)$ as a function of $\theta$ and $\phi$, omitting everything constant *w.r.t* $x$ in the $\propto$ :

$$q_z^*(z; \theta, \phi) \propto \exp\left(\mathbb{E}_{q_x^*}[\log\left(p_{z|\theta}(z)p_{x|z,\theta}(x)\exp(\psi(x;\phi))\right)]\right)$$

Using equation (1), we have, keeping only what depends of $z$:

$$q_z^*(z; \theta, \phi) \propto \exp\left(\mathbb{E}_{q_x^*}\left[\langle\eta_z^0(\theta)], t_z(z)\rangle + \langle t_z(z), \eta_x^0(\theta)^\top \begin{bmatrix} t_x(x) \\ 1 \end{bmatrix}\rangle\right]\right)$$

$$\propto \exp\left(\langle\mathbb{E}_{q_x^*}[\eta_x^0(\theta)^\top\left(t_x(x), x\right)] + \eta_z^0(\theta), t_z(z)\rangle\right)$$

and thus:

$$\eta_z^*(\phi) = \mathbb{E}_{q_x^*}\left[\eta_x^0(\theta)^\top \begin{bmatrix} t_x(x) \\ 1 \end{bmatrix}\right] + \eta_z^0(\theta) \tag{8}$$

We can inject the expressions (7) and (8) in (6) to get:

$$KL(q_z^* q_x^* || p_{z|\theta} p_{x|z,\theta}) = \left[\left\langle r_\phi(y), \mathbb{E}_{q_x^*}[t_x(x)]\right\rangle - \log\left(Z_x(\eta_x^*(\phi))\right)\right]$$

$$+ \left[\left\langle\mathbb{E}_{q_z^*}[\eta_x^0(\theta, z)], \mathbb{E}_{q_x^*}[t_x(x)]\right\rangle \right. \tag{9}$$

$$\left. - \left(\log\left(Z_z(\eta_z^*(\phi))\right) - \log\left(Z_z(\eta_z^0(\theta))\right)\right)\right]$$

We denote the first term of the *RHS* of (9) the *gaussian KL* and the second term of the *RHS* of (9) the *label KL*:

$$gaussian\_kl = \left\langle r_\phi(y), \mathbb{E}_{q_x^*}[t_x(x)]\right\rangle - \log\left(Z_x(\eta_x^*(\phi))\right) \tag{10}$$

$$label\_kl = \left\langle\mathbb{E}_{q_z^*}[\eta_x^0(\theta, z)], \mathbb{E}_{q_x^*}[t_x(x)]\right\rangle - \left(\log\left(Z_z(\eta_z^*(\phi))\right) - \log\left(Z_z(\eta_z^0(\theta))\right)\right) \tag{11}$$

We can then compute the local KL and perform the block ascent algorithm using (10) and (11).

# 4   PSEUDO CODE

---

**Algorithm 1** Implementation of SVAE

---

1: **function** SVAE
2:  Initialization of the parameters
3:   Mean parameters $\theta = (\pi^0, \{\mu_k^0, \Sigma_k^0\}_{k \in [0,K]})$
4:   Recognition network: $\gamma$
5:   Generative network: $\phi$
6:
7:  Process mean parameters
8:   $\psi \leftarrow r_\phi(y)$            $\triangleright$ Get the node potential from the recognition network
9:   $\eta_z^0 \leftarrow natpar(\pi^0)$           $\triangleright$ Get natural parameters for Discrete distribution
10:   $\eta_x^0 \leftarrow natpar(\{\mu_k^0, \Sigma_k^0\}_k)$      $\triangleright$ Get natural parameters for Gaussian distribution
11:
12:   $(\eta_x^*, \bar{t_x}, \eta_z^*, \bar{t_z}) \leftarrow$ **FIXEDmeanfield**$(\psi, \eta_z^0, \eta_x^0)$ $\triangleright$ Coordinate block ascent algorithm
13:   $\mathrm{KL}^{local} \leftarrow$ **LOCALmeanfield**$(\psi, \eta_z^0, \eta_x^0, \bar{t_z}, \bar{t_x})$        $\triangleright$ Compute local meanfield
14:   $\hat{x} \sim q_x^*$                    $\triangleright$ Sample mean parameter of observations
15:   loglike $\leftarrow$ **Xentropy**$(y, f_\gamma(\hat{x}))$           $\triangleright$ Estimate loglikelihood term
16:  **return** loglike$-\mathrm{KL}^{local}$

---