VAE with a VampPrior

Jakub Tomczak, Max Welling

Tübingen, 22nd of March 2018

Tomczak, J. M., & Welling, M. (2018). VAE with a VampPrior. AISTATS 2018 (oral presentation)

Modeling in a high-dimensional space is difficult.

Modeling in a high-dimensional space is difficult.





Modeling in a high-dimensional space is difficult.





Modeling in a high-dimensional space is difficult.

 \rightarrow modeling all dependencies among pixels.

$$p(x) = \prod_{d=1}^{c} \psi_c(x_c)$$

Modeling in a high-dimensional space is difficult.

 \rightarrow modeling all dependencies among pixels.



Modeling in a high-dimensional space is difficult.

 \rightarrow modeling all dependencies among pixels.



A possible **solution**? →**Models with latent variables**

Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ \mathrm{d}\mathbf{z}$$



Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$

First sample z.
Second sample x for given z.



Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$

First sample z.
Second sample x for given z.



Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$

How to calculate this integral?

If $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$ and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$, then we get **Factor Analysis**.

What if we take a **non-linear transformation** of **z**? →an infinite mixture of Gaussians



Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ \mathrm{d}\mathbf{z}$$

If $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$ and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$, then we get Factor Analysis.

What if we take a **non-linear transformation** of **z**? →**an infinite mixture of Gaussians**



Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ \mathrm{d}\mathbf{z}$$

If
$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$$
 and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$,
then we get **Factor Analysis**.

What if we take a **non-linear transformation** of **z**? →<mark>an infinite mixture of Gaussians</mark>



Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ \mathrm{d}\mathbf{z}$$

If $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$ and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$, then we get Factor Analysis.

What if we take a **non-linear transformation** of **z**? →an infinite mixture of Gaussians



Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ \mathrm{d}\mathbf{z}$$

If
$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$$
 and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$,
then we get Factor Analysis.

What if we take a **non-linear transformation** of z? \rightarrow an infinite mixture of Gaussians

MacKay, D. J., & Gibbs, M. N. (1999). Density networks. Statistics and neural networks: advances at the interface. Oxford University Press, Oxford, 129-144.

Latent variable model:

$$p(\mathbf{x}) = \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ \mathrm{d}\mathbf{z}$$

If
$$p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{W}\mathbf{z} + \mathbf{b}, \Psi)$$
 and $p_{\lambda}(\mathbf{z}) = \mathcal{N}(\mu_0, \Sigma_0)$,
then we get Factor Analysis.

What if we take a **non-linear transformation** of z? \rightarrow an infinite mixture of Gaussians



Not scalable...

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$
$$= \log \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$
$$\geq \int q_{\phi}(\mathbf{z}|\mathbf{x}) \ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ d\mathbf{z}$$
$$= \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathrm{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\lambda}(\mathbf{z})]$$

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z} \qquad \text{Variational posterior}$$
$$= \log \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$
$$\geq \int q_{\phi}(\mathbf{z}|\mathbf{x}) \ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ d\mathbf{z}$$
$$= \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathrm{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\lambda}(\mathbf{z})]$$

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$

$$= \underbrace{\log} \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$
Jensen's inequality
$$\stackrel{\geq}{=} \int q_{\phi}(\mathbf{z}|\mathbf{x}) \underbrace{\log} \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ d\mathbf{z}$$

$$= \mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] - \mathrm{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\lambda}(\mathbf{z})]$$

$$\log p(\mathbf{x}) = \log \int p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$
$$= \log \int \frac{q_{\phi}(\mathbf{z}|\mathbf{x})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z}) \ d\mathbf{z}$$
$$\geq \int q_{\phi}(\mathbf{z}|\mathbf{x}) \ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z})}{q_{\phi}(\mathbf{z}|\mathbf{x})} \ d\mathbf{z}$$
$$= \underbrace{\mathbb{E}_{\mathbf{z} \sim q_{\phi}(\mathbf{z}|\mathbf{x})}[\log p_{\theta}(\mathbf{x}|\mathbf{z})]}_{\mathsf{Reconstruction error}} - \underbrace{\mathrm{KL}[q_{\phi}(\mathbf{z}|\mathbf{x})||p_{\lambda}(\mathbf{z})]}_{\mathsf{Regularization}}$$

Let us assume the following distributions:

 $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \operatorname{diag}(\boldsymbol{\sigma}^2))$ encoder

 $p_{\theta}(\mathbf{x}|\mathbf{z}) = \operatorname{Bern}(\theta(\mathbf{z}))$ decoder

 $p_{\lambda}(\mathbf{z}) = \mathcal{N}(0, \mathbf{I})$ prior





Let us assume the following distributions:

 $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \operatorname{diag}(\boldsymbol{\sigma}^2))$ encoder



 $p_{\lambda}(\mathbf{z})$ $q_{\phi}(\mathbf{z}|\mathbf{x})$ $p_{\theta}(\mathbf{x}|\mathbf{z})$

Let us assume the following distributions:

 $q_{\phi}(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \operatorname{diag}(\boldsymbol{\sigma}^2))$ encoder

 $p_{\theta}(\mathbf{x}|\mathbf{z}) = \operatorname{Bern}(\theta(\mathbf{z}))$ decoder





$q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) \ p_{\lambda}(\mathbf{z})$



 $p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})$ $q_{\phi}(\mathbf{z}|\mathbf{x})$ **Fully-connected** ConvNets **PixelCNN**



 $\begin{array}{c} q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z}) \\ \\ \text{Normalizing flows} \\ \text{Volume-preserving flows} \end{array} \qquad \begin{array}{c} \text{Fully-connected} \\ \text{ConvNets} \\ \text{PixelCNN} \end{array}$







 $q_{\phi}(\mathbf{z}|\mathbf{x}) \propto p_{\theta}(\mathbf{x}|\mathbf{z}) p_{\lambda}(\mathbf{z})$

Autoregressive Prior Objective Prior Stick-Breaking Prior VampPrior



• Let's re-write the ELBO:

$$\mathbb{E}_{\mathbf{x}\sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \geq \mathbb{E}_{\mathbf{x}\sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \right] + \mathbb{E}_{\mathbf{x}\sim q(\mathbf{x})} \left[\mathbb{H}[q_{\phi}(\mathbf{z}|\mathbf{x})] \right] + \mathbb{E}_{\mathbf{x}\sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right]$$

• Let's re-write the ELBO:

$$\begin{split} \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \geq \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \right] + \\ + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{H}[q_{\phi}(\mathbf{z}|\mathbf{x})] \right] + \\ \\ \mathsf{Empirical distribution} - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right] \end{split}$$

• Let's re-write the ELBO:

 $\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \underbrace{\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \right]}_{+ \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{H} \left[q_{\phi}(\mathbf{z}|\mathbf{x}) \right] \right]_{+} \\ - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[- \ln p_{\lambda}(\mathbf{z}) \right]}$

• Let's re-write the ELBO:

$$\begin{split} \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \geq & \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})] \right] + \\ & \left(+ \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{H}[q_{\phi}(\mathbf{z}|\mathbf{x})] \right] + \right] \\ & - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [-\ln p_{\lambda}(\mathbf{z})] \end{split}$$
Encoder's entropy

• Let's re-write the ELBO:

$$\begin{split} \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \geq & \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \right] + \\ & + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{H}[q_{\phi}(\mathbf{z}|\mathbf{x})] \right] + \\ & \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right] \end{split}$$
Cross Entropy

• Let's re-write the ELBO:

$$\mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \geq \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \right] + \\ + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{H}[q_{\phi}(\mathbf{z}|\mathbf{x})] \right] + \\ - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right]$$

Aggregated posterior

$$q(\mathbf{z}) = \mathbb{E}_{q(\mathbf{x})}[q_{\phi}(\mathbf{z}|\mathbf{x})]$$
$$= \frac{1}{N} \sum_{n=1}^{N} q_{\phi}(\mathbf{z}|\mathbf{x}_n)$$

• Let's re-write the ELBO:

$$\begin{array}{l} \max. \ \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \geq & \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} [\ln p_{\theta}(\mathbf{x}|\mathbf{z})] \right] + \quad \text{Variance} \\ & + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{H}[q_{\phi}(\mathbf{z}|\mathbf{x})] \right] + \quad \text{Variance} \\ & - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [-\ln p_{\lambda}(\mathbf{z})] \end{array}$$

• Let's re-write the ELBO:

$$\begin{split} \max & \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\ln p(\mathbf{x}) \right] \geq \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{E}_{q_{\phi}(\mathbf{z}|\mathbf{x})} \left[\ln p_{\theta}(\mathbf{x}|\mathbf{z}) \right] \right] + \quad \text{Variance} \\ & + \mathbb{E}_{\mathbf{x} \sim q(\mathbf{x})} \left[\mathbb{H}[q_{\phi}(\mathbf{z}|\mathbf{x})] \right] + \quad \text{Variance} \\ & - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right] \\ \end{split} \\ \\ \min & \mathbb{H}[q(\mathbf{z})] + KL[q(\mathbf{z})||p_{\lambda}(\mathbf{z})] \end{split}$$

min. $\mathbb{H}[q(\mathbf{z})] + KL[q(\mathbf{z})||p_{\lambda}(\mathbf{z})]$



Prior

Aggregated posterior

 $\min\left(\mathbb{H}[q(\mathbf{z})] + KL[q(\mathbf{z})||p_{\lambda}(\mathbf{z})]\right)$





min. $\mathbb{H}[q(\mathbf{z})] + KL[q(\mathbf{z})||p_{\lambda}(\mathbf{z})]$



min. $\mathbb{H}[q(\mathbf{z})] + KL[q(\mathbf{z})||p_{\lambda}(\mathbf{z})]$



Standard prior is too strong and overregularizes the encoder.

min. $\mathbb{H}[q(\mathbf{z})] + KL[q(\mathbf{z})||p_{\lambda}(\mathbf{z})]$



Standard prior is too strong and overregularizes the encoder.

What is the "optimal" prior?

• We look for **the optimal prior** using the Lagrange function:

$$\max_{p_{\lambda}(\mathbf{z})} - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [-\ln p_{\lambda}(\mathbf{z})] + \beta \Big(\int p_{\lambda}(\mathbf{z}) d\mathbf{z} - 1\Big)$$

- The solution is simply **the aggregated posterior**.
- We approximate it using K pseudo-inputs instead of N observations:

$$p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z} | \mathbf{u}_k)$$

• We look for the optimal prior using the Lagrange function:

$$\max_{p_{\lambda}(\mathbf{z})} - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right] + \beta \left(\int p_{\lambda}(\mathbf{z}) d\mathbf{z} - 1 \right)$$

• The solution is simply the aggregated posterior.

$$p_{\lambda}^{*}(\mathbf{z}) = \frac{1}{N} \sum_{n=1}^{N} q_{\phi}(\mathbf{z}|\mathbf{x}_{n})$$

We approximate it using K pseudo-inputs instead of N observations:

$$p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z} | \mathbf{u}_k)$$

• We look for the optimal prior using the Lagrange function:

$$\max_{p_{\lambda}(\mathbf{z})} - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right] + \beta \left(\int p_{\lambda}(\mathbf{z}) d\mathbf{z} - 1 \right)$$

• The solution is simply the aggregated posterior.

$$p_{\lambda}^{*}(\mathbf{z}) = \frac{1}{N} \sum_{n=1}^{N} q_{\phi}(\mathbf{z} | \mathbf{x}_{n})$$

We approximate it using K pseudo-inputs instead of N observations:

infeasible

$$p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z} | \mathbf{u}_k)$$

• We look for the optimal prior using the Lagrange function:

$$\max_{p_{\lambda}(\mathbf{z})} - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} [-\ln p_{\lambda}(\mathbf{z})] + \beta \Big(\int p_{\lambda}(\mathbf{z}) d\mathbf{z} - 1\Big)$$

- The solution is simply the aggregated posterior.
- We approximate it using *K* pseudo-inputs instead of *N* observations:

$$p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z} | \mathbf{u}_k)$$

• We look for the optimal prior using the Lagrange function:

$$\max_{p_{\lambda}(\mathbf{z})} - \mathbb{E}_{\mathbf{z} \sim q(\mathbf{z})} \left[-\ln p_{\lambda}(\mathbf{z}) \right] + \beta \left(\int p_{\lambda}(\mathbf{z}) d\mathbf{z} - 1 \right)$$

- The solution is simply the aggregated posterior.
- We approximate it using *K* pseudo-inputs instead of *N* observations:

$$p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z} \mathbf{u}_{k})$$

they are trained from scratch

- Is the VampPrior different than the Mixture of Gaussians? $p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} \mathcal{N}(\mu_k, \operatorname{diag}(\sigma_k^2))$
- VampPrior: the prior and the posterior must "cooperate" during training.

VampPrior

$$\begin{split} &\frac{1}{K}\sum_{k=1}^{K} \Big\{ \left(\frac{q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \ \frac{\partial}{\partial \phi_{i}}q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) - q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) \ \frac{\partial}{\partial \phi_{i}}q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})}{\frac{1}{K}\sum_{k=1}^{K} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) \ q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})} \right) + \\ &+ \Big(\frac{\left(q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \ \frac{\partial}{\partial \mathbf{z}_{\phi}}q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) - q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) \ \frac{\partial}{\partial \mathbf{z}_{\phi}}q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \right) \ \frac{\partial}{\partial \phi_{i}}\mathbf{z}_{\phi}^{(l)}}{\frac{1}{K}\sum_{k=1}^{K} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) \ q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})} \end{split} \Big) \Big\} \end{split}$$

standard/ MoG

$$\frac{1}{p_{\lambda}(\mathbf{z}_{\phi}^{(l)}) \ q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})} \Big(q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \frac{\partial}{\partial \mathbf{z}_{\phi}} p_{\lambda}(\mathbf{z}_{\phi}^{(l)}) - p_{\lambda}(\mathbf{z}_{\phi}^{(l)}) \frac{\partial}{\partial \mathbf{z}_{\phi}} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \Big) \frac{\partial}{\partial \phi_{i}} \mathbf{z}_{\phi}^{(l)}$$

- Is the VampPrior different than the Mixture of Gaussians? $p_{\lambda}(\mathbf{z}) = \frac{1}{K} \sum_{k=1}^{K} \mathcal{N}(\mu_k, \operatorname{diag}(\sigma_k^2))$
- VampPrior: the prior and the posterior must "cooperate" during training.

VampPrior

$$\frac{1}{K} \sum_{k=1}^{K} \left\{ \left(\frac{q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \frac{\partial}{\partial \phi_{i}} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) - q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) \frac{\partial}{\partial \phi_{i}} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})}{\frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})} \right) + \left(\frac{\left(q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \frac{\partial}{\partial \mathbf{z}_{\phi}} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) - q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) \frac{\partial}{\partial \mathbf{z}_{\phi}} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \right) \frac{\partial}{\partial \phi_{i}} \mathbf{z}_{\phi}^{(l)}}{\frac{1}{K} \sum_{k=1}^{K} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{u}_{k}) q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})} \right) \right\}$$

standard/ MoG

$$\frac{1}{p_{\lambda}(\mathbf{z}_{\phi}^{(l)}) \ q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x})} \Big(q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \frac{\partial}{\partial \mathbf{z}_{\phi}} p_{\lambda}(\mathbf{z}_{\phi}^{(l)}) - p_{\lambda}(\mathbf{z}_{\phi}^{(l)}) \frac{\partial}{\partial \mathbf{z}_{\phi}} q_{\phi}(\mathbf{z}_{\phi}^{(l)}|\mathbf{x}) \Big) \frac{\partial}{\partial \phi_{i}} \mathbf{z}_{\phi}^{(l)}$$

- VampPrior is closely related to the **Empirical Bayes**.
 - We propose a new approach that learns parameters of the prior and combines the variational

inference with the EB approach.

- VampPrior is closely related to the **Information Bottleneck**.
 - The aggregated posterior naturally plays the role of the prior.
 - The VampPrior brings the VAE and the IB formulations together.

Hierarchical VampPrior VAE

Typical issue in hierarchical VAE: inactive stochastic units

 $p(\mathbf{z}_2) = \frac{1}{K} \sum_{k=1}^{K} q_{\psi}(\mathbf{z}_2 | \mathbf{u}_k),$ $p_{\lambda}(\mathbf{z}_1 | \mathbf{z}_2) = \mathcal{N}(\mathbf{z}_1 | \mu_{\lambda}(\mathbf{z}_2), \operatorname{diag}(\sigma_{\lambda}^2(\mathbf{z}_2))),$ $q_{\phi}(\mathbf{z}_1 | \mathbf{x}, \mathbf{z}_2) = \mathcal{N}(\mathbf{z}_1 | \mu_{\phi}(\mathbf{x}, \mathbf{z}_2), \operatorname{diag}(\sigma_{\phi}^2(\mathbf{x}, \mathbf{z}_2))))$ $q_{\psi}(\mathbf{z}_2 | \mathbf{x}) = \mathcal{N}(\mathbf{z}_2 | \mu_{\psi}(\mathbf{x}), \operatorname{diag}(\sigma_{\psi}^2(\mathbf{x})))$



Hierarchical VampPrior VAE

Typical issue in hierarchical VAE: inactive stochastic units

$$p(\mathbf{z}_{2}) = \frac{1}{K} \sum_{k=1}^{K} q_{\psi}(\mathbf{z}_{2} | \mathbf{u}_{k}),$$

$$p_{\lambda}(\mathbf{z}_{1} | \mathbf{z}_{2}) = \mathcal{N}(\mathbf{z}_{1} | \mu_{\lambda}(\mathbf{z}_{2}), \operatorname{diag}(\sigma_{\lambda}^{2}(\mathbf{z}_{2})))),$$

$$q_{\phi}(\mathbf{z}_{1} | \mathbf{x}, \mathbf{z}_{2}) = \mathcal{N}(\mathbf{z}_{1} | \mu_{\phi}(\mathbf{x}, \mathbf{z}_{2}), \operatorname{diag}(\sigma_{\phi}^{2}(\mathbf{x}, \mathbf{z}_{2})))),$$

$$q_{\psi}(\mathbf{z}_{2} | \mathbf{x}) = \mathcal{N}(\mathbf{z}_{2} | \mu_{\psi}(\mathbf{x}), \operatorname{diag}(\sigma_{\psi}^{2}(\mathbf{x})))$$

1





| | VAE | (L = 1) | HVAE $(L=2)$ | | CONVHVAE $(L=2)$ | | PIXELHVAE $(L = 2)$ | |
|----------------|---------------------------|-----------|---------------------------|-----------|---------------------------|-----------|---------------------------|-----------|
| DATASET | $\operatorname{standard}$ | VampPrior | $\operatorname{standard}$ | VampPrior | $\operatorname{standard}$ | VampPrior | $\operatorname{standard}$ | VampPrior |
| staticMNIST | -88.56 | -85.57 | -86.05 | -83.19 | -82.41 | -81.09 | -80.58 | -79.78 |
| dynamicMNIST | -84.50 | -82.38 | -82.42 | -81.24 | -80.40 | -79.75 | -79.70 | -78.45 |
| Omniglot | -108.50 | -104.75 | -103.52 | -101.18 | -97.65 | -97.56 | -90.11 | -89.76 |
| Caltech 101 | -123.43 | -114.55 | -112.08 | -108.28 | -106.35 | -104.22 | -85.51 | -86.22 |
| Frey Faces | 4.63 | 4.57 | 4.61 | 4.51 | 4.49 | 4.45 | 4.43 | 4.38 |
| Histopathology | 6.07 | 6.04 | 5.82 | 5.75 | 5.59 | 5.58 | 4.84 | 4.82 |

Table 2: Test LL for static MNIST.

| Model | LL |
|---|---------------------|
| VAE $(L = 1) + NF$ 32 | -85.10 |
| VAE $(L = 2)$ 6 | -87.86 |
| IWAE $(L=2)$ 6 | -85.32 |
| $\mathrm{HVAE}~(L=2)~+~\mathrm{SG}$ | -85.89 |
| $\mathrm{HVAE}\ (L=2)\ +\ \mathrm{MoG}$ | -85.07 |
| HVAE $(L = 2)$ + VAMPPRIOR data | -85.71 |
| HVAE $(L = 2)$ + VampPrior | -83.19 |
| $AVB + AC \ (L = 1)$ 28 | -80.20 |
| VLAE 7 | -79.03 |
| VAE + IAF 18 | -79.88 |
| CONVHVAE (L = 2) + VAMPPRIOR | -81.09 |
| PixelHVAE $(L = 2)$ + VampPrior | -79.78 |



Figure 3: A comparison between two-level VAE and IWAE with the standard normal prior and theirs Vamp-Prior counterpart in terms of number of active units for varying number of pseudo-inputs on static MNIST.

| Table 3: | Test | LL | for | dynamic | MNIST. |
|----------|------|----|-----|---------|--------|
| | | | | •/ | |

| Model | LL |
|---------------------------------|--------|
| VAE $(L=2)$ + VGP 40 | -81.32 |
| CAGEM-0 $(L = 2)$ 25 | -81.60 |
| LVAE $(L = 5)$ 36 | -81.74 |
| HVAE $(L = 2)$ + VAMPPRIOR data | -81.71 |
| HVAE $(L = 2)$ + VampPrior | -81.24 |
| VLAE 7 | -78.53 |
| VAE + IAF 18 | -79.10 |
| PixelVAE 15 | -78.96 |
| CONVHVAE $(L = 2)$ + VAMPPRIOR | -79.78 |
| PixelHVAE $(L = 2)$ + VampPrior | -78.45 |

Table 4: Test LL for OMNIGLOT.

| Model | $\mathbf{L}\mathbf{L}$ |
|---------------------------------|------------------------|
| VR-max $(L=2)$ 24 | -103.72 |
| IWAE $(L=2)$ 6 | -103.38 |
| LVAE $(L = 5)$ 36 | -102.11 |
| HVAE $(L=2)$ + VampPrior | -101.18 |
| VLAE 7 | -89.83 |
| CONVHVAE $(L = 2)$ + VAMPPRIOR | -97.56 |
| PixelHVAE $(L = 2)$ + VampPrior | -89.76 |

Table 5: Test LL for Caltech 101 Silhouettes.

| Model | LL |
|---|---------------------|
| IWAE $(L = 1)$ 24 | -117.21 |
| VR-max $(L=1)$ 24 | -117.10 |
| HVAE $(L = 2)$ + VAMPPRIOR | -108.28 |
| VLAE 7 | -78.53 |
| $\operatorname{convHVAE}(L=2) + \operatorname{VampPrior}$ | -104.22 |
| PIXELHVAE $(L = 2) + VAMPPRIOR$ | -86.22 |



Figure 4: $(top \ row)$ Images generated by PIXELHVAE + VAMPPRIOR for chosen pseudo-input in the left top corner. $(bottom \ row)$ Images represent a subset of trained pseudo-inputs for different datasets.



Figure 5: (a) Real images from test sets and images generated by (b) the vanilla VAE, (c) the HVAE (L = 2) + VampPrior, (d) the convHVAE (L = 2) + VampPrior and (e) the PixelHVAE (L = 2) + VampPrior.



Figure 6: (a) Real images from test sets, (b) reconstructions given by the vanilla VAE, (c) the HVAE (L = 2) + VampPrior, (d) the convHVAE (L = 2) + VampPrior and (e) the PixelHVAE (L = 2) + VampPrior.



Figure 6: (a) Real images from test sets, (b) reconstructions given by the vanilla VAE, (c) the HVAE (L = 2) + VampPrior, (d) the convHVAE (L = 2) + VampPrior and (e) the PixelHVAE (L = 2) + VampPrior.

| The prior in VAE is extremely important. | VampPrior = approximated aggregated posterior as the optimal prior |
|---|--|
| | Multimodal prior \rightarrow better generative process |

The **prior** in VAE is extremely important.

VampPrior = approximated aggregated posterior as the optimal prior

Hierarchical VampPrior VAE → **less** inactive stochastic units.

Multimodal prior \rightarrow **better** generative process

The **prior** in VAE is extremely important.

VampPrior = approximated aggregated posterior as the optimal prior

Hierarchical VampPrior VAE \rightarrow less inactive stochastic units.

Multimodal prior \rightarrow **better** generative process

The **prior** in VAE is extremely important.

VampPrior = approximated aggregated posterior as the optimal prior

Hierarchical VampPrior VAE \rightarrow less inactive stochastic units.

Multimodal prior \rightarrow **better** generative process

Future directions

VampPrior + Normalizing flows



Future directions

VampPrior for other data (sequential, sound, text, genomics, etc.)

 \rightarrow RNN posteriors





he tilving wint with Thighter left and Galad wiffed At foolohold In Septen Thysing realifier and generalist flours piflopen, hop the Heiperfeis Rollable vier me Sugar general wirt, in bill lett he Highage -Anlifistying Theffe im Rellacile in going al low men , fis de Rolonifles aif der Server foll ein Pfre Jamia Uni fillemente i telester begefte mertra, Saif Mr Bedjett Halefilytagen in elle helfe werheile in the yo seem geminister times for daystake weather in this fin son motores have the interes. for foldles hanged by fine here give there iting on 25.5. levil if de Alerreigtang be vorgansfaranka Weaking his pregantagen gi sterrespices to takale fingled the fire portionatorile and hill wit, hell at sigh mighting new hand were Rollabler mile 35 portinha.

Future directions

How to (better) learn pseudoinputs?

 \rightarrow MCMC?

→Wake-Sleep?



Webpage: https://jmtomczak.github.io/

Code on github: https://github.com/jmtomczak/

Contact: jakubmkt@gmail.com



Marie Skłodowska-Curie Actions

The research conducted by Jakub M. Tomczak was funded by the European Commission within the Marie Skłodowska-Curie Individual Fellowship (Grant No. 702666, "Deep learning and Bayesian inference for medical imaging").