

Machine Learning 2020

Volker Roth

Department of Mathematics & Computer Science
University of Basel

25th May 2020

Section 11

Non-linear latent variable models

Non-linear latent variable models

Latent variable $\mathbf{z} \rightsquigarrow$ Gaussian likelihood with

nonlinearly transformed mean

$$\mu = f(\mathbf{z}, \phi).$$

Prior and likelihood:

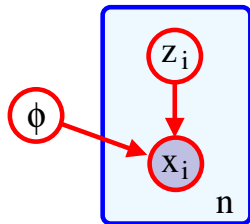
$$p(\mathbf{z}) = N(\mathbf{0}, I)$$

$$p(\mathbf{x}|\mathbf{z}, \phi) = N(f(\mathbf{z}, \phi), \sigma^2 I).$$

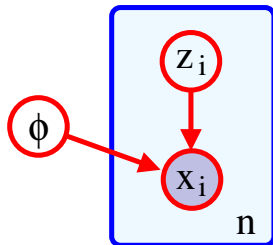
- Given observed \mathbf{x} , we want to understand what possible values of the hidden variable \mathbf{z} were responsible for it:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}.$$

- No closed form expression available. Cannot evaluate denominator $p(\mathbf{x})$ and so we can't even compute the numerical value of the posterior for a given pair \mathbf{z} and \mathbf{x} .



Sampling



- ...but it is easy to generate a new sample \mathbf{x}^* using sampling:
 - ▶ Draw \mathbf{z}^* from the prior $p(\mathbf{z})$, pass this through $f(\mathbf{z}^*, \phi)$
 \rightsquigarrow mean of likelihood $p(\mathbf{x}^* | \mathbf{z}^*)$,
 - ▶ then draw \mathbf{x}^* from this distribution.
- Prior and likelihood are normal distributions \rightsquigarrow sampling is easy.

Evaluating marginal likelihood (evidence)

$$\begin{aligned} p(\mathbf{x}|\phi) &= \int p(\mathbf{x}, \mathbf{z}|\phi) d\mathbf{z} \\ &= \int p(\mathbf{x}|\mathbf{z}, \phi) p(\mathbf{z}) d\mathbf{z} \\ &= \int N(\mathbf{f}[\mathbf{z}, \phi], \sigma^2 I) \cdot N(\mathbf{0}, I) d\mathbf{z}. \end{aligned}$$

No closed form for this integral \rightsquigarrow lower bound (Jensen's inequality):

$$\begin{aligned} \log[p(\mathbf{x}|\phi)] &= \log \left[\int p(\mathbf{x}, \mathbf{z}|\phi) d\mathbf{z} \right] \\ &= \log \left[\int q(\mathbf{z}) \frac{p(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} d\mathbf{z} \right] \\ &\geq \int q(\mathbf{z}) \log \left[\frac{p(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z})} \right] d\mathbf{z}, \end{aligned}$$

Known as the evidence lower bound ELBO, because $p(\mathbf{x}|\phi)$ is the evidence (= marginal likelihood) in the context of Bayes' rule.

- In practice, the distribution $q(\mathbf{z})$ will have some parameters θ :

$$\text{ELBO}[\theta, \phi] = \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z}.$$

- To learn the non-linear latent variable model, we'll maximize this quantity as a function of both ϕ and θ .
- We will see: the maximum is obtained (theoretically) if the variational distribution is the true posterior, $q(\mathbf{z}|\theta) = p(\mathbf{z}|\mathbf{x}, \phi)$.
- In practice, we maximize ELBO over some tractable family of distributions $q(\mathbf{z}|\mathbf{x}, \theta)$ to obtain an approximation of the intractable posterior.
- The neural architecture that computes this is the **variational autoencoder**.

Tightness of ELBO

$$\begin{aligned}\text{ELBO}[\theta, \phi] &= \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\&= \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{z}|\mathbf{x}, \phi)p(\mathbf{x}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\&= \int q(\mathbf{z}|\theta) \log [p(\mathbf{x}|\phi)] d\mathbf{z} + \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{z}|\mathbf{x}, \phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\&= \log[p(\mathbf{x}|\phi)] + \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{z}|\mathbf{x}, \phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\&= \log[p(\mathbf{x}|\phi)] - D_{KL} [q(\mathbf{z}|\theta) \| p(\mathbf{z}|\mathbf{x}, \phi)].\end{aligned}$$

ELBO is the log marginal likelihood minus $D_{KL} [q(\mathbf{z}|\theta) \| p(\mathbf{z}|\mathbf{x}, \phi)]$.

D_{KL} zero when $q(\mathbf{z}|\theta) = p(\mathbf{z}|\mathbf{x}, \phi) \rightsquigarrow \text{ELBO} = \log[p(\mathbf{x}|\phi)]$.

ELBO as reconstruction loss minus KL to prior

$$\begin{aligned}\text{ELBO}[\theta, \phi] &= \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{x}, \mathbf{z}|\phi)}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\ &= \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{x}|\mathbf{z}, \phi)p(\mathbf{z})}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\ &= \int q(\mathbf{z}|\theta) \log [p(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z} + \int q(\mathbf{z}|\theta) \log \left[\frac{p(\mathbf{z})}{q(\mathbf{z}|\theta)} \right] d\mathbf{z} \\ &= \int q(\mathbf{z}|\theta) \log [p(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z} - D_{\text{KL}}[q(\mathbf{z}|\theta), p(\mathbf{z})]\end{aligned}$$

- First term measures the average agreement $p(\mathbf{x}|\mathbf{z}, \phi)$ of the hidden variable and the data (reconstruction loss)
- Second one measures the degree to which the auxiliary distribution $q(\mathbf{z}, \theta)$ matches the prior.

The variational approximation

- ELBO is tight when we choose $q(\mathbf{z}|\boldsymbol{\theta}) = p(\mathbf{z}|\mathbf{x}, \phi)$.
- Intractable \rightsquigarrow variational approximation: choose simple parametric form for $q(\mathbf{z}|\boldsymbol{\theta})$, use it as an approximation to the true posterior.
- Choose a normal distribution with parameters $\boldsymbol{\mu}$ and $\Sigma = \sigma^2 I$.
- Optimization \rightsquigarrow find normal distribution closest to true posterior $p(\mathbf{z}|\mathbf{x})$. Corresponds to minimizing the KL divergence.
- True posterior $p(\mathbf{z}|\mathbf{x})$ depends on \mathbf{x}
 \rightsquigarrow variational approximation should also depend on \mathbf{x} :

$$q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{x}) = N(g_{\boldsymbol{\mu}}[\mathbf{x}|\boldsymbol{\theta}], g_{\sigma^2}[\mathbf{x}|\boldsymbol{\theta}]),$$

where $g[\mathbf{x}, \boldsymbol{\theta}]$ is a neural network with parameters $\boldsymbol{\theta}$.

The variational autoencoder

Recall

$$\text{ELBO}[\theta, \phi] = \int q(\mathbf{z}|\mathbf{x}, \theta) \log [p(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z} - D_{KL}[q(\mathbf{z}|\mathbf{x}, \theta), p(\mathbf{z})]$$

Involves an intractable integral, but it is an expectation

↪ approximate with samples:

$$E_{q(\mathbf{z}|\mathbf{x}, \theta)}[\log [p(\mathbf{x}|\mathbf{z}, \phi)]] \approx \frac{1}{N} \sum_{n=1}^N \log [p(\mathbf{x}|\mathbf{z}_n^*, \phi)]$$

where \mathbf{z}_n^* is the n -th sample from $q(\mathbf{z}|\mathbf{x}, \theta)$. Limit: use a single sample:

$$\text{ELBO}[\theta, \phi] \approx \log [p(\mathbf{x}|\mathbf{z}^*, \phi)] - D_{KL}[q(\mathbf{z}|\mathbf{x}, \theta), p(\mathbf{z})]$$

The second term is just the KL divergence between two Gaussians and is available in closed form.

The reparameterization trick

Recall: Want to sample from

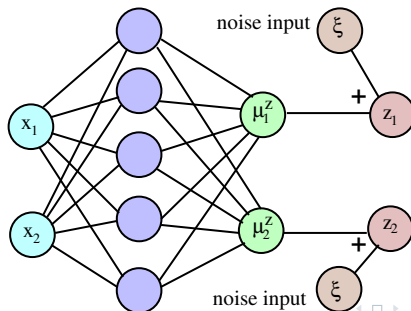
$$q(\mathbf{z}|\boldsymbol{\theta}, \mathbf{x}) = N(g_{\boldsymbol{\mu}}[\mathbf{x}|\boldsymbol{\theta}], g_{\sigma^2}[\mathbf{x}|\boldsymbol{\theta}]),$$

To let PyTorch / Tensorflow perform automatic differentiation via backpropagation, we must avoid the sampling step.

Simple solution: draw a sample $\boldsymbol{\xi} \sim N(0, I)$ and use

$$\mathbf{z}^* = g_{\boldsymbol{\mu}} + \sigma^{1/2} \boldsymbol{\xi}.$$

Now “the gradient can flow through the network”. **Encoder network:**



- Finally, minimize negative expectation of ELBO over $p(\mathbf{x})$:

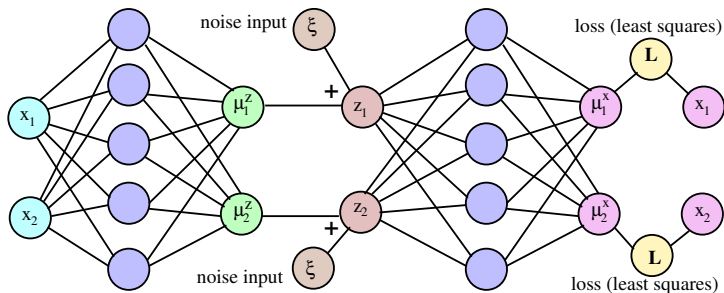
$$\min_{\phi, \theta} -E_{p(\mathbf{x})} E_{q(\mathbf{z}|\mathbf{x}, \theta)} [\log [p(\mathbf{x}|\mathbf{z}, \phi)]] + E_{p(\mathbf{x})} D_{KL}[q(\mathbf{z}|\mathbf{x}, \theta), p(\mathbf{z})]$$

- The first term is approximated as

$$E_{p(\mathbf{x})} E_{q(\mathbf{z}|\mathbf{x}, \theta)} [\log [p(\mathbf{x}|\mathbf{z}, \phi)]] \approx \frac{1}{n} \sum_{i=1}^n \log [p(\mathbf{x}_i|\mathbf{z}_i^*, \phi)].$$

We assume $p(\mathbf{x}_i|\mathbf{z}_i^*, \phi) = \mathcal{N}(f_\phi(\mathbf{z}_i^*), \sigma^2)$,

where f is implemented via a neural net: \rightsquigarrow **Decoder network**

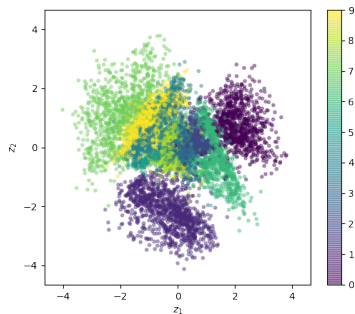
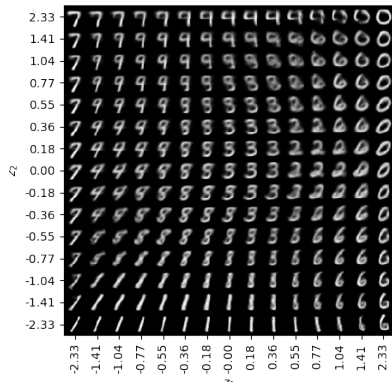


Further Variations

- For maximizing ELBO, we jointly optimize over the parameters of encoder and decoder network.
- When adjusting the decoder, we also change the “true” posterior that we are going to approximate!
- So approximation quality should not be our only goal... need “tuning knob” for steering the model into a desired direction.
- Solution: introduce parameter $\beta > 0$ that controls the relative importance of the two loss terms:

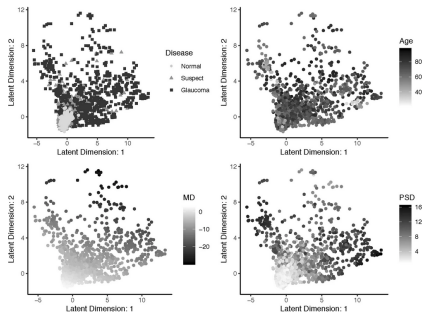
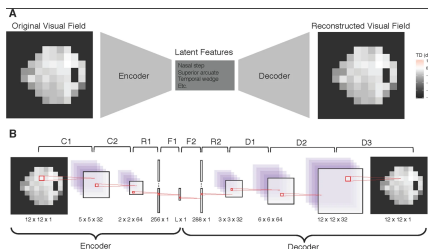
$$\min_{\theta, \phi} \frac{1}{n} \sum_{i=1}^n D_{KL}[q(\mathbf{z}|\mathbf{x}_i, \theta), p(\mathbf{z})] - \beta \frac{1}{n} \sum_{i=1}^n \log [p(\mathbf{x}_i|\mathbf{z}_i^*, \phi)]$$

Applications: MNIST example



Taken from Louis Tiago: A Tutorial on Variational Autoencoders with a Concise Keras Implementation

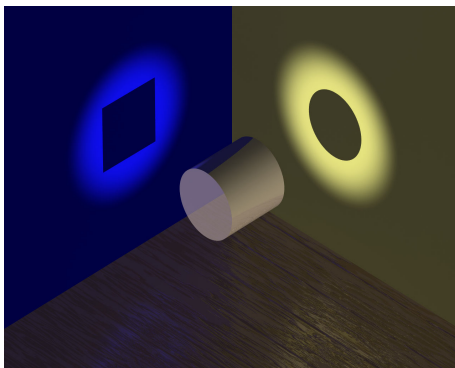
Applications: Medical example



Berchuck, S.I., Mukherjee, S. & Medeiros, F.A. Estimating Rates of Progression and Predicting Future Visual Fields in Glaucoma Using a Deep Variational Autoencoder. Sci Rep 9, 18113 (2019). <https://doi.org/10.1038/s41598-019-54653-6>

Multiple Views: Deep Information Bottleneck

- Consider paired samples from different views.
- What is the dependency structure between the views ?
- Nonlinear model: dependency detected by **deep IB**.



Two-view version: The deep information bottleneck

- So far we argued that since the true posterior $p(\mathbf{z}|\mathbf{x})$ depends on \mathbf{x} , the variational approximation should also depend on \mathbf{x} .
- Restricted setting: explain posterior **only by external variable $\tilde{\mathbf{x}}$** :

$$q = q(\mathbf{z}|\boldsymbol{\theta}, \tilde{\mathbf{x}}).$$

$$\begin{aligned}\text{ELBO}[\boldsymbol{\theta}, \phi] &= \int q(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta}) \log [p(\mathbf{x}|\mathbf{z}, \phi)] d\mathbf{z} - D_{KL}[q(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta}), p(\mathbf{z})] \\ &= E_{q(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta})} \log [p(\mathbf{x}|\mathbf{z}, \phi)] - D_{KL}[q(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta}), p(\mathbf{z})]\end{aligned}$$

- Connection to IB:

- ▶ Assume (or define) $q(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta}) := p(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta})$
- ▶ Take expectation w.r.t. joint data distribution $p(\tilde{\mathbf{x}}, \mathbf{x})$:

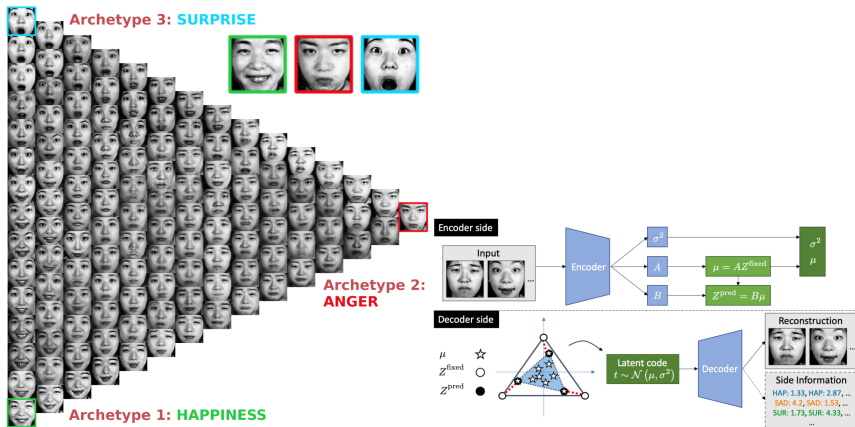
$$E_{p(\tilde{\mathbf{x}}, \mathbf{x})} E_{p(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta})} \log [p(\mathbf{x}|\mathbf{z}, \phi)] - E_{p(\tilde{\mathbf{x}})} D_{KL}[p(\mathbf{z}|\tilde{\mathbf{x}}, \boldsymbol{\theta}), p(\mathbf{z})]$$

- ▶ First term $\leq \mathcal{I}_{\boldsymbol{\theta}, \phi}(\mathbf{z}; \mathbf{x}) + \text{const.}$ Second term $= \mathcal{I}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}; \mathbf{z}),$

- This defines the deep information bottleneck (with weight β)

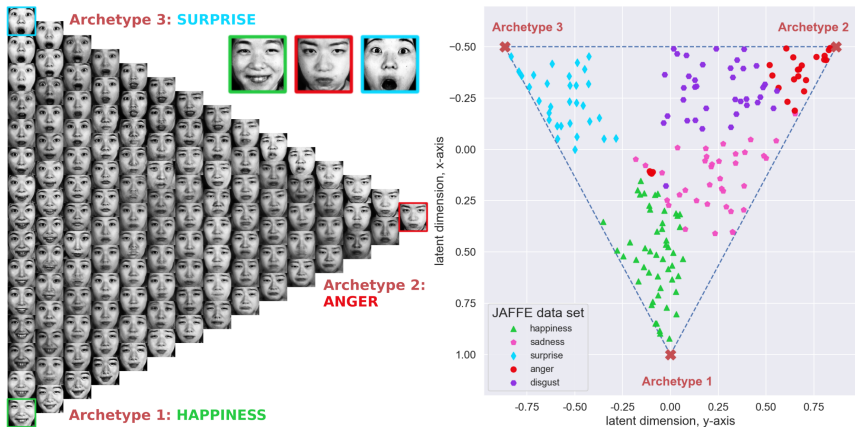
$$\min_{\phi, \boldsymbol{\theta}} \mathcal{I}_{\boldsymbol{\theta}}(\tilde{\mathbf{x}}; \mathbf{z}) - \beta \mathcal{I}_{\boldsymbol{\theta}, \phi}^{\text{low}}(\mathbf{z}; \mathbf{x}), \quad \text{where } \mathcal{I}^{\text{low}} \text{ is a lower bound of } \mathcal{I}.$$

Applications: Face images



Keller et al. 2020: Learning Extremal Representations with Deep Archetypal Analysis

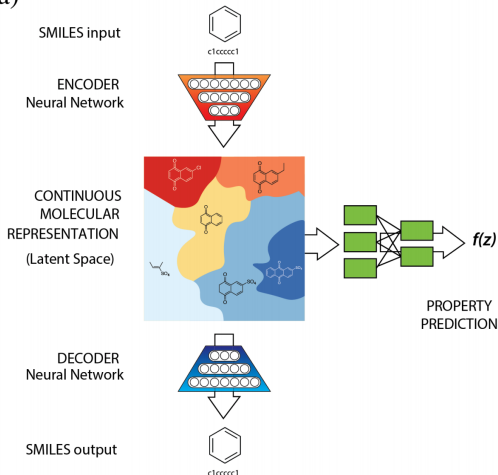
Applications: Face images



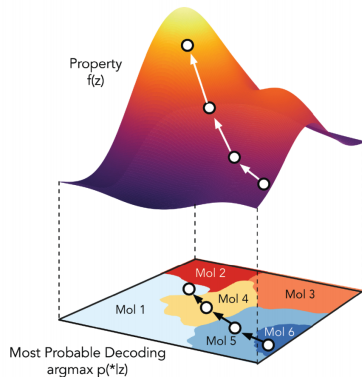
Keller et al. 2020: Learning Extremal Representations with Deep Archetypal Analysis

Applications: Deep Chemical Variational Autoencoders

(a)



(b)



(Gomez-Bombarelli et al., ACS Cent Sci, 2018)