

Variational Auto-Encoder

Yingzi Bu



Department of Pharmaceutical Sciences
University of Michigan

Aug 21, 2023

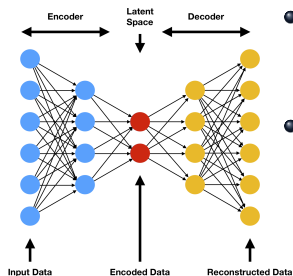
References

- 1 Kingma, Diederik P., and Max Welling. "Auto-encoding variational bayes." arXiv preprint arXiv:1312.6114 (2013). Cited by 28916
- 2 Kingma, Durk P., et al. "Semi-supervised learning with deep generative models." Advances in neural information processing systems 27 (2014). Cited by 3102
- 3 Gomez-Bombarelli, Rafael, et al. "Automatic chemical design using a data-driven continuous representation of molecules." ACS central science 4.2 (2018): 268-276. Cited by 2569
- 4 Zhavoronkov, Alex, et al. "Deep learning enables rapid identification of potent DDR1 kinase inhibitors." Nature biotechnology 37.9 (2019): 1038-1040. Cited by 769

Probability is all you need

$p(\text{The Sun will rise tomorrow})$	≈ 1
$p(\text{cure cancer} \mid \text{compounds})$	$= ???$
$p(\text{compounds} \mid \text{cure cancer})$	$= ???$
$p(\text{approximation, interesting} \mid \text{Biologist, Chemist})$	≈ 1
$p(\text{Talks in math = boring} \mid \text{Biologist, Chemist})$	> 0.5
$p(\text{Math proof} \mid \text{you want to know ML})$	$= 1$

Autoencoders



- Two sets:
 - the space of decoded messages $\mathcal{X} \subseteq \mathbb{R}^n$;
 - the space of encoded messages $\mathcal{Z} \subseteq \mathbb{R}^m$
- Two parametrized families of functions:
 - the encoder family: $E_\phi : \mathcal{X} \rightarrow \mathcal{Z}$, parametrized by ϕ
 - the decoder family: $D_\theta : \mathcal{Z} \rightarrow \mathcal{X}$, parametrized by θ

for $x \in \mathcal{X}$, latent variable $z = E_\phi(x) \in \mathcal{Z}$

for $z \in \mathcal{Z}$, decoded message $x' = D_\theta(z) \in \mathcal{X}$

Training Autoencoders

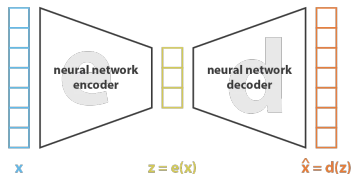
A reference probability distribution μ_{ref} over \mathcal{X} ,
a function $d: \mathcal{X} \times \mathcal{X} \rightarrow [0, \infty)$ measures how much x' differs from x

$$L(\theta, \phi) := \mathbb{E}_{x \sim \mu_{\text{ref}}} [d(x, D_{\theta}(E_{\phi}(x)))]$$

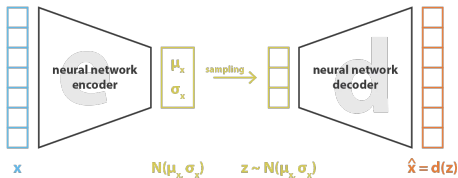
The least-squares optimization:

$$\min_{\theta, \phi} L(\theta, \phi), \text{ where } L(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N \|x_i - D_{\theta}(E_{\phi}(x_i))\|_2^2$$

Variational Autoencoder (VAE)



$$\text{loss} = \|x - \hat{x}\|^2 = \|x - d(z)\|^2 = \|x - d(e(x))\|^2$$



$$\text{loss} = \|x - \hat{x}\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)] = \|x - d(z)\|^2 + \text{KL}[N(\mu_x, \sigma_x), N(0, I)]$$

Evidence lower bound (ELBO)

p^* true distribution of x , p_θ estimated distribution of x

Objective: $\max \mathbb{E}_{x \sim p^*(\cdot)} [\log p_\theta(x)]$

Alternative objective: $\max_{\theta, \phi} ELBO = \max_{\theta, \phi} \mathbb{E}_{z \sim q_\phi(\cdot|x)} \log \frac{p_\theta(x, z)}{q_\phi(z|x)}$

$$\begin{aligned} \log p_\theta(x) &= \mathbb{E}_{x \sim p^*(\cdot)} [\log p_\theta(x)] \\ &= KL(q_\phi(z|x) || p_\theta(z|x)) + \mathbb{E}_{q_\phi(z|x)} [-\log q_\phi(z|x) + \log p_\theta(x, z)] \\ &\geq ELBO = -KL(q_\phi(z|x) || p_\theta(z)) + \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x|z)] \end{aligned}$$

$$\tilde{\mathcal{L}}_{VAE}(x; \theta, \phi) = -KL(q_\phi(z|x) || p_\theta(z)) + \frac{1}{L} \sum_{l=1}^L \log p_\theta(x|z^{(l)})$$

Loss function

$$\begin{aligned}\min \text{Loss}(x; \theta, \phi) &= \min -\tilde{\mathcal{L}}_{\text{VAE}}(x; \theta, \phi) \\&= \min KL(q_{\phi}(z|x) || p_{\theta}(z)) - \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(x|z^{(l)}) \\&= \min \left[KL(\mathcal{N}(\mu, \Sigma) || \mathcal{N}(0, I)) - \frac{1}{L} \sum_{l=1}^L \log p_{\theta}(x|z^{(l)}) \right] \\&= \min \left[\left(-\frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2 - \mu_j^2 - \sigma_j^2) \right) - \sum_{l=1}^L \log p_{\theta}(x|z^{(l)}) \right] \\&= \min(\text{KLD} + \text{BCE} / \text{MSE})\end{aligned}$$

```
def loss_function(recon_x, x, mu, logvar):  
    BCE = F.binary_cross_entropy(recon_x.cuda(),  
                                  x.view(-1, 784).cuda(),  
                                  reduction='sum')  
    KLD = -0.5 * torch.sum(1 + logvar - mu.pow(2) - logvar.exp())  
    return BCE + KLD
```


AE vs VAE generation, 100 epochs on MNIST data set

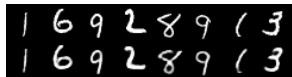


Figure 3: AE $x' \sim p_{\theta}(*|x)$

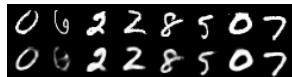


Figure 5: VAE $x' \sim p_{\theta}(*|x)$

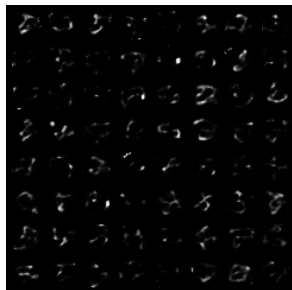


Figure 4: AE $x' \sim p_{\theta}(*|z)$,
 $z \sim \mathcal{N}(0, I)$



Figure 6: VAE $x' \sim p_{\theta}(*|z)$,
 $z \sim \mathcal{N}(0, I)$

ML question: AE latent space does not follow normal distribution anyway.

What if I want compounds that have certain properties?

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

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

$$\mathbf{x} \sim p_{\theta}(\cdot|\mathbf{z}, y)$$

Question for audience: What will the architecture look like?

Math proof is left for your practice.

Stacked generative semi-supervised model

Only a subset of the observations have corresponding class labels.
How to make use of observations without labels?

Latent-feature discriminative model (M1):

VAE $p(z) = \mathcal{N}(z|0, I), p_{\theta}(x|z) = f(x; z, \theta)$

Generative semi-supervised model (M2)

VAE $p(y) = \text{Cat}(y|\pi), p(z) = \mathcal{N}(z|0, I), p_{\theta}(x|y, z) = f(x; y, z, \theta)$

Stacked generative semi-supervised model (M1 + M2)

z_1 from M1, learn model M2 using embeddings from z_1 instead of x .

Two layers of stochastic variables:

$$p_{\theta}(x, y, z_1, z_2) = p(y)p(z_2)p_{\theta}(z_1|y, z_2)p_{\theta}(x|z_1)$$

M1: $p(z_1|x) = \mathcal{N}(z|\mu, \text{diag}(\sigma^2)), p(z_1) = \mathcal{N}(z|0, I)$, decoder: $x \sim p_{\theta}(x|z_1)$

M2: $p(z_2|z_1, y) = \mathcal{N}(z_2|\mu_1, \text{diag}(\sigma_1^2)), p(z_2) = \mathcal{N}(z_2|0, I)$, decoder: $z_1 \sim p_{\theta}(z_1|z_2, y)$

Stacked generative semi-supervised model: M1 loss

M1 loss:

$$p(z|x) = \mathcal{N}(z|\mu, \text{diag}(\sigma^2)), p(z) = \mathcal{N}(z|0, I), \text{ decoder: } x \sim p_\theta(x|z)$$

$$\log p_\theta(x) \geq \mathbb{E}_{q_\theta(z|x)}[\log p_\theta(x|z) - KL(q_\phi(z|x)||p_\theta(z))] = -\mathcal{J}(x)$$

Stacked generative semi-supervised model: M2 Loss

M2 loss:

$$p(z_2|z_1, y) = \mathcal{N}(z_2|\mu_1, \text{diag}(\sigma_1^2)), p(z_2) = \mathcal{N}(z_2|0, I),$$

decoder: $z_1 \sim p_\theta(z_1|z_2, y)$, (x means z_1 , z means z_2) for eqs. below:

$$\begin{aligned}\log p_\theta(x, y) &\geq \mathbb{E}_{q_\phi(z|x, y)}[\log p_\theta(x|y, z) + \log p_\theta(y) + \log p(z) - \log q_\phi(z|x, y)] \\ &= -\mathcal{L}(x, y) \text{ data with labels}\end{aligned}$$

$$\begin{aligned}\log p_\theta(x) &\geq \mathbb{E}_{q_\phi(y, z|x)}[\log p_\theta(x|y, z) + \log p_\theta(y) + \log p(z) - \log q_\phi(y, z|x)] \\ &= \sum_y q_\phi(y|x)(-\mathcal{L}(x, y)) + \mathcal{H}(q_\phi(y|x)) = -\mathcal{U}(x) \text{ data no labels}\end{aligned}$$

$$\mathcal{J} = \sum_{(x, y) \sim \tilde{p}_l} \mathcal{L}(x, y) + \sum_{x \sim \tilde{p}_u} \mathcal{U}(x)$$

$$\mathcal{J}^\alpha = \mathcal{J} + \alpha \cdot \mathbb{E}_{\tilde{p}_l(x, y)}[-\log q_\phi(y|x)]$$

Derive ELBO

KL divergence: $KL(q||p) = \int q \log \frac{q}{p} = \mathbb{E}_q[\log \frac{q}{p}] \geq 0$, $=0$ holds iff $q = p$.

We want $p_\theta(x) \approx p^*(x)$, which is equivalent of $\min KL(p^*(x)||p_\theta(x))$

$$\begin{aligned}\log p_\theta(x) &= \mathbb{E}_{x \sim p^*(\cdot)}[\log p_\theta(x)] = \mathbb{E}_{x \sim p^*(\cdot)}[\log \frac{p^*(x)}{p^*(x)} p_\theta(x)] \\ &= \mathbb{E}_{x \sim p^*(\cdot)} \log p^*(x) + \mathbb{E}_{x \sim p^*(\cdot)}[\log \frac{p_\theta(x)}{p^*(x)}] \\ &= -H(p^*) - KL(p^*(x)||p_\theta(x))\end{aligned}$$

$\max \log p_\theta(x)$ is $\min KL(p^*(x)||p_\theta(x))$ as $H(p^*) = \text{const.}$

Then we need to find a way of calculating $\log p_\theta(x)$

$$\begin{aligned}p_{\theta}(x) &= \int p_{\theta}(x|z)p(z)dz = \int p_{\theta}(x,z)dz \\&= \int \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)}q_{\phi}(z|x)dz = \mathbb{E}_{z \sim q_{\phi}(z|x)} \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \\ \log p_{\theta}(x) &= \log \mathbb{E}_{z \sim q_{\phi}(z|x)} \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \geq \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \right] \\ ELBO &:= \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \right]\end{aligned}$$

You can also see that:

$$\begin{aligned}ELBO &= \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log \frac{p_{\theta}(x,z)}{q_{\phi}(z|x)} \right] = \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log \frac{p_{\theta}(z|x)p_{\theta}(x)}{q_{\phi}(z|x)} \right] \\&= \mathbb{E}_{z \sim q_{\phi}(z|x)} \left[\log \frac{p_{\theta}(z|x)}{q_{\phi}(z|x)} \right] + \mathbb{E}_{z \sim q_{\phi}(z|x)} [\log p_{\theta}(x)] \\&= -KL(q_{\phi}(z|x) || p_{\theta}(z|x)) + \log p_{\theta}(x) \leq \log p_{\theta}(x)\end{aligned}$$

Appendix

$\min KL(p * (x) || p_\theta(x)) \rightarrow \max \log p_\theta(x) \rightarrow \max ELBO$. How to max $ELBO$?

$$\begin{aligned} ELBO &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \frac{p_\theta(x, z)}{q_\phi(z|x)} \right] = \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \frac{p_\theta(x|z)p_\theta(z)}{q_\phi(z|x)} \right] \\ &= \mathbb{E}_{z \sim q_\phi(z|x)} \left[\log \frac{p_\theta(z)}{q_\phi(z|x)} \right] + \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] \\ &= -KL(q_\phi(z|x) || p_\theta(z)) + \mathbb{E}_{z \sim q_\phi(z|x)} [\log p_\theta(x|z)] \\ &\approx -KL(q_\phi(z|x) || p_\theta(z)) + \frac{1}{L} \sum_{l=1}^L \log p_\theta(x|z^{(l)}) \end{aligned}$$

Thus $\max ELBO = \min -ELBO = \min KL(q_\phi(z|x) || p_\theta(z)) + ||x - \hat{x}||^2$

How to calculate $KL(q_\phi(z|x) || p_\theta(z))$?

Appendix One dimension: $z \in \mathbb{R}$

$$q_\phi(z|x) = \mathcal{N}(\mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} e^{\left[-\frac{(z-\mu)^2}{2\sigma^2}\right]}$$

$$\mathbb{E}_{z \sim q_\phi(\cdot|x)}(z) = \mu, \mathbb{E}_{z \sim q_\phi(\cdot|x)}(z^2) = \mu^2 + \sigma^2$$

$$p_\theta(z) = \mathcal{N}(0, 1) = \frac{1}{\sqrt{2\pi}} \exp\left[-\frac{z^2}{2}\right]$$

$$\begin{aligned} KL(q_\phi(z|x) || p_\theta(z)) &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[\log \frac{q_\phi(z|x)}{p_\theta(z)} \right] \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \log \frac{\frac{1}{\sqrt{2\pi}\sigma} e^{\left[-\frac{(z-\mu)^2}{2\sigma^2}\right]}}{\frac{1}{\sqrt{2\pi}} \exp\left[-\frac{z^2}{2}\right]} \\ &= \mathbb{E}_{z \sim q_\phi(\cdot|x)} \left[-\log \sigma - \frac{(z-\mu)^2}{2\sigma^2} + \frac{z^2}{2} \right] \\ &= -\frac{1}{2} [\log \sigma^2 + 1 - \mu^2 - \sigma^2] \end{aligned}$$

Appendix Multi dimension: $\mathbf{z} \in \mathbb{R}^J$

$$\begin{aligned}\int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{z}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I}) d\mathbf{z} \\ &= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (\mu_j^2 + \sigma_j^2)\end{aligned}$$

$$\begin{aligned}\int q_\phi(\mathbf{z}|\mathbf{x}) \log q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} &= \int \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) \log \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}, \boldsymbol{\sigma}^2) d\mathbf{z} \\ &= -\frac{J}{2} \log(2\pi) - \frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2)\end{aligned}$$

Appendix Multi dimension: $\mathbf{z} \in \mathbb{R}^J$

$$\begin{aligned} KL(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z})) &= \mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot|\mathbf{x})} \left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p_\theta(\mathbf{z})} \right] \\ &= \mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot|\mathbf{x})} [\log q_\phi(\mathbf{z}|\mathbf{x})] - \mathbb{E}_{\mathbf{z} \sim q_\phi(\cdot|\mathbf{x})} [\log p_\theta(\mathbf{z})] \\ &= \int q_\phi(\mathbf{z}|\mathbf{x}) \log q_\phi(\mathbf{z}|\mathbf{x}) d\mathbf{z} - \int q_\phi(\mathbf{z}|\mathbf{x}) \log p_\theta(\mathbf{z}) d\mathbf{z} \\ &= -\frac{1}{2} \sum_{j=1}^J (1 + \log \sigma_j^2 - \mu_j^2 - \sigma_j^2) \end{aligned}$$