

Introduction to **Deep Learning**

Variational Autoencoder



Move beyond associating inputs and outputs

Classification: mapping a distribution to a **certain** output



Happy cat



Angry cat

Generation: mapping an input to a probability **distribution**









Relaxed cat

Happy cat



Angry cat



Relaxed cat

What is a Generative Model?



The underlying process of p(x) is unknown (and usually intractable). We approximate p(x) with $\hat{p}_{\theta}(x)$.



Autoencoder



It looks like an autoencoder trivially copies X to X, **but**:

- compressed.
- A good autoencoder selects relevant information and discards irrelevant information.

• The latent vector z has lower dimension than the input X, and therefore information are

Hinton & Salakhutdinov (2006, <u>Science</u>)



Express complex distribution by marginalising simple distributions







х

$$p(x) = \sum_{i=1}^{k} p(\mathbf{x}, h = i)$$



$$p(x) = \int p(\mathbf{x}, \mathbf{h}) d\mathbf{h}$$

Credit: Borealis Al







2D case:

Each value of h is a spherical gaussian A non-linear combination described by $f(h, \theta)$ (blue curve) allows us to approximate a complex distribution p(x).







Computing the posterior



 $p(\mathbf{x} | \mathbf{h}, \theta)$ tells us how to compute the observed distribution given hidden variable **h**. However, we are also interested in what \mathbf{h} is responsible for \mathbf{x} . Therefore, we are interested in $p(\mathbf{h} \mid \mathbf{x})$.



Image credit: Borealis Al







Plug \mathbf{h}^* into $p(\mathbf{x} | \mathbf{h}, \phi)$ Draw \mathbf{h}^* from the prior $p(\mathbf{h})$

Generation

Repeat many times until an approximation of $p(\mathbf{x})$ is properly constructed

Image credit: Borealis Al



The likelihood is:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{h} | \phi) d\mathbf{h}$$
$$= \int p(\mathbf{x} | \mathbf{h}, \phi) p(\mathbf{h}) d\mathbf{h}$$
$$= \int \mathcal{N} \left[f(\mathbf{h}, \phi), \sigma^2 \mathbf{I} \right] \mathcal{N}(\mathbf{0}, \mathbf{1}) d\mathbf{h}$$

More frequently the log likelihood is used.

Problem is: There is no close form for this integral.

Since we can't directly compute p(x), can we at least **constrain its range?**

Evaluating the likelihood

Jensen's Inequality

For a concave function $g[\bullet]$, Jensen's inequality states that $g[\mathbf{E}(y)] \ge \mathbf{E}[g(y)]$

Since log is a concave function, so

$$\log \left[\mathbf{E}(y) \right] \ge \mathbf{E} \left[\log(y) \right]$$

Which means

$$\log\left[\int p(y)dy\right] \le \int p(y)\log(y)dy$$





Evidence Lower Bound (ELBO)

If $q_{\theta}(\mathbf{h})$ is an arbitrary PDF depending on parameters θ ,

$$\log\left[p_{\phi}(x)\right] = \log\left[\int p_{\phi}(\mathbf{x}, \mathbf{h}) d\mathbf{h}\right] = \log\left[\int q_{\theta}(\mathbf{h}) \frac{p_{\phi}(\mathbf{x}, \mathbf{h})}{q(\mathbf{h})} d\mathbf{h}\right] \ge \int q_{\theta}(\mathbf{h}) \log\left[\frac{p_{\phi}(\mathbf{x}, \mathbf{h})}{q_{\theta}(\mathbf{h})} d\mathbf{h}\right]$$
ELBO



By adjusting the ELBO parameters θ , it is possible to find a reasonable approximation to the original PDF.

Image credit: Borealis AI







Evidence Lower Bound (ELBO)

ELBO[
$$\theta, \phi$$
] = $\int q_{\theta}(\mathbf{x} | \mathbf{h}) \log \left[\frac{p_{\phi}(\mathbf{x}, \mathbf{h})}{q_{\theta}(\mathbf{x} | \mathbf{h})} \right] d\mathbf{h}$
= $\int q_{\theta}(\mathbf{x} | \mathbf{h}) \log \left[\frac{p_{\phi}(\mathbf{h} | \mathbf{x}) p_{\phi}(\mathbf{x})}{q_{\theta}(\mathbf{x} | \mathbf{h})} \right] d\mathbf{h}$
= $\int q_{\theta}(\mathbf{x} | \mathbf{h}) \log \left[p_{\phi}(\mathbf{x}) \right] d\mathbf{h} + \int q_{\theta}(\mathbf{x} | \mathbf{h}) \log \left[\frac{p_{\phi}(\mathbf{h} | \mathbf{x})}{q_{\theta}(\mathbf{x} | \mathbf{h})} \right] d\mathbf{h}$
= $\log[p_{\phi}(\mathbf{x})] + \int q_{\theta}(\mathbf{x} | \mathbf{h}) \log \left[\frac{p_{\phi}(\mathbf{h} | \mathbf{x})}{q_{\theta}(\mathbf{x} | \mathbf{h})} \right] d\mathbf{h}$
= $\log[p_{\phi}(\mathbf{x})] - \mathbf{D}_{\mathrm{KL}} \left[q_{\theta}(\mathbf{x} | \mathbf{h}) | p_{\phi}(\mathbf{h} | \mathbf{x}) \right]$
Reconstruction loss

Tightness of the bound (i.e., ELBO gap)

Autoencoder

How can we let an autoencoder to generate **new** contents?

We can sample the latent vector z, but

- We don't know how to sample *z*;
- *z* can be very irregular.



We need to regularize the distribution (i.e., arrangement) of the latent space!

X





Latent Space Distribution







A regular latent space allows us to do **interpolation**. New contents are generated by interpolating the latent space.







Image: Joseph Rocca (TowardDataScience)



Variational Autoencoder

Kingma & Welling (2013, <u>arXiv:1312.6114</u>)



Latent space: tight, smooth, and complete

Without constraining μ and σ :



Solution:

- Use **KL divergence** to force clusters to get **close to each other**;
- Use reconstruction losses (e.g., MSE) to avoid clusters from overlapping.











Latent vector interpolation

Ð ю to D -5

FIXED CONTENT



Bouchacourt, Tomioka & Nowozin (2017).



Application: VAE for Drug Discovery

Proposing candidate molecules by interpolating the latent space.







Fig.2 of Gómez-Bombarelli et al. (2016, arXiv:<u>1610.02415</u>)

