

## Variational Autoencoder

map input to a distribution (rather than a fixed vector, like in autoencoders for dimensionality reduction)

This distribution  $p_{\theta}$  is parametrized by neural network weights  $\theta$ .

Relationship between input  $x$  and latent encode vector  $z$  fully described by:

prior  $p_{\theta}(z)$

likelihood  $p_{\theta}(x|z)$

posterior  $p_{\theta}(z|x)$

If we know  $\theta^*$  (which we learned), to generate samples:

① sample a  $z^{(i)}$  from prior  $p_{\theta^*}(z)$

② generate value  $x^{(i)}$  from conditional  $p_{\theta^*}(x|z=z^{(i)})$

$$\text{where } \theta^* = \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^n p_{\theta}(x^{(i)})$$

For computational tractability, find an approximate function  $q_{\phi}(z|x)$  to output what is likely latent code  $z$  given input  $x$ .

- $p_{\theta}(x|z)$  is generative model (probabilistic decoder)
- $q_{\phi}(z|x)$  is restricted posterior (probabilistic encoder)

We want  $q_{\phi}(z|x)$  to be close to  $p_{\theta}(z|x)$  so

minimize  $D_{KL}(q_{\phi}(z|x) \parallel p_{\theta}(z|x)) \rightarrow$  led to the ELBO objective.

To encourage less blurring, more use of the latent code, focus disentanglement of the latent space, we can do:

- ① introduce a mutual information term: into VAE
- ② allow some Lagrangian weighting parameter in the 2 ELBO terms  $\beta$ -VAE.

ELBO:

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

$\beta$ -VAE:

$$\mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x|z)] - \beta D_{KL}(q_{\phi}(z|x) || p(z))$$

when  $\beta > 1$  is usually chosen, to force more disentanglement in latent  $z$  representation.

The  $\beta$ -VAE objective is closely related to the information bottleneck

$$\max [I(z; Y) - \beta I(x; z)]$$

where goal was to preserve as much information  $I(z; Y)$  about task  $Y$  while reducing rate  $I(x; z)$  of representation.

Constrain latent space to be discrete

# VQ-VAE (vector quantization - VAE)

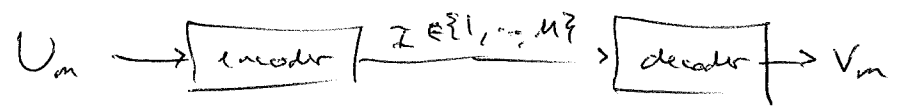
## Intro to Quantization Theory.

Consider an iid sequence of analog random variables

$$U_1, U_2, \dots \sim f_U(u)$$

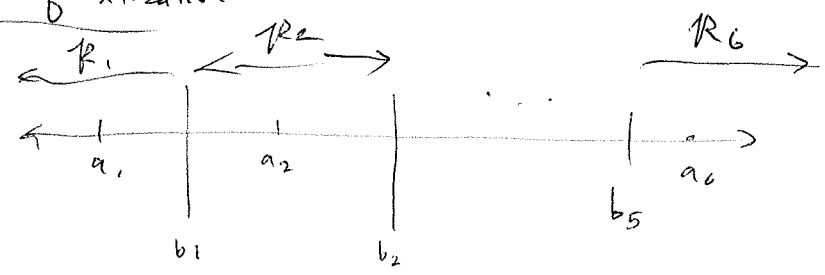
A quantizer maps this sequence into a sequence of discrete r.v.s  $V_1, V_2, \dots$  where  $V_m$  should represent  $U_m$ , for each  $m$  with as little distortion as possible.

If restricted to an alphabet of  $M$  possible values, quantizer output  $V_m$  can't represent losslessly, but more  $M$  is generally better



a common distortion criterion is mean-square error  $E[(U-V)^2]$

### Scalar quantization



- ① given representation points  $\{a_j\}$ , how should intervals  $\{R_j\}$  be chosen?  
nearest neighbor rule:  $b_j = \frac{a_j + a_{j+1}}{2}$
  - ② given intervals  $\{R_j\}$ , how to choose representation  $\{a_j\}$ ?  
centroid:  $a_j = E[U | U \in R_j]$
- Lloyd-Max algorithm alternates between them.  
Trushkin (1982)

Lloyd-Max: given  $M, f_V(u)$

- ① arbitrarily initialize
- ② for each  $j$ , set  $b_j = \frac{a_j + a_{j+1}}{2}$
- ③ for each  $j$ , set  $a_j = \mathbb{E}[U | U \in [b_{j-1}, b_j]]$

repeat until convergence.

Under some convexity conditions (Trushkin, 1982), finds global optimum.

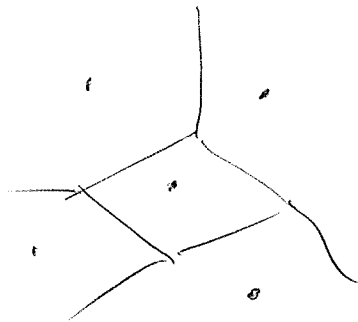
[Sharma, 1978 has dynamic programming formulation].

---

Vector quantization: consider longer vectors. (i.e.  $n=2$ )

① region  $R_j$  must be set of points  $(u, u')$  that are closer to  $(u_j, u'_j)$  than any other representation point.

②  $R_j$  are minimum distance region (i.e. Voronoi regions for given representation points)



Convex polygonal  
regions where boundaries  
are perpendicular bisectors.

② centroids still  $\mathbb{E}[U | U \in R_j]$

Vector Quantization Advantages (Gray & Loutkabaugh, 1989)

- ① memory advantage
- ② space-filling advantage
- ③ shape advantage.

VAE-VAE

using a discrete latent space should prevent "posterior collapse", i.e. blurring.

$$\begin{cases} \text{encoder} & q_\phi(z|x) \\ \text{decoder} & p_\theta(x|z) \\ \text{prior} & p(z) \end{cases}$$

Define latent embedding space as subset of  $\mathbb{R}^D$  with  $K$  representative points:

$$e_i \in \mathbb{R}^D, i=1, 2, \dots, K \text{ which are possible embedding vectors.}$$

The encoder output  $E(x) = z_e$  goes through nearest-neighbor lookup to match one of the  $K$  embedding vectors, and then this is input to decoder.

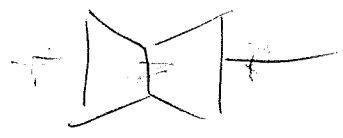
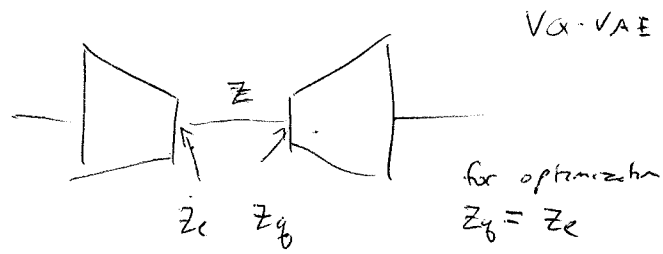
$$z_q(x) = Q(E(x)) = e_k \text{ where } k = \underset{i}{\operatorname{argmin}} \|E(x) - e_i\|_2.$$

Another view (posterior is categorical distribution that is one-hot):

$$q(z=k|x) = \begin{cases} 1 & \text{for } k = \underset{i}{\operatorname{argmin}} \|E(x) - e_i\|_2 \\ 0 & \text{else} \end{cases} = \begin{cases} 1 & \text{for } k = \underset{i}{\operatorname{argmin}} \|z_e(x) - e_i\|_2 \\ 0 & \text{else} \end{cases}$$

For training, since  $\operatorname{argmin}$  is non-differentiable on discrete space, gradient of loss function  $\mathcal{L}$  (e.g. ELBO or something like that)

$$\nabla_z \mathcal{L} \text{ from decoder input } z_q \text{ copied to encoder output } z_e.$$



Consider 3 terms in  $\mathcal{L}$  for VQ-VAE.

- ① reconstruction loss which optimizes decoder, encoder.
- ② due to copying of gradients from  $z_e(x)$  to  $z_d(x)$ , need some way to learn embedding space, i.e. vector quantizer under MSE.
- ③ to ensure encoder commits to an embedding and its output doesn't grow, have "commitment loss".

Let  $sg(\cdot)$  be stop-gradient operator that is defined as identity at forward computation and has zero partial derivatives, effectively constraining its operand to be a non-updated constant.

$$\mathcal{L} = \underbrace{\|x - D(z_k)\|_2^2}_{\text{reconstruction loss}} + \underbrace{\|sg[E(x)] - z_k\|_2^2}_{\text{VQ loss}} + \underbrace{\beta \|E(x) - sg(z_k)\|_2^2}_{\text{commitment loss}}$$

encourages encoder output to stay close to embedding space and prevents it from fluctuating too frequently from one code vector to ~~the~~ another

$$L = \log p(x|z_g(x)) + \|sg[z_e(x)] - e\|_2^2 + \beta \|z_e(x) - sg[e]\|_2^2$$

Also learn a prior over the ~~latent~~ latent discrete codebook so we can sample from it

For generation: decoder can receive input vectors sampled from similar distribution as in training.

[fancy in VO-VAE 2]

output-interpretable VAE (oi-VAE)

introduce structured sparsity into the generator.

- factors generator across observation dimensions with separate generator per group of variable.
  - coupled through shared latent space and joint training

