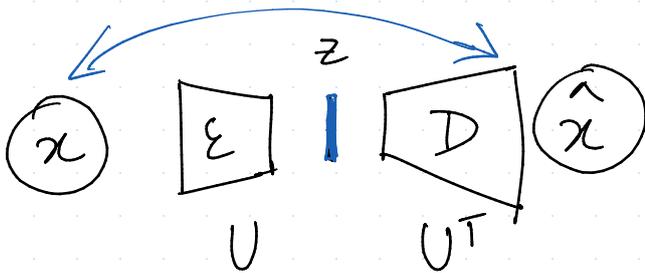# Variational Autoencoders.

## Autoencoders. (A.E.'s)



$\rightarrow$ Meaningful information captured in the latent space.

$z = U^T x$.

$\hat{x} = U z$

$\quad = U \cdot U^T x$.

### Objective function / loss function or Criterion.

$$\min \| \hat{x} - U U^T x \|^2$$

$$\text{s.t.} \quad U U^T = I.$$

\# with one layer neurons and without any non-linearity the A.E. behaves as PCA.

## There are different types of autoencoders.
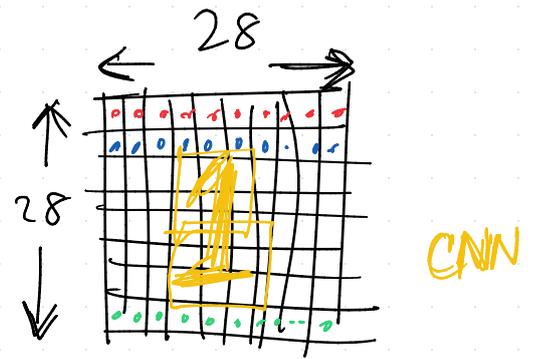
$x \in \mathbb{R}^m$

$z \in \mathbb{R}^n$

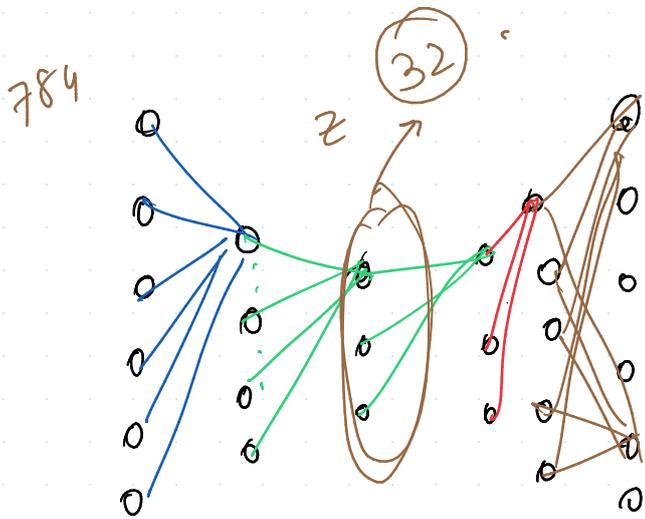$\boxed{m \gg n}$.

$x \boxed{E} z \boxed{D} \hat{x} \simeq x$.

$\rightarrow$ o Undercomplete autoencoders.

# Overcomplete Autoencoder.

$$x \rightarrow \boxed{E} \;|\; \boxed{D} \rightarrow x.$$

$$z$$

$x \in \mathbb{R}^{\textcircled{m}}$

$z \in \mathbb{R}^n$ $\therefore$ $n \gg m.$

28

28

CNN

MNIST.

$\rightarrow$ stacked.

$x \in \mathbb{R}^{\textcircled{1} \times 784}$

$$\boxed{x \in \cdot \mathbb{R}^{784}}$$

784

$z$ $\textcircled{32}$

Using fully connected layers.

$$\| M(x) - x \|_2^2$$

Inception distance.

Img 1          Img 2.

MSE

$\rightarrow$  (Large  MSE.)

ImageNet-1K

In.

MSE

- blurry artifacts.

- edgy artifacts.

$x$ — $E$ — $z$ — $D$ — $\hat{x}$

$+ \, \mathcal{E}$.

Noise term.

meaningful representation

Sparse autoencoders.

LLMs $\rightarrow$ TBs

150B. params. $\rightarrow$ good interpolation.

Variational AutoEncoders $\rightarrow$ Stochastic.

$z \sim$ comes from a certain distribution.

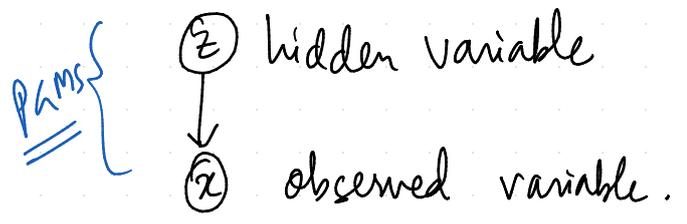$z \sim N(0, I)$.     some gaussian (say)

Variational Inference.

params $\{$

$\textcircled{z}$ hidden variable

$\downarrow$

$\textcircled{x}$ observed variable.

$p(z|x) \rightarrow$ Latent variable.
given an input.

# Probabilistic Graphical Models → representation.

- Topic Modelling (spam/non-spam). 
- Classification
- Encoding to lower dimens.

Decoder (Model).

likelihood. → prior.

→ latent variable.

posterior.

$$p(z|x) = \frac{p(z,x)}{p(x)} = \frac{p(x|z) \cdot p(z)}{p(x)}$$

encoder (Model)

$p(x)$ →

evidence.

intractable quantity

$x \boxed{E} \boxed{\hat{Z}} \boxed{D} x$

$$\frac{p(z|x)}{\downarrow} \qquad \frac{p(x|z)}{Decoder}$$

encoder.

256 28.

256.

28

$(256)^{28 \times 28}$ → $10^{1889} \approx 10^{940}$

$10^5$ ⇒ $0.000 \cdots 1\%$

MNIST - sample?

$1884 \rightarrow 0's$

$70K$ → $100 000$ →

Grayscale → (0 - 255).

# Variational Inference

Turn this intractable quantity to an optimization problem, by assuming there is another distribution which is tractable. Now, find the parameters of that distribution that is very close to this one. That distribution is used as a surrogate to the current intractable distribution.

$q_\theta(z) \leftarrow$ comes from a well behaved family of distributions (e.g. Gaussian).

$$\min \ KL \left( q(z) \| p(z|x) \right) \quad \Leftarrow \text{Minimize the}$$
KL divergence b/w these two distributions.

known.    unknown.

$\approx q_\theta(z|x) \approx q_\theta(z)$.     Since $\theta$ is dependent on $n$.
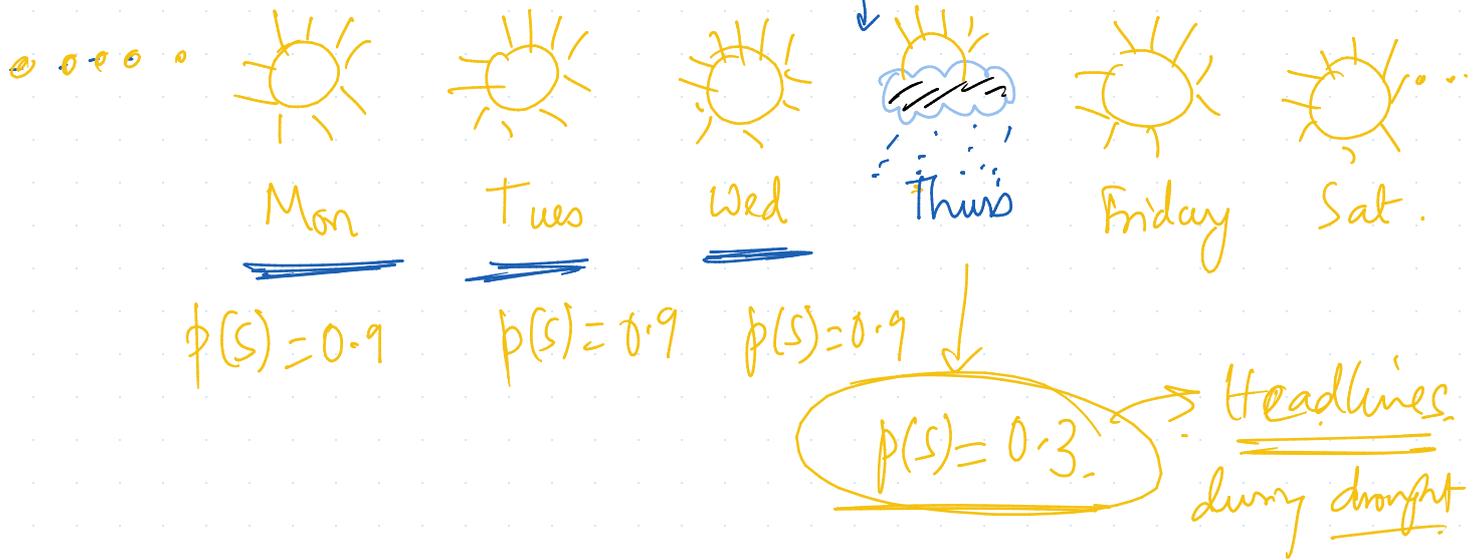
$x \in X$     $G(0, I)$.

# Information

$$I = -\log(p(x))$$

$x \rightarrow$ event.

$p(x) \rightarrow$ probability of that event occuring ( not occuring )

(Higher) probability means (lower) information.

new information.

| Mon | Tues | Wed | Thurs | Friday | Sat. |

$p(s) = 0.9$     $p(s) = 0.9$     $p(s) = 0.9$

$p(s) = 0.3$   $\rightarrow$ Headlines during drought

In the top right:

$x \rightarrow \boxed{\varepsilon} \,\Big|^{z}\, \boxed{D} \rightarrow x.$

$\varepsilon \sim p(z|x)$        $D \sim p(x|z).$

$q_\theta(z|x) \sim q(z).$

$z \sim G(0, I).$

Entropy $\rightarrow$ Expectation of the information.

$$H = E[I] = -\sum p(x) \log(p(x))$$

$$\boxed{E[x] = \sum x\, p(x)}$$

<u>KL - divergence → more like</u> :: Entropy of $p$ − Entropy of $q$

$$KL(p||q) = -\sum \boxed{p(x)} \log p(x) + \sum q(x) \log q(x)$$

But in KL we compute the expectation w.r.t. certain quantities, like eg., if the expectation is w.r.t. $q$, then it is KL divergence.

$$KL(q||p) = -\sum q(x) \log p(x) + \sum q(x) \log q(x).$$

$$\boxed{KL(q||p) = -\sum q(x) \log \frac{p(x)}{q(x)}}$$

KL can also be written as the information loss if we want to transfer from one distribution to another, hence, this is a measure b/w two distributions.

property $\begin{cases} \bullet\ KL(p||q) \neq KL(q||p) \rightarrow \text{hence} \\ \qquad\qquad\qquad\qquad\qquad \text{divergence & not distance} \\ \bullet\ KL(p||q)\ \text{or}\ KL(\cdot||\cdot) \geqslant 0 \end{cases}$

Hence KL is the measure of dissimilarity b/w the two distributions.

So we are minimizing the KL divergence b/w $q_\theta(z)$ & $p(z|x)$ → intractable. Here $q_\theta(z)$ → is from a family of well behaved distribution.

$$\min_\theta \quad KL\left(q_\theta(z) \| p(z|x)\right) \quad \rightarrow \text{continuous form.}$$

$$\Rightarrow -\int_z q(z) \log \frac{p(z|x)}{q(z)} \, dz.$$

$$\Rightarrow -\int_z q(z) \log \frac{p(z,x)}{q(z) \cdot p(x)} \, dz$$

$$\Rightarrow -\int_z q(z) \log \left\{ \frac{p(z,x)}{q(z)} \cdot \frac{1}{p(x)} \right\} \, dz$$

$$\Rightarrow -\int_z q(z) \log \frac{p(z,x)}{q(z)} \, dz + \int_z q(z) \log p(x) \, dz$$

$$\log(A \cdot B) = \log A + \log B.$$

$$\Rightarrow - \int_z q(z) \log \frac{p(z,x)}{q(z)} dz + \underline{\log p(x)} \underbrace{\int_z q(z) dz}_{1.}$$

Since we are integrating on z and $p(x)$ observation which is a constant, and it doesn't depends on z or anything, hence we are taking it out. $\longrightarrow$ Nothing to do <u>with</u> <u>$\theta$ either.</u>

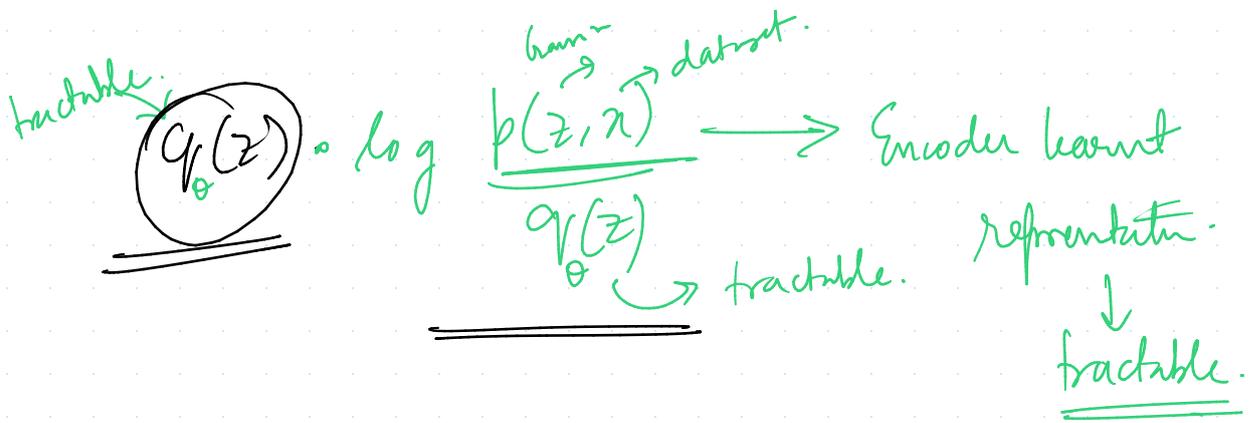$$\min \; KL\left(q(z) \| p(z|x)\right) = - \int q(z) \log \frac{p(z,x)}{q(z)} + $$

maximize this quantity

$$\underbrace{\log p(x)}$$
constant, tractable
from dataset.

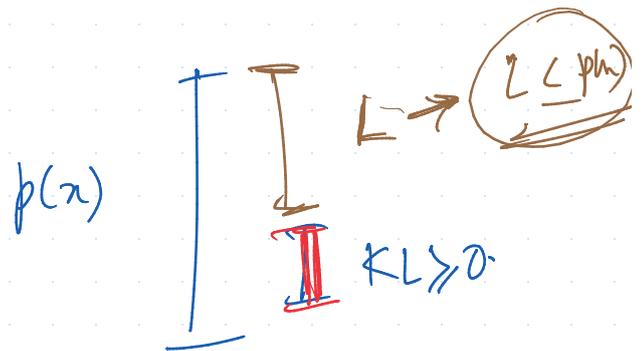ELBO - Evidence Lower Bound.

VLB - Variational Lower Bound.

$q_\theta(z)$ $\circ$ log $\dfrac{p(z, x)}{q_\theta(z)}$ $\longrightarrow$ Encoder learnt

learn → dataset.

tractable.

representation.

↓

tractable.

$$\log p(x) = KL\left(q(z) \| p(z|x)\right) + \int q(z) \log \frac{p(z, x)}{q(z)}.$$

Constant.                    $\geq 0$                    $\mathcal{L}$ or lower bound.

$$\mathcal{L} \neq \log p(x) \quad \text{unless} \quad KL = 0.$$

Hence, $\mathcal{L}$ is the lower bound of the $p(x)$.

$p(x)$ $\mathcal{L} \rightarrow$ $\mathcal{L} < p(x)$

$KL \geq 0$

## Lower Bound.

$$\mathcal{L} = \int q(z) \log \frac{\boxed{p(z,x)}}{q(z)}$$

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

$$p(x,z) = p(x|z) \cdot p(z)$$

$$= \int q(z) \log \frac{p(x|z) \cdot p(z)}{q(z)}$$

← Decoder.

↖ well defined Gaussian.

$$= \int q(z) \log p(x|z) + \int q(z) \log \frac{p(z)}{q_\theta(z)}$$
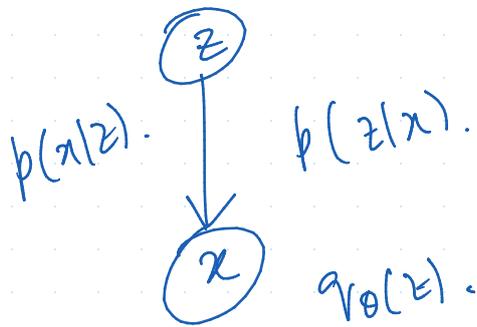
maximize $\mathcal{L}$.

$$- KL\left(q(z) \| p(z)\right).$$

$$= \max \int q(z) \log p(x|z)$$

$$= \max \, E\left[\log p(x|z)\right]. \longrightarrow \text{likelihood of the dataset.}$$

$$\log p(x) = KL\left(q(z) \| p(z|x)\right) + E\left[p(x|z)\right] - KL\left(q(z) \| p(z)\right)$$



$p(x|z).$     $p(z|x).$

$q_\theta(z).$

Maximize likelihood:-

→ Gaussian: minimize MSE

→ Bernoulli — minimize CE loss.

Gaussian

$$|x - \hat{x}| + KL\left(q(z) \| N(0,0)\right)$$

$\underbrace{\phantom{|x-\hat{x}|}}_{AE}$     $\underbrace{\phantom{KL(q(z)\|N)}}_{VAE}$
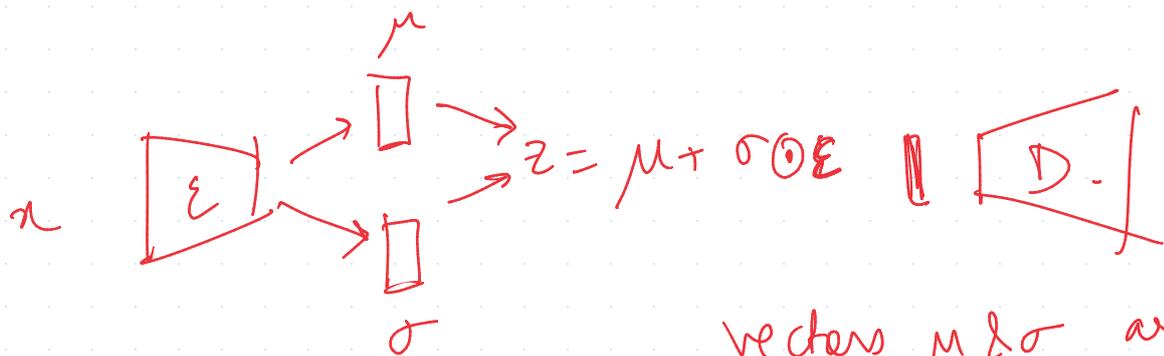
This additional loss in VAEs ensures that the z is Gaussian, since z → stochastic, hence no backpropagation.
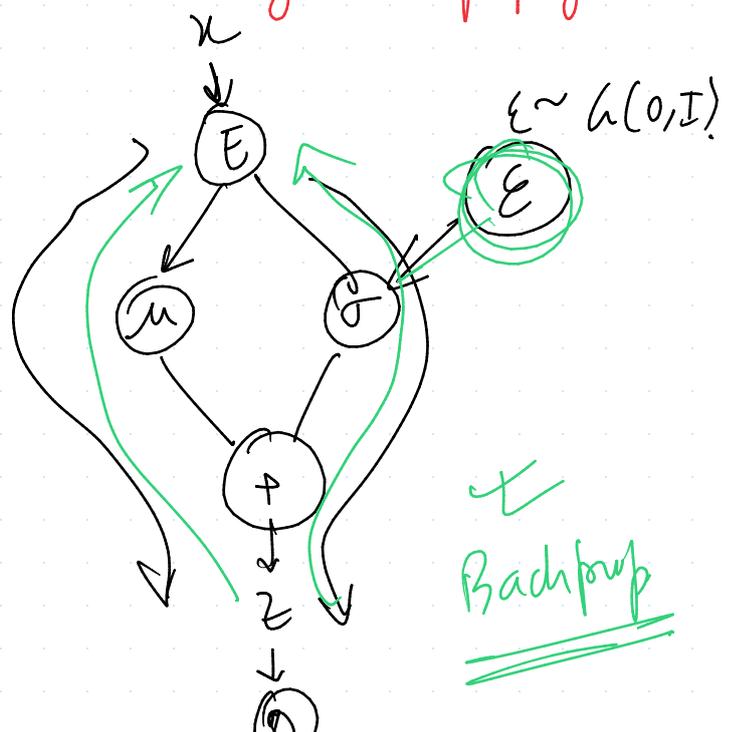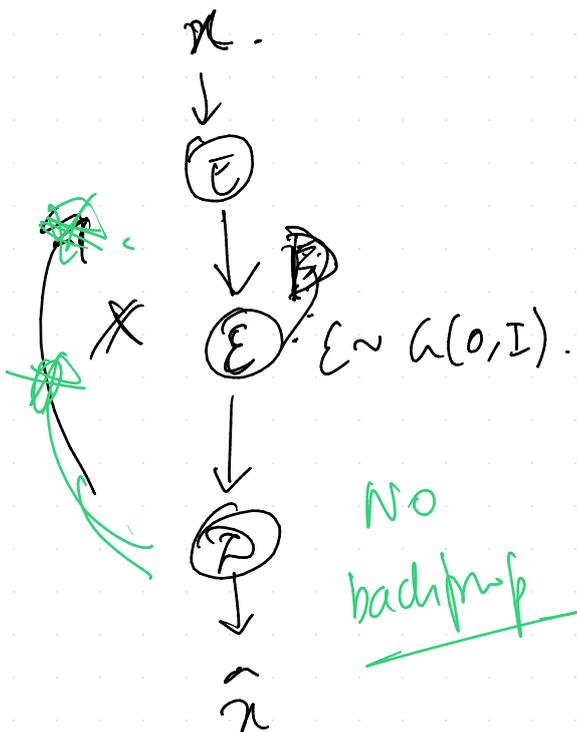
Reparameterization trick. :- find the mean & variance
of the dist. via the neural network.

(mean, variance) $\rightarrow$ deterministic.

Through this Gaussian - sample something random,
representation of $z \rightarrow$ parameters of $z$ in the model.



$$z = \mu + \sigma \odot \varepsilon$$

vectors $\mu$ & $\sigma$ are learnt
using backpropagation.

$\varepsilon \sim \mathcal{N}(0, I)$

$\varepsilon \sim \mathcal{N}(0, I)$.

No
backprop

Backprop

$$\tilde{x}$$

- Inpainting .

- Grayscale — colored images (colorization).