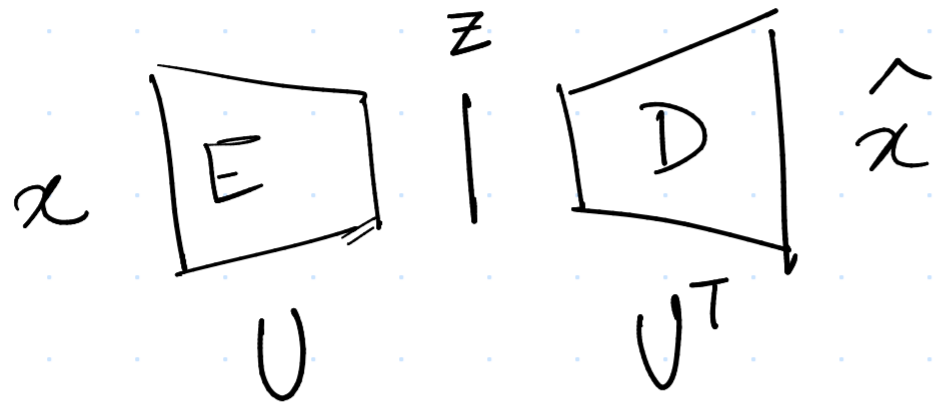


# VAEs (Variational Auto-Encoders)

26/02/2025

## Autoencoders (A.E's)

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



$$z = U^T x$$

$$\hat{x} = U z$$

$$= U U^T x.$$

Objective function :-

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

s.t.  $U \cdot U^T = I.$

$I I^T = I$  ✓

Autoencoders → used for compression. → undercomplete autoencoder.

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

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

$m \gg n.$

$$x \begin{bmatrix} E \end{bmatrix} z \begin{bmatrix} D \end{bmatrix} \hat{x} \approx x.$$

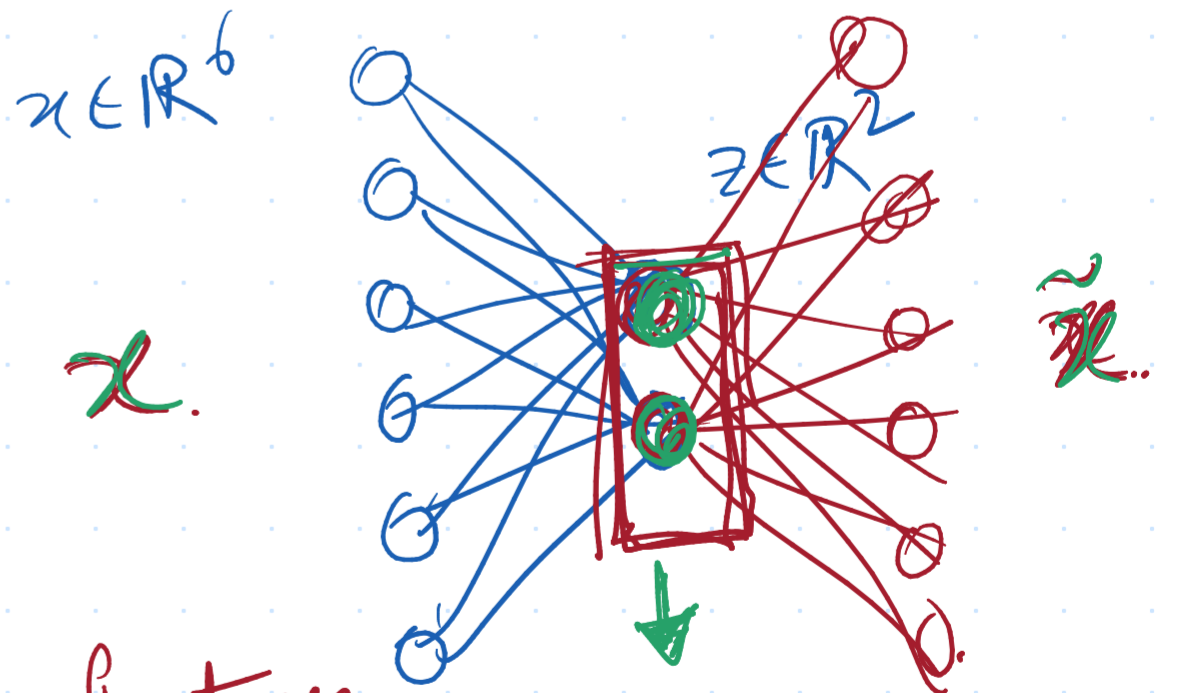
undercomplete autoencoder.



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

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

$n \gg m.$



redundant features

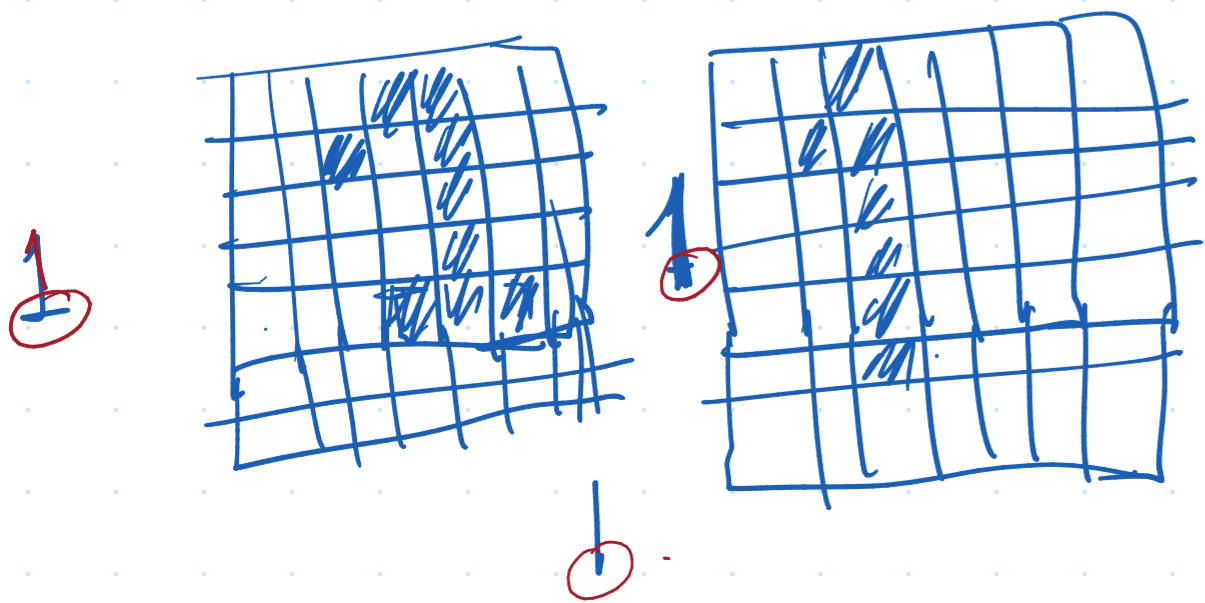
$$28 \begin{bmatrix} 1 \end{bmatrix} \rightarrow \textcircled{1}$$

$$x'$$

$$y' \in \{0, 1, 2, \dots, 9\}$$

$$28 \times 28 \Rightarrow 784.$$

7x7



$$\begin{array}{r} 28 \\ \times 28 \\ \hline 224 \\ 560 \\ \hline 784 \end{array}$$

$x \in \mathbb{R}^{m \times m}$  higher dimensional.

$z \in \mathbb{R}^n$  where  $m \gg n$ , is always lower dimensional.

□ □ □ □ → no labels.



MSE loss.

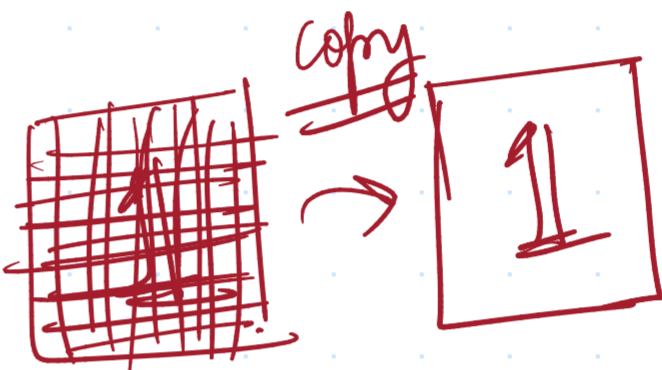
L1 loss.

for similarity

(5, 7, 3, 2, 1)

(5, 7, 5, 2, 2) (5, 7, 3, 2, 1)

(5, 7, 1, 2, 1)



$$\cos \theta = \frac{a \cdot b}{\|a\| \|b\|} = \frac{25 + 49 + 9 + 4 + 1}{\sqrt{25 + 49 + 9 + 4 + 1}} = \frac{1}{1} = 1$$

$$\theta = \cos^{-1}(1) = 0$$

[-1 to 1]

(1, 2, 3) (4, 5, 6)

$$\Rightarrow \sqrt{(4-1)^2 + (5-2)^2 + (6-3)^2} \approx \text{high}$$

but cosine similarity = 0.

$$A = [1, 0]$$

$$B = [-1, 0]$$

$$x = [3, 4]$$

$$\hat{x} = [0; 0]$$

$$\theta = 90$$

$$\underline{\underline{\text{Cosine}}}$$
  
$$\underline{\underline{(-1, 1)}}$$

$$\underline{\underline{\text{MSE}}}$$
  
$$\underline{\underline{(0, \infty)}}$$

$$x = (3, 4)$$
  
$$\hat{x} = (-4, 3)$$

$$\underline{\underline{L1}}$$
  
$$\underline{\underline{(0, \infty)}}$$

$$\text{C.S.} = \frac{x \cdot \hat{x}}{\|x\| \|\hat{x}\|} =$$

$$\frac{-24}{25} = \underline{\underline{-0.96}}$$

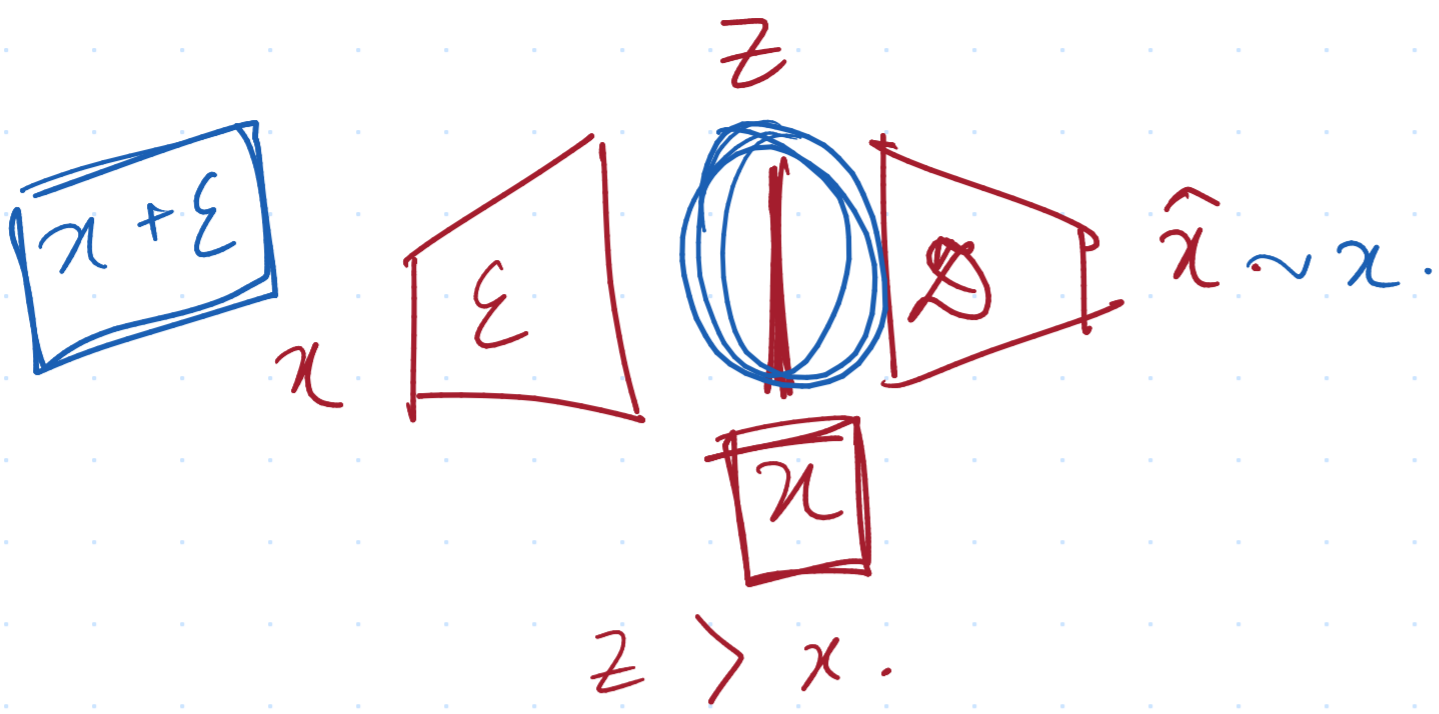
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{x}_i)^2$$

$$\sqrt{(4)^2 + (3)^2}$$

$$= \frac{1}{2} \left( (3 - (-4))^2 + (4 - (-3))^2 \right) = \frac{49 + 49}{2}$$

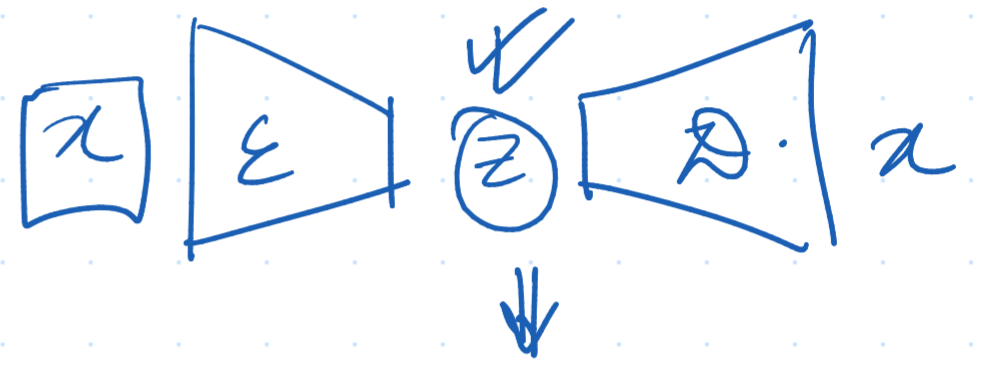
$$= \underline{\underline{49}}$$

MSE, L1 good for better representation.



$$\|x - \hat{x}\|_2$$

identity mapping is being learnt.



Noise Removal.

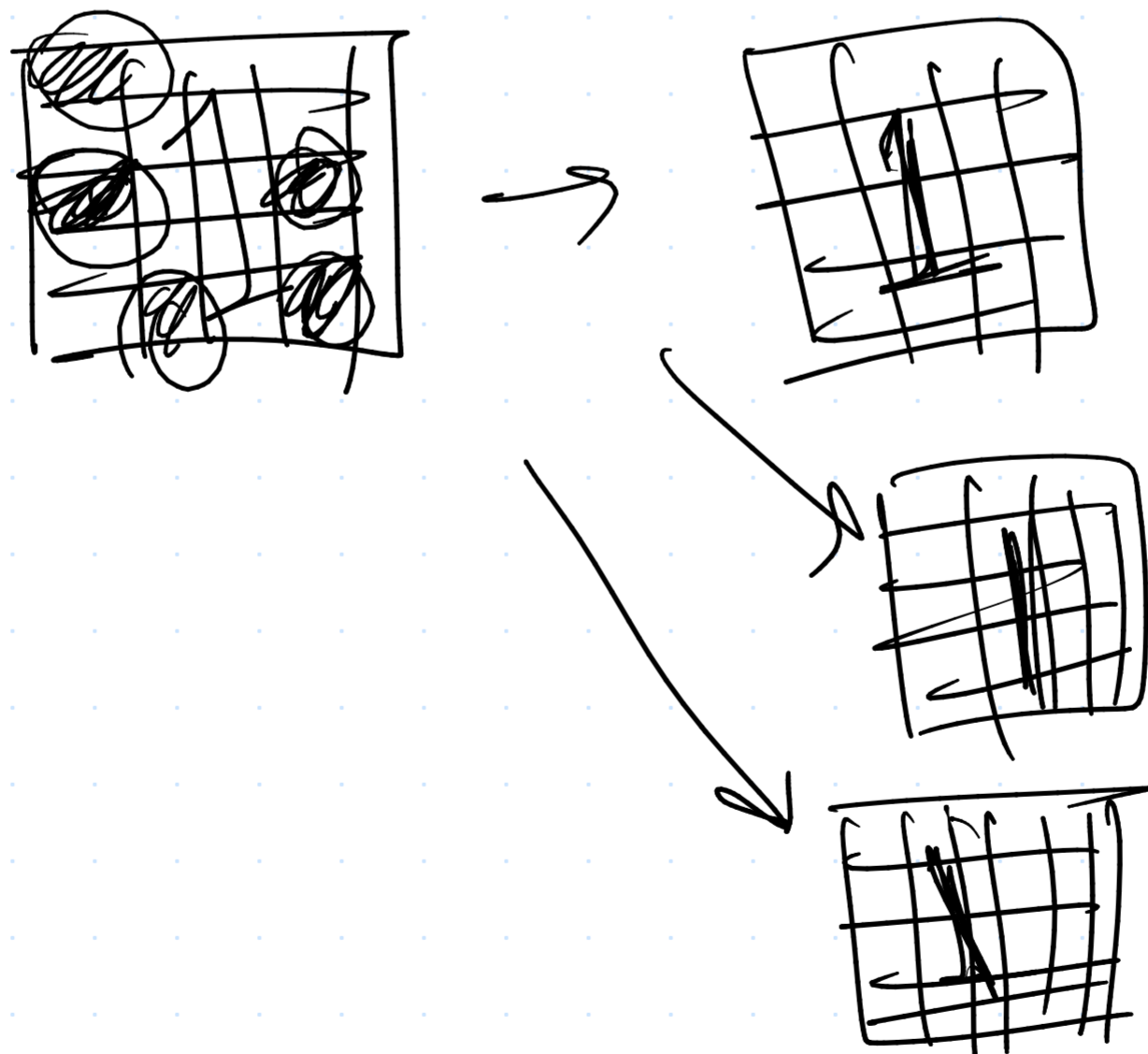


Denoising autoencoder.

$$x + \epsilon$$

identity mapping.

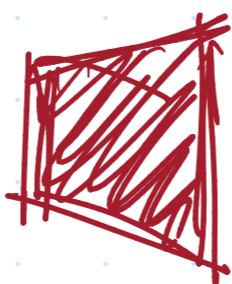
$$U = I$$



100%

outliers.

95% → wood image + Noise



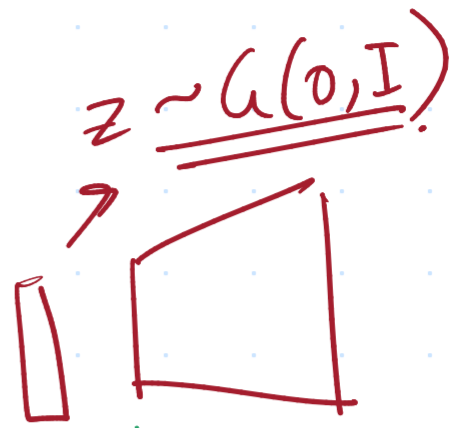


# Variational AE - stochastic.

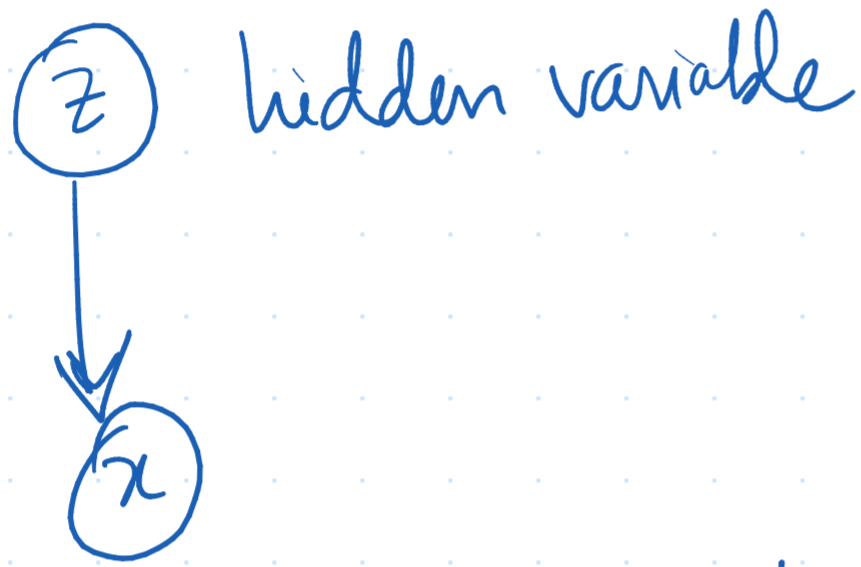
$z \rightarrow$  comes from a certain distribution.

$z \sim \mathcal{N}(0, I)$  some Gaussian, say.

o o o o o o o o



Variational Inference  $\rightarrow$  Bayesian Statistics



$p(z|x)$   $\rightarrow$  latent var. given an input.

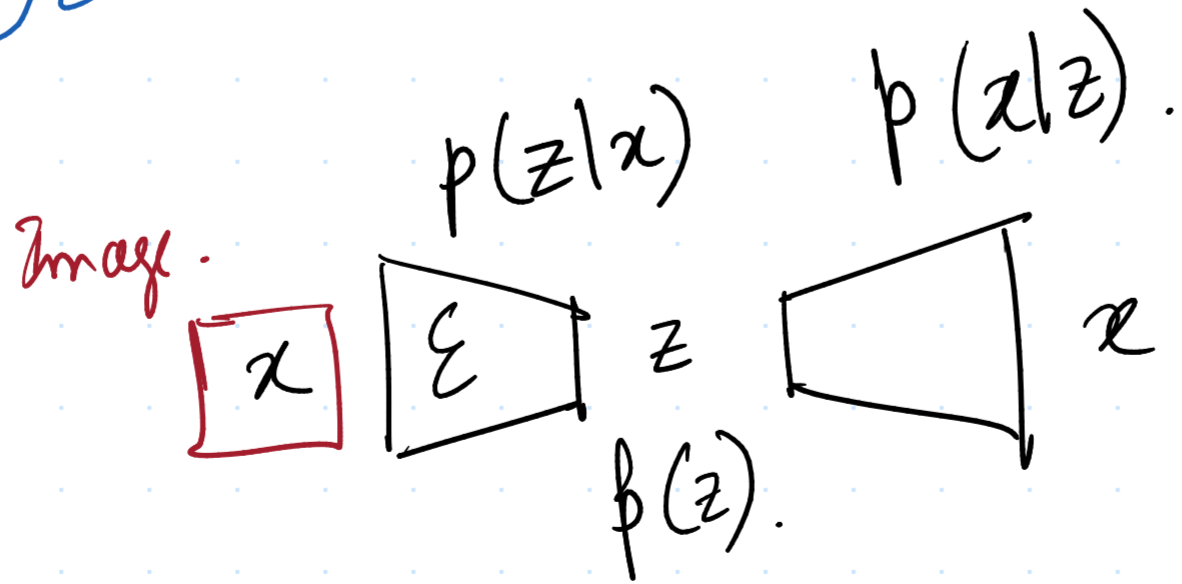
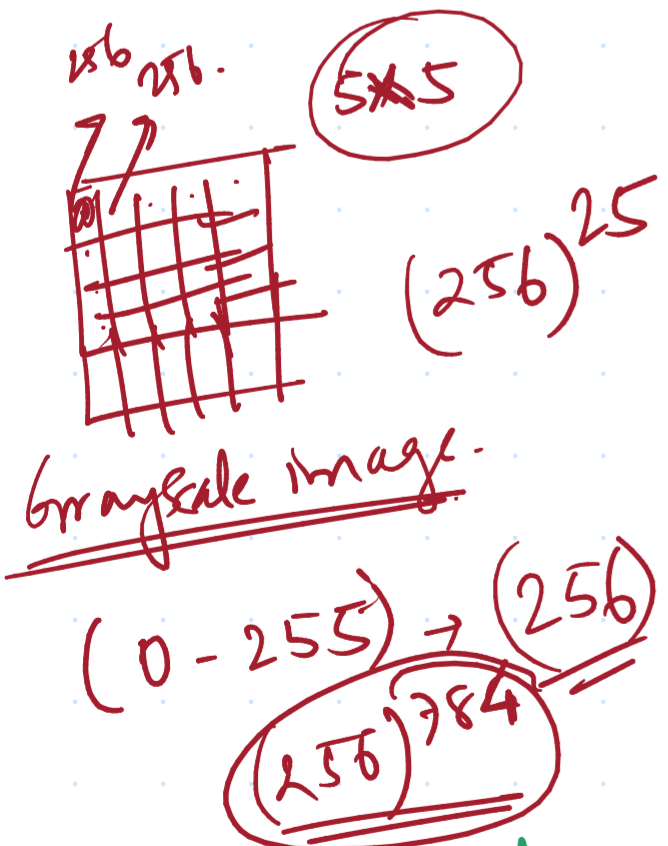
- Topic Modelling
- Classification
- Encoding to lower dim, etc.

$$p(z|x) = \frac{p(x|z) p(z) \rightarrow \text{prior}}{p(x)}$$

$$\int_z p(x|z) \cdot p(z) dz$$

Intractable quantity

$\rightarrow$  when  $z \rightarrow$  higher dimension, then



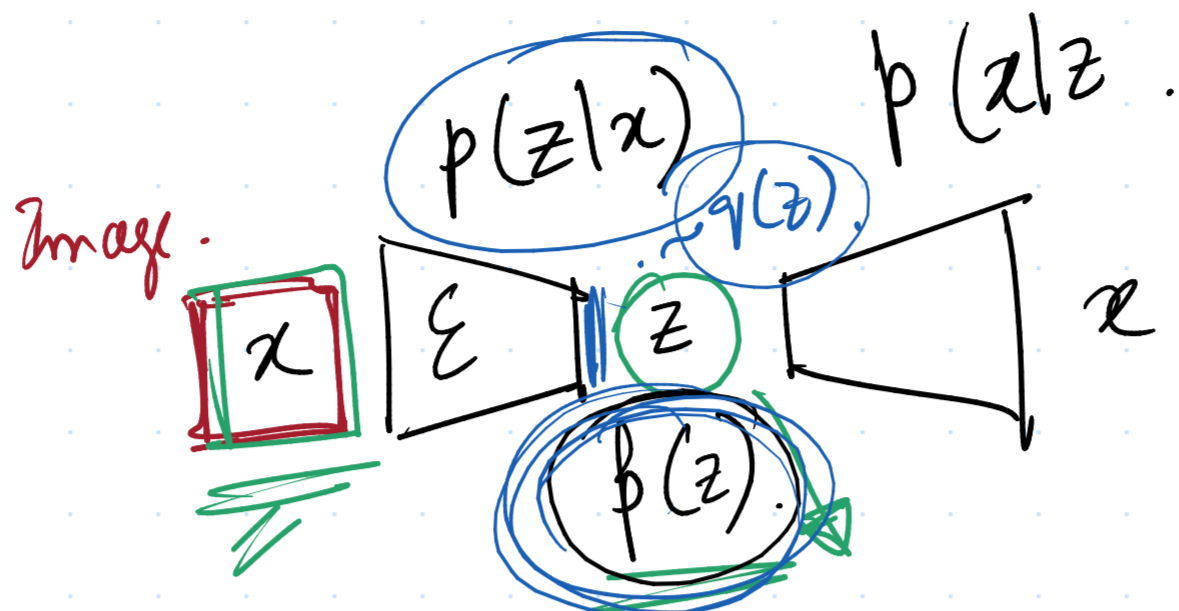
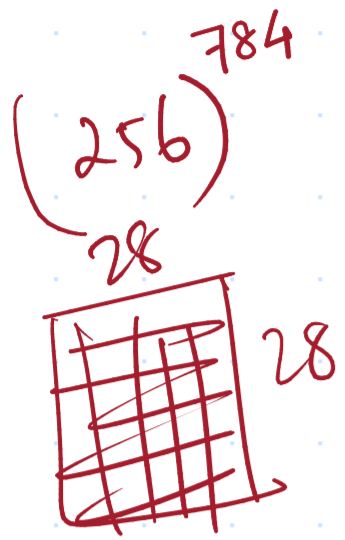
$$\int \int \int \dots$$

$z_1 z_2 z_3 z_4$  this

Sampling techniques  $\rightarrow$  one time of technique to address this intractability.

becomes complicated integral.

$$p(z|x) = \frac{p(x|z) p(z)}{p(x)} \rightarrow \text{prior.}$$

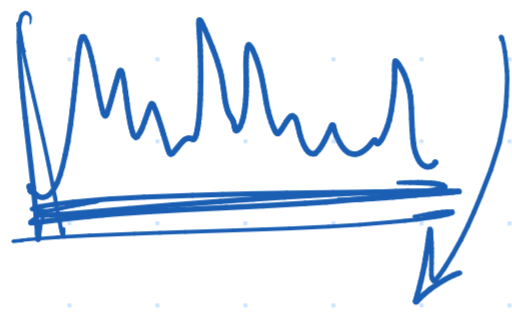


so this is intractable

since we are not getting  $x$  fully like  $0.00 \dots 01\%$  of the 1998 of 0's

(1998) 0's  
7 (1998)

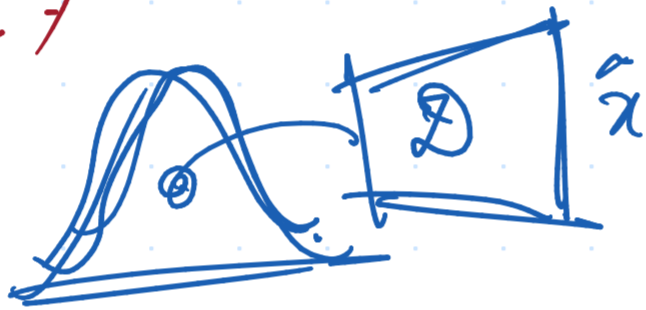
60K  
10-1998



$$q_0(z) \sim \mathcal{G}(0, \Sigma)$$

real data in even  $28 \times 28$  pixel which is barely recognizable & even grayscale (not even colored)

$0.00 \dots 1\%$



blurry average.

$$\hat{x} \sim x.$$

## ⊗ 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.

$q_0(z) \leftarrow$  comes from a well behaved family of distributions. (like gaussian)

$$\min KL(q(z) \parallel p(z|x))$$

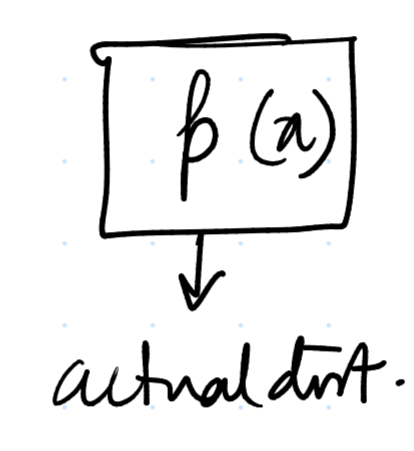
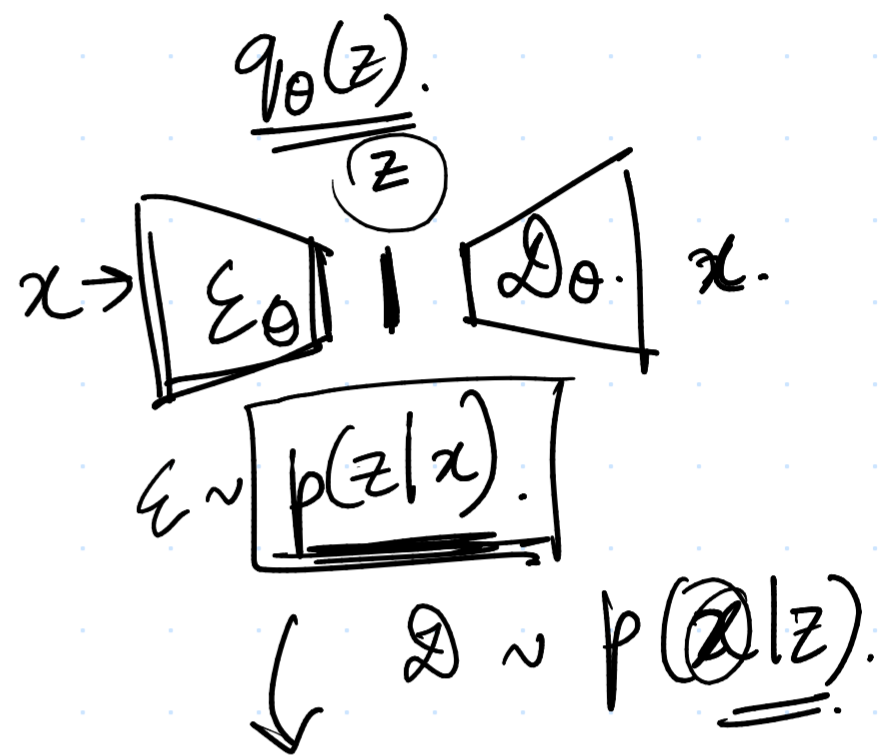
known.
unknown.

← Minimize the KL divergence b/w these two distributions

$$\approx q_{\theta}(z|x) \approx q_{\theta}(z)$$

$x \in X$ 
 $\sim \mathcal{U}(0,1)$

since  $\theta$  is dependent on  $x$ .



$$q_{\theta}(z|x) \approx q(z)$$

$$z \sim \mathcal{U}(0,1)$$

Information

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

$x \rightarrow$  event.

Higher probability means lower information

Entropy  $\rightarrow$  Expectation of the information.

$$E(x) = \sum x p(x)$$

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

KL-divergence  $\rightarrow$  more like: Entropy of  $p$  - Entropy of  $q$ .

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

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



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

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

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

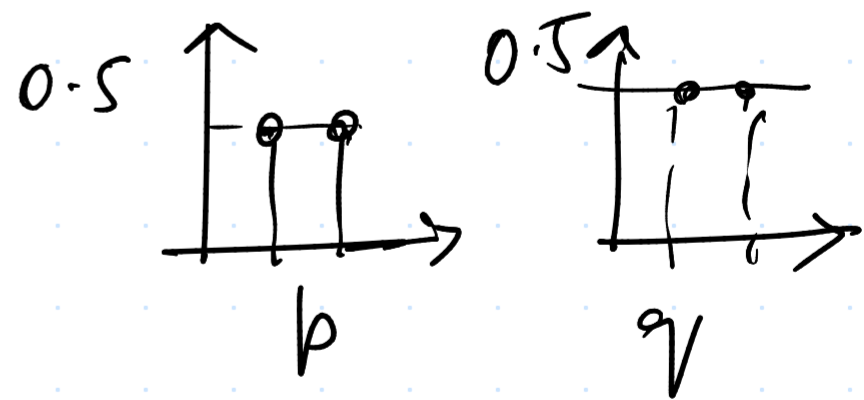
- Property of KL divergence
- $KL(p||q) \neq KL(q||p)$ .  $\rightarrow$  hence divergence & not distance
  - $KL(p||q)$  or  $KL(\cdot||\cdot) \geq 0$ .

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

$$KL(p||q)$$

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

$$= -\sum q \log \frac{0.5}{0.5} = -\sum q \log(1) = 0$$

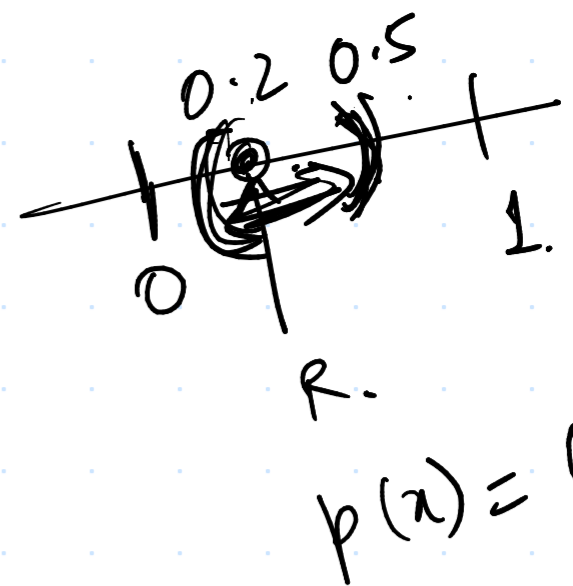


So, we were minimizing the KL divergence b/w  $q_\theta(z)$  and  $p(z|x)$   $\rightarrow$  intractable. Here  $q_\theta(z) \rightarrow$  well behaved family



q distribution.

$$\min_{\theta} KL(q_{\theta}(z) || p(z|x))$$



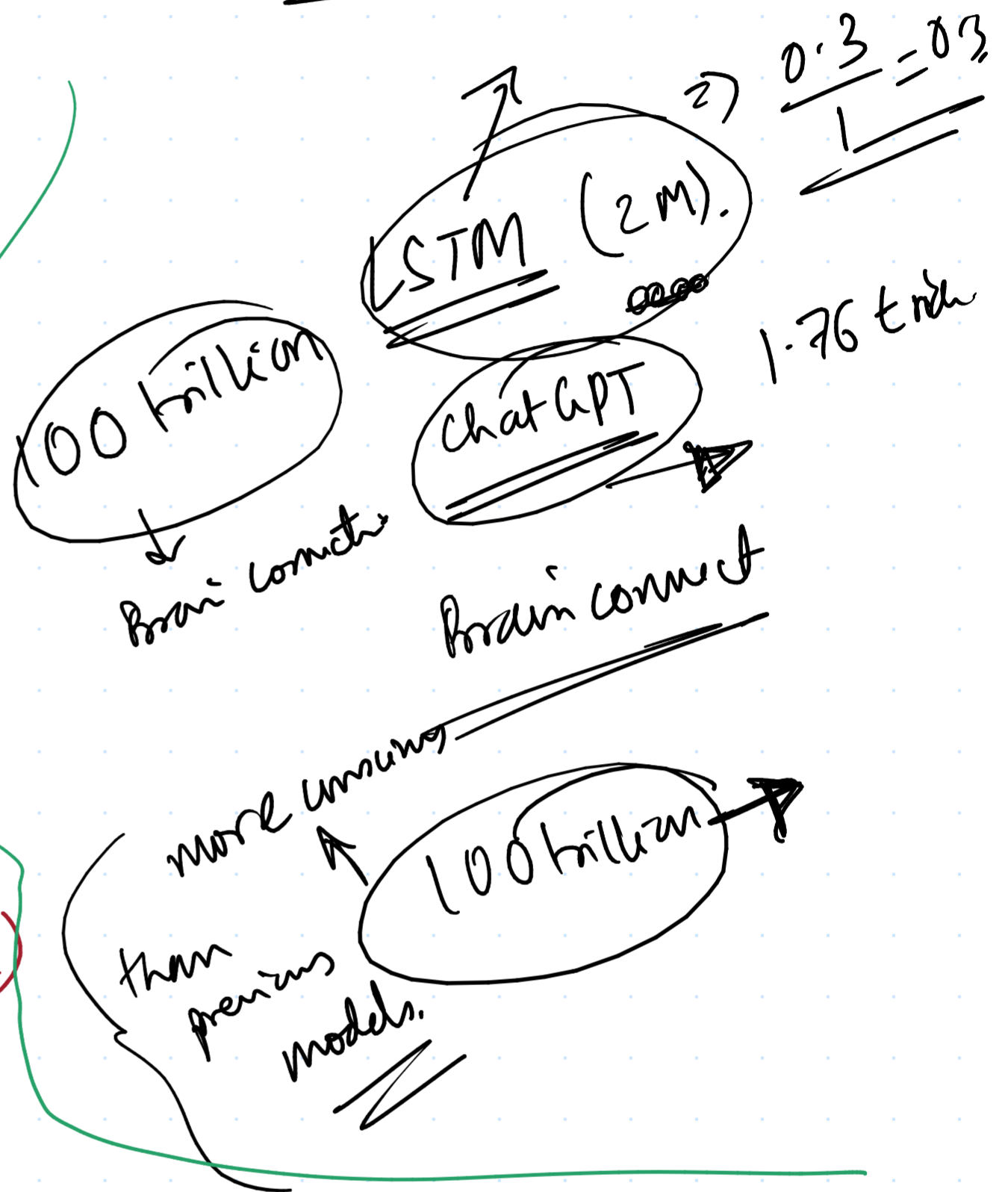
KL(q(z) || p(z|x)) → KLD in continuous space

$$= - \int_z q(z) \log \frac{p(z|x)}{q(z)}$$

$$= - \int_z q(z) \log \frac{p(z,x)}{q(z)p(x)}$$

$$= - \int_z q(z) \log \frac{p(z,x)}{q(z)} \cdot \frac{1}{p(x)}$$

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



$$\therefore \min KL(q(z) || p(z|x)) = - \int q(z) \log \frac{p(z|x)}{q(z)}$$

Since we are integrating on  $z$  and  $p(x)$  is observation, which is a constant, and it doesn't depend on  $z$  or anything, hence we are taking it out. → nothing to do with  $\theta$  either.

$$\min_{q(z)} \text{KL}(q(z) \parallel p(z|x)) = - \int q(z) \log \frac{p(z|x)}{q(z)} + \log p(x).$$

$\underbrace{\log p(x)}_{\text{Constant.}} \downarrow \text{tractable.}$

maximize this quantity.

ELBO - Evidence Lower Bound.

VLB - Variational Lower Bound.

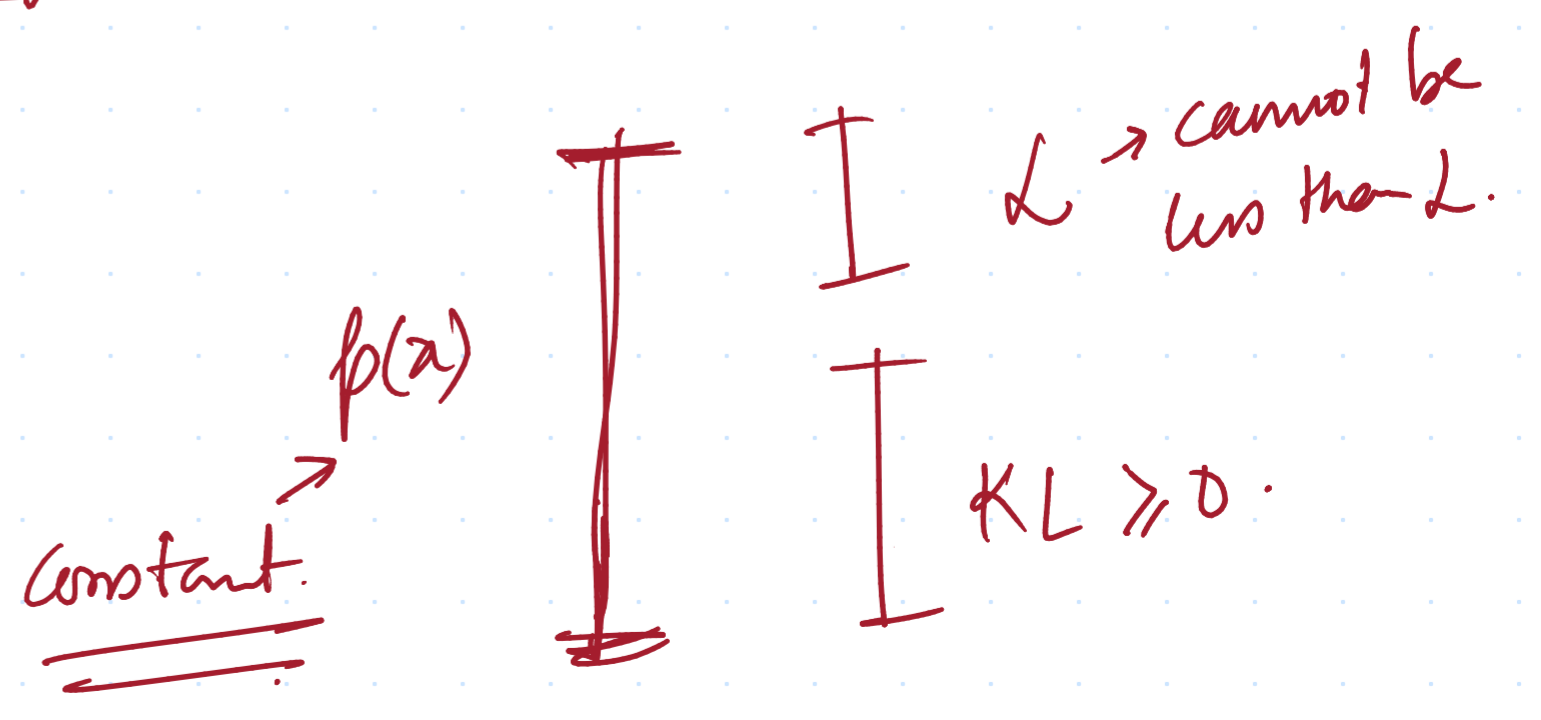
$q(z)$  tractable.  $\log \frac{p(z|x)}{q(z)}$  tractable.  $\rightarrow$  Encoder learnt.  
 $q(z) \leftarrow$  tractable, gaussian

$$\log p(x) = \text{KL}(q(z) \parallel p(z|x)) + \int q(z) \log \frac{p(z|x)}{q(z)}$$

$\log p(x)$  Constant.  
 $\text{KL}(q(z) \parallel p(z|x)) \geq 0$   
 $\int q(z) \log \frac{p(z|x)}{q(z)}$  L or lower bound.

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

Hence,  $L = \text{lower bound of this } \log p(x)$ .



Lower Bound :-

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

$q(z)$

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

$$= \int q(z) \log \frac{p(x|z) p(z)}{q(z)} \leftarrow \text{well defined Gaussian.}$$

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

$$= \underbrace{\int q(z) \log p(x|z)}_{=} + \underbrace{\int q(z) \log \frac{p(z)}{q(z)}}_{\substack{\text{well behaved gaussian} \\ \text{distribution.}}}$$

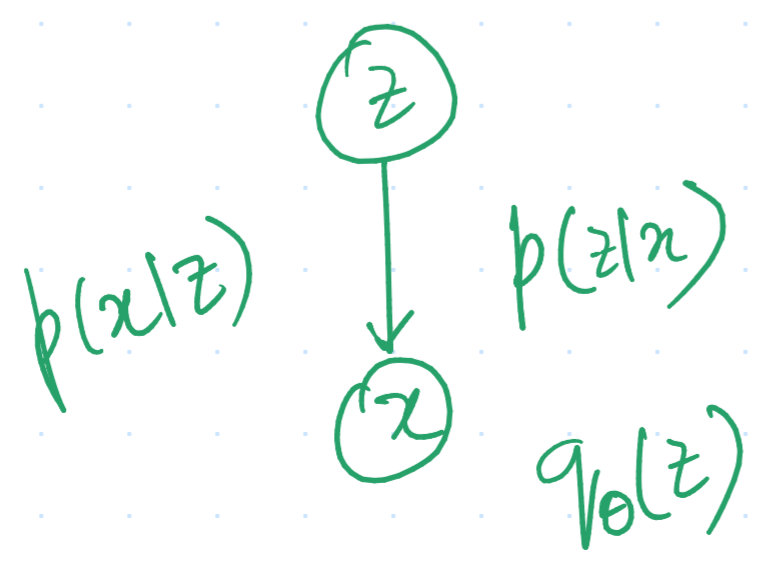
$\rightarrow$  KL( $q||p$ ).

maximize  $\mathcal{L}$

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

$$= \max E[\log p(x|z)] \rightarrow \text{likelihood of the data.}$$

$$\log p(x) = \text{KL}(q(z) || p(z|x)) + E[\log p(x|z)] = \text{KL}(q(z) || p(z))$$



Maximize Likelihood :-

$\rightarrow$  Gaussian - minimize MSE

$\rightarrow$  Bernoulli  $\rightarrow$  minimize CE loss.

Gaussian :-

$$\underbrace{|x - \hat{x}|}_{A \cdot E} + \text{KL} \left( q(z) \parallel \underbrace{N(\cdot, \cdot)}_{\text{VAE}} \right)$$

This additional loss in VAEs ensures that the  $z$  is Gaussian.

Since  $z \rightarrow$  stochastic hence no backpropagation.

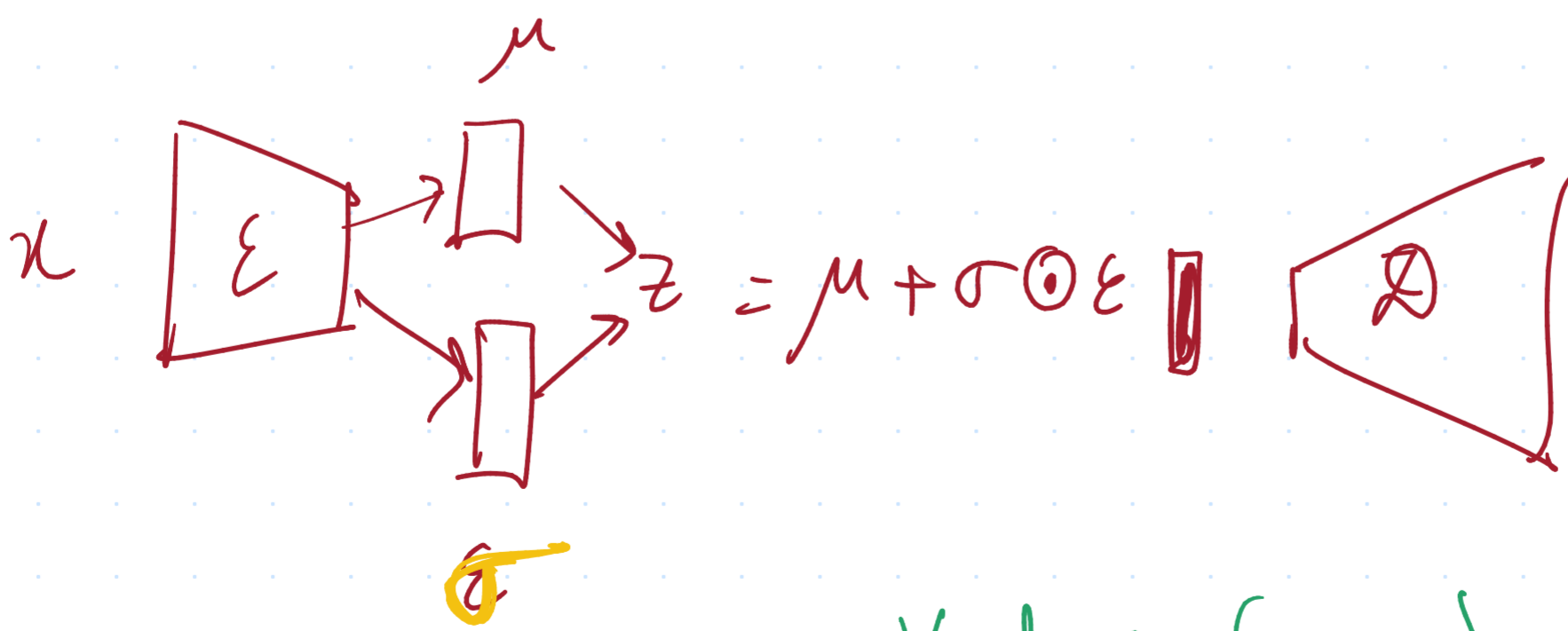
Reparameterization Trick :- Find the mean & variance of the dist. via the neural network.

(mean, variance)



deterministic.

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



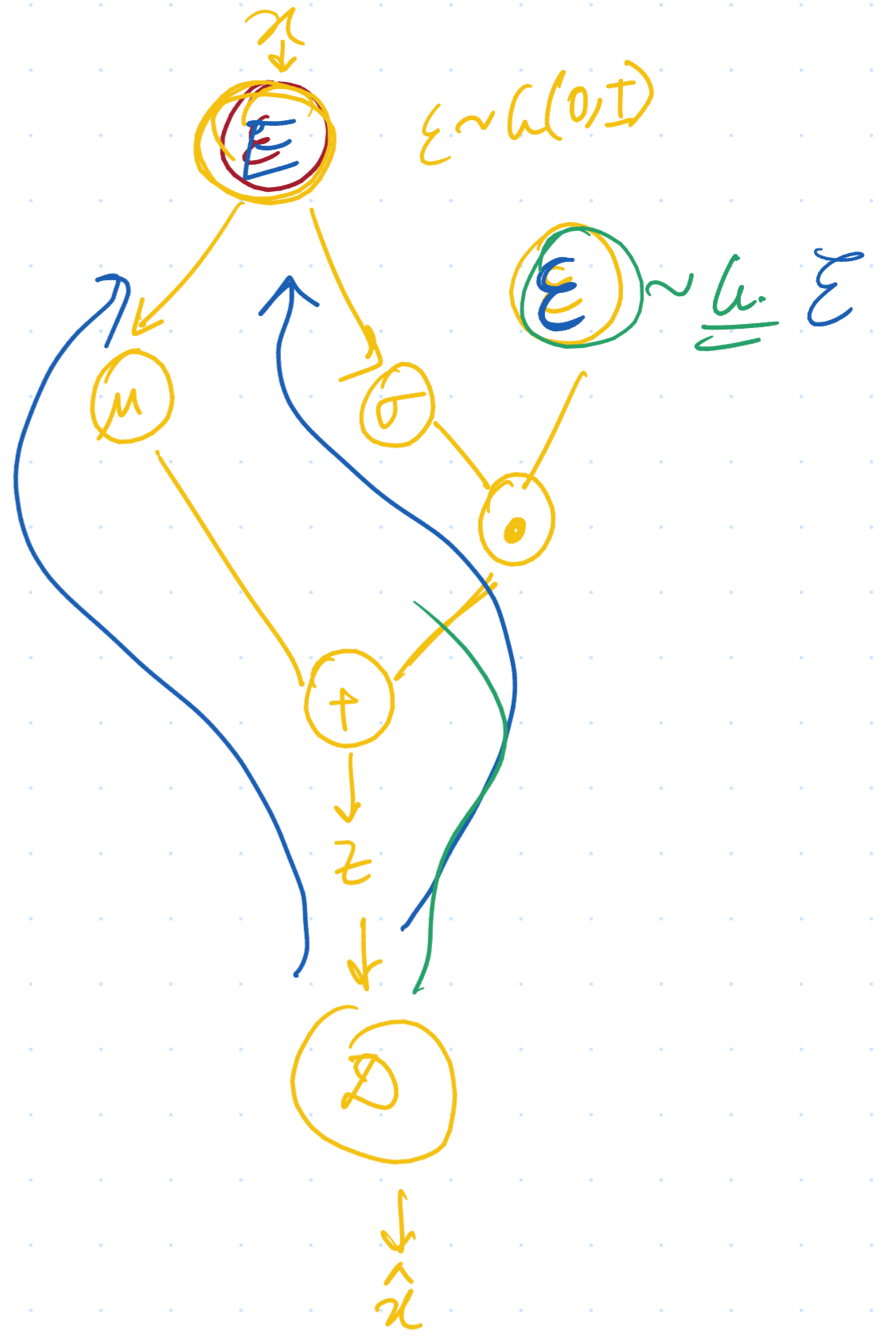
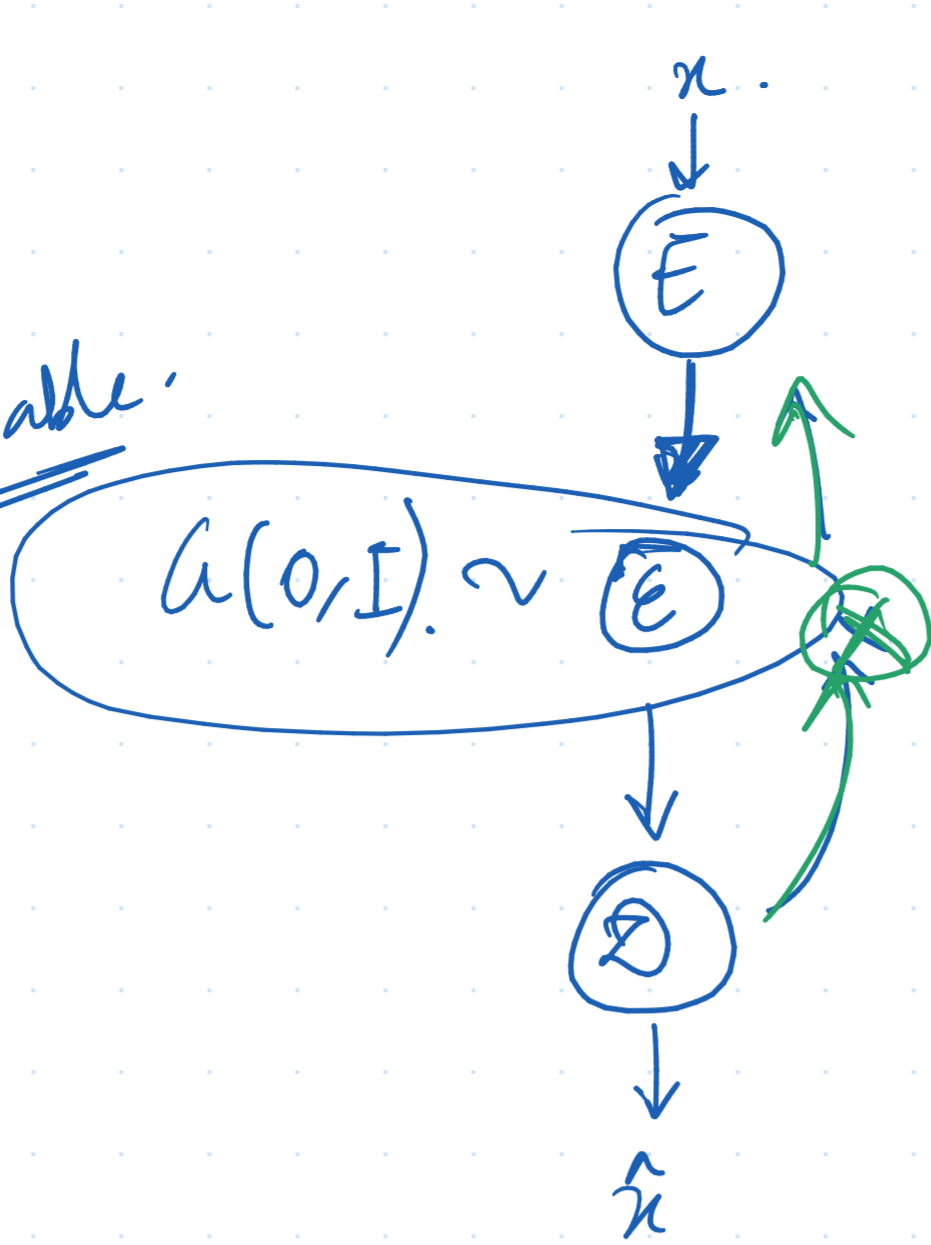
Vectors  $(\mu, \sigma)$  are learned by

(Idea of reparameterization)

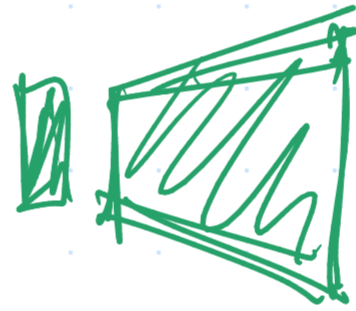
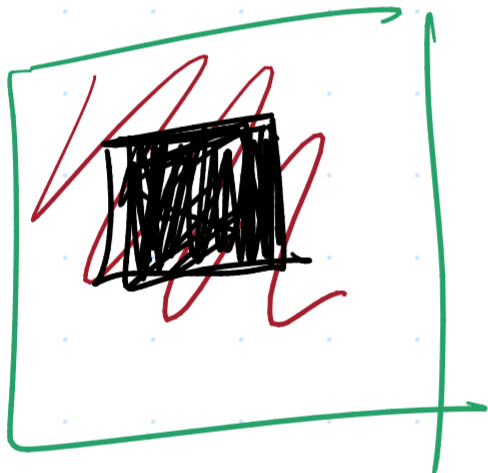
backpropagation.



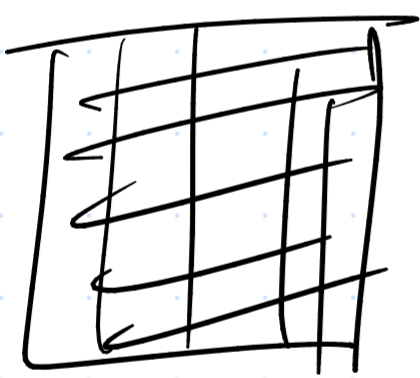
Non-differentiable.



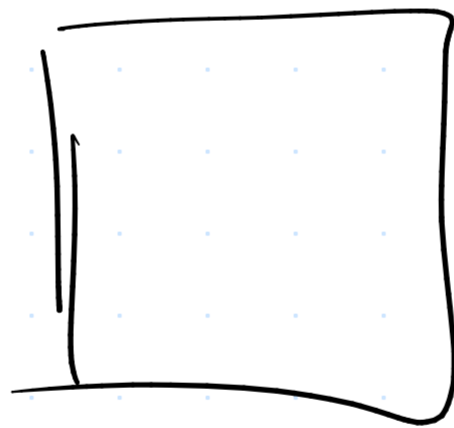
Inpainting



Colorization

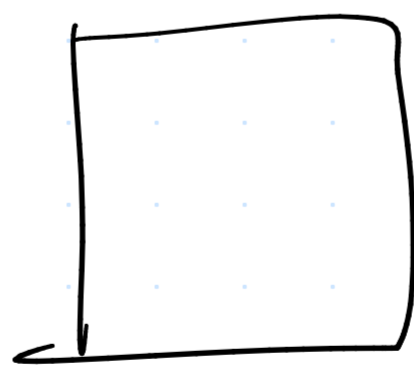


Grayscale  $z \sim$



Colored image.

$z \sim$

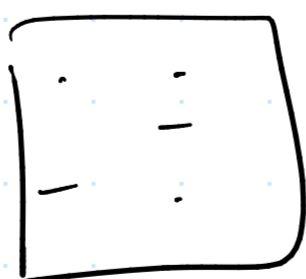


colored.

Old photos  $\rightarrow$



$\rightarrow$



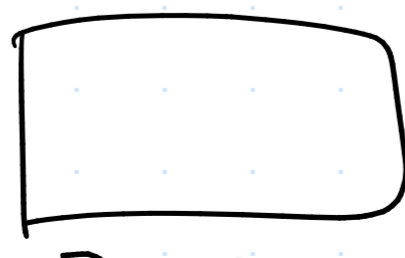
New photos

Super-resolution



144px

$\rightarrow$



720px

High resolution

Depth map. → realistic images.

De-blurring images.

unblurred images

