

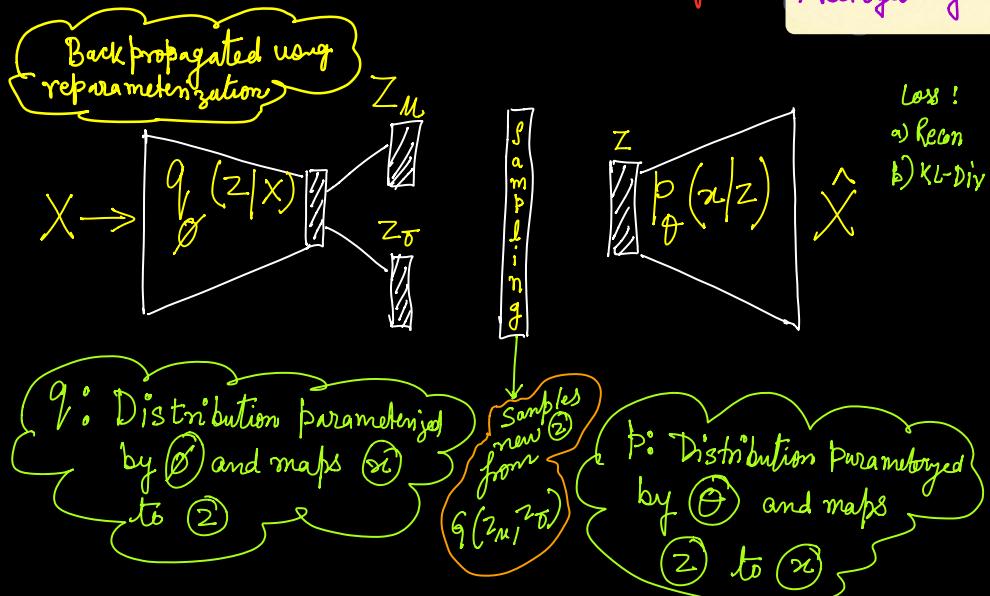
# GAN's

## Variational Auto Encoder (Summary)

Our YouTube Channel →

Deep learning for all at Manas  
lab at IIT Mandi.

Aditya Nigam - I.I.T Mandi.



**ELBO**

$$L(\theta, \phi; x, z) = E_{q_\phi(z|x)} \log p_\theta(x|z) - D_{KL}(q_\phi(z) \| p(z))$$

**Reconstruction**

This need to be minimized the ELBO

maximized

minimized

Known tractable distri bution

$N(0, 1)$

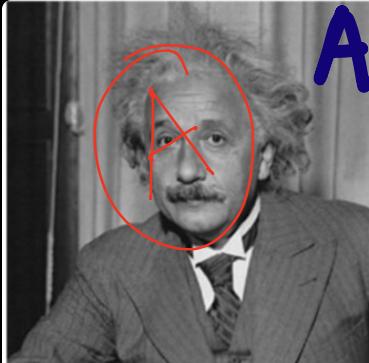
Minimizing this KL-Div makes  $q_\phi(z)$  as close to  $p(z) = N(0, 1)$

But these  $z$ 's are coming from Encoder hence averaged over  $q_\phi(z|x)$

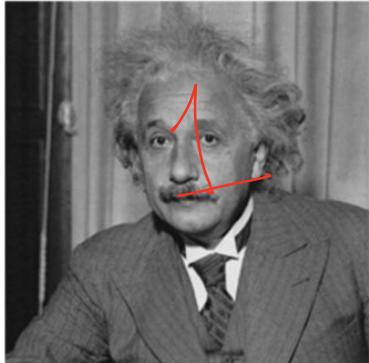
This is the log likelihood of  $x$  given  $z$ . Maximizing it ensures  $z$  should have to generate  $x$ .

(A)

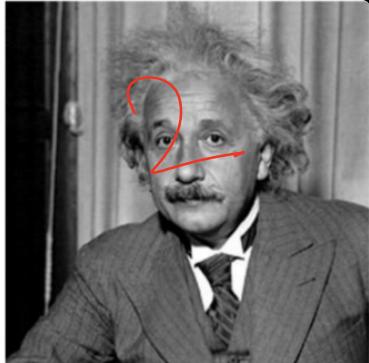
## Problems with $L_2$ / MSE loss



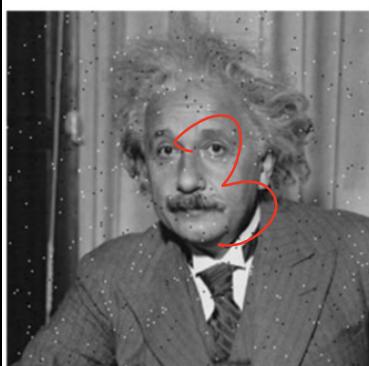
(a)



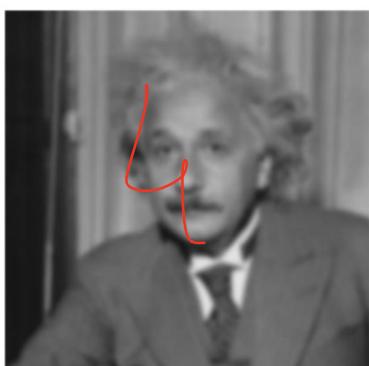
(b)



(c)



(d)



(e)



(f)

Image (A) is the original image. All the remaining images are transformed images in such a way that the MSE remains to be similar / same.

$$\mathcal{L} = \frac{1}{2 \cdot M \cdot N} \sum_i^M \sum_j^N \left( \hat{x}_{ij} - x_{ij} \right)^2$$

(Basic MSE/L2 loss.)

$$\mathcal{L} \propto \frac{1}{2} \| \hat{x} - x \|$$

Let us look at Gaussian distribution

$$p(x | \underline{\mu}, \underline{\sigma^2}) = \frac{1}{Z} e^{-\frac{(\underline{\mu} - x)^2}{2 \sigma^2}}$$

Normalization Constant      Mean ( $\sigma^2$ ) Variance

Let us assume  $\underline{\underline{M = \hat{x}}}$  and  $\underline{\underline{\sigma^2 = 1}}$

$$\therefore p(x|\hat{x}) \propto e^{-\frac{1}{2} \|\hat{x}-x\|_2^2}$$

Taking log both side

$$\log p(x|\hat{x}) \propto -\frac{1}{2} \|\hat{x}-x\|_2^2$$

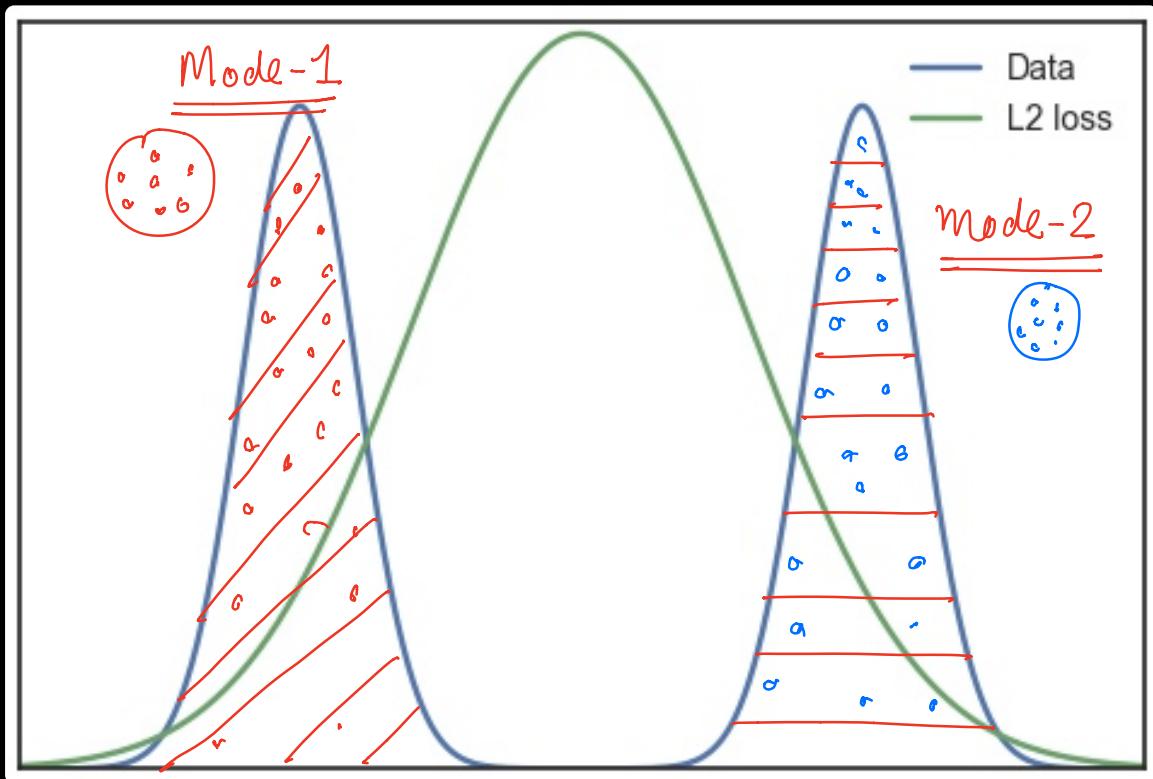
$\Rightarrow$  Minimizing  $L_2/\text{MSE}$  loss is equivalent to maximizing the log-likelihood of Gaussian.

This assumes that our real data is coming from Gaussian distribution.

But Gaussian is simple and a single peak distribution may not be able to represent complex image distribution.

In other words : What happen that our data is not following unimodal distribution but a Bi-modal distribution and we try to minimize  $(L_2)$  loss. This uses Uni-modal gaussian to fit a Bi-modal gaussian distribution.

Basically data is coming from 2 modes (Red) and (Blue)  $\Rightarrow$  2 types of data.



$L_2$  have to minimize the reconstruction loss for both types of samples during the (optimization) training.

→ Gaussian got centered right in the middle

∴ If we sample from Gaussian, sample will come from the middle/mean of both  
 $\Rightarrow$  Average of samples of 2 modes  
 Blurry image.

B) KL-Div : This is a divergence and not symmetric  
 (Issues) [Target] [Source]  $P$ : Unknown distribution (we need to estimate)  
 $Q$ : We try to estimate  $P$  using  $Q$

KL-Div → KL-Div (Reverse) :  $KL(Q \parallel P) = \sum Q(x) \log \frac{Q(x)}{P(x)}$   
 This was also utilized as a loss function for VAE.  
 $\min(Q(z) \parallel P(z|x))$   
 Parameterized by  $\phi$  [Source]  $\hookrightarrow$  Unknown [Target]

$$KL(P \parallel Q) = \sum P(x) \log \frac{P(x)}{Q(x)}$$

Question: Why do VAE uses Reverse-KL Div not Forward KL-Div.

Forward KL-Div:

$$\min_{\phi} \left\{ KL(P \parallel Q_{\phi}) \right\}$$

(Unknown)

Estimating  $P$  via  $Q$  parameterized using  $\phi$ .  
 $Q \Rightarrow$  Tractable and Parameterized.  
 Say unimodal Gaussian.

$$\min_{\phi} \left( \sum_x P(x) \log \frac{P(x)}{Q(x)} \right)$$

*Parameterization of neural network to minimize the objective function.*

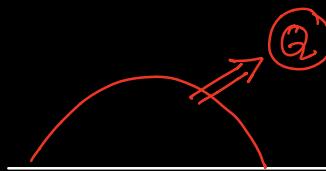
What happen when  $Q(x) = 0$  and  $P(x) \neq 0$  (some finite value)

F-KL  $\rightarrow$  tends to  $(\infty)$  not defined.

Bimodal (Unknown)



Unimodal (assumed)



$$\min_{\phi} [KL(P \parallel Q_{\phi})]$$

Tends to infinity

Here  $Q(x) = 0$

But  $P(x) \neq 0$

$F_{KL} \rightarrow \infty$   
 (Not defined)

This will never let our neural network (NN) parameterized over  $\phi$  to learn a mode of  $P$ .

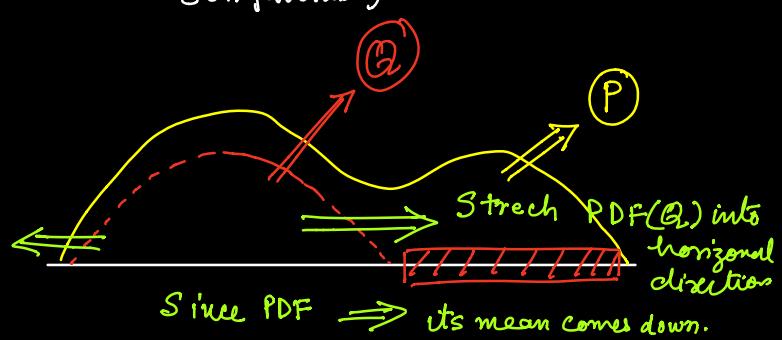
Using F-KLDIV

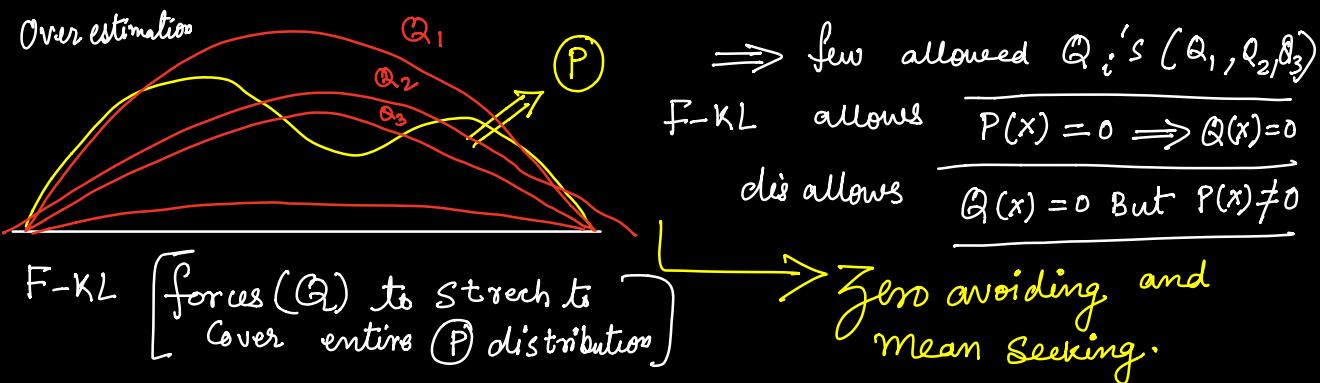
It avoids such conditions when  $Q(x) = 0$  &  $P(x) \neq 0$  therefore individual mode learning is avoided.

[Zero avoiding situation]

Very high (KL)

(unfavorable)

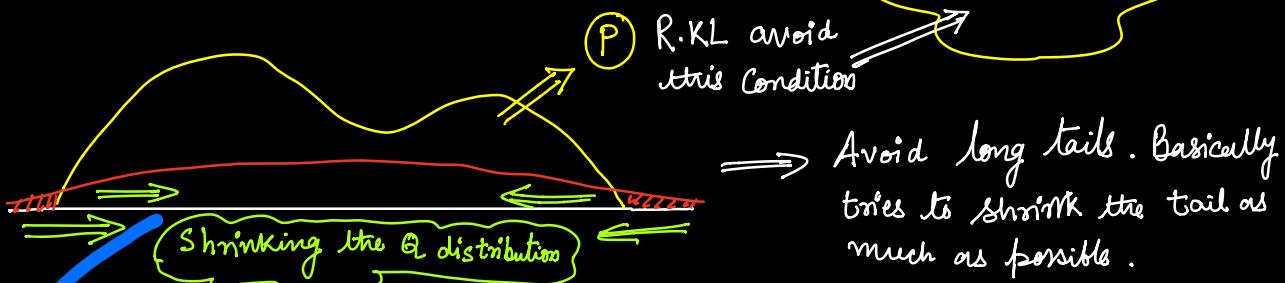




$$\text{Reverse KL-Div : } \min_{\emptyset} \left\{ \text{KL} (Q_{\emptyset} \| P) \right\}$$

$$\min_{\emptyset} \left\{ \sum_x Q_{\emptyset}(x) \log \left( \frac{Q_{\emptyset}(x)}{P(x)} \right) \right\}$$

$Q_{\emptyset}$  = non zero  
and  
 $P(x)$  = zero  
Blow it up



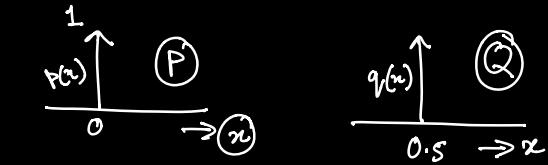
$\therefore$  F-KL  $\Rightarrow$  Stretch  $Q$  and support long tail (Mean seeking) [Over-estimating]  
 (General average behavior)

R-KL  $\Rightarrow$  Shrink  $Q$  and avoid long tail (Mode seeking) [under-estimating]  
 (estimating only some mode)

$\Rightarrow$  Mode seeking

Due to this mode seeking behaviour VAE uses R-KL. (F-KL) introduces more and more blur.

But KL-Div Suffers from an inherent problem (both F/R).  
We have two discrete PDF  $\textcircled{P}$  and  $\textcircled{Q}$ .



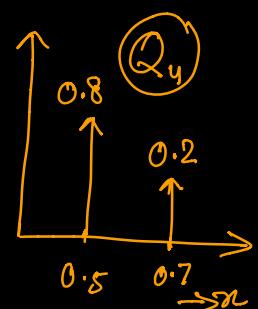
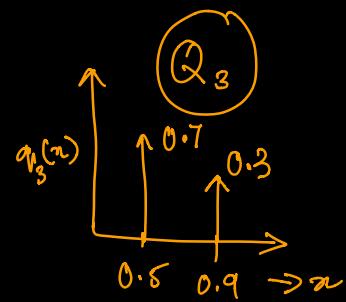
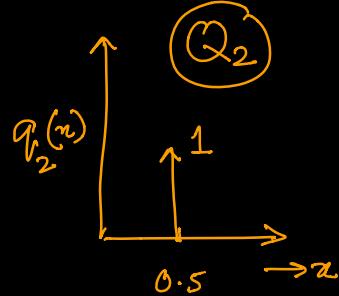
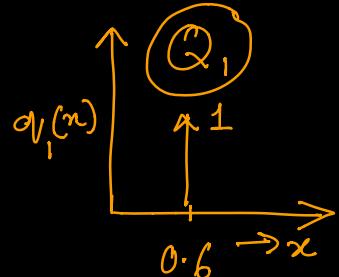
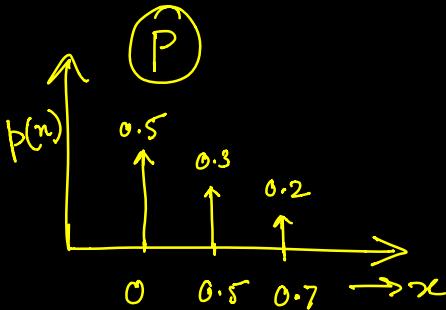
$$\begin{aligned} \text{KL}(P||Q) &= P(x=0) \log \frac{P(x=0)}{Q(x=0)} = \infty \\ \text{KL}(Q||P) &= Q(x=0.5) \log \frac{Q(x=0.5)}{P(x=0.5)} = \infty \end{aligned}$$

$\Rightarrow$  KL-Div fails to measure when there is NO-OVERLAP.

This is most of the times true especially at the start of the optimization.  $P \equiv$  unknown

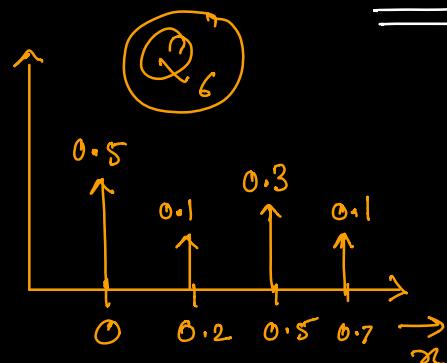
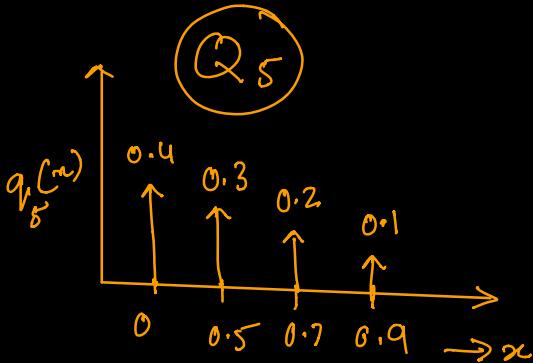
$Q \equiv$  Tractable / Parameterized

Basically  $\textcircled{Q}$  is very far from the real distribution



Compute all  
 $\text{KL}(P||Q_i)$  &  
 $\text{KL}(Q_i||P)$

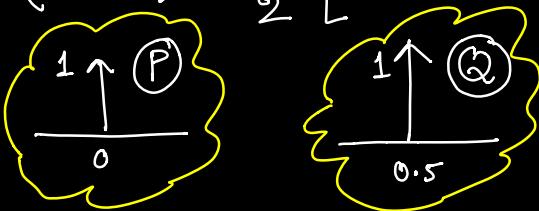
12 Values



Inference about Containment.

## JS - divergence :

$$JSD(P \parallel Q) = \frac{1}{2} [KL(P \parallel M) + KL(Q \parallel M)]$$

Say for 

$$\text{where } M = \left( \frac{P+Q}{2} \right)$$

$$KL(P \parallel M) = \sum P(x) \log \frac{P(x)}{M(x)} = P(x=0) \log \frac{P(x=0)}{\frac{P(x=0) + Q(x=0)}{2}} = 1 \cdot \log \frac{1}{\frac{1}{2}} = \log 2$$

$$KL(Q \parallel M) = \sum Q(x) \log \frac{Q(x)}{M(x)} = Q(x=0.5) \log \frac{Q(x=0.5)}{\frac{P(x=0.5) + Q(x=0.5)}{2}} = 1 \cdot \log \frac{1}{\frac{1}{2}} = \log 2$$

$$\therefore JSD(P \parallel Q) = \frac{1}{2} [\log 2 + \log 2] = \log 2$$

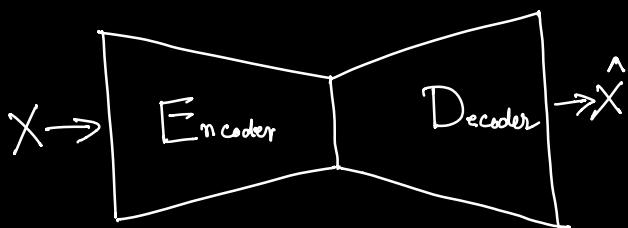
This behaviour is required when we estimate  $P$  via  $Q$  as it rarely happens that  $P$  &  $Q$  overlap right from the starting of optimization.

Moving towards generative adversarial networks (GAN).

\* Instead of MSE they utilize adversarial training.

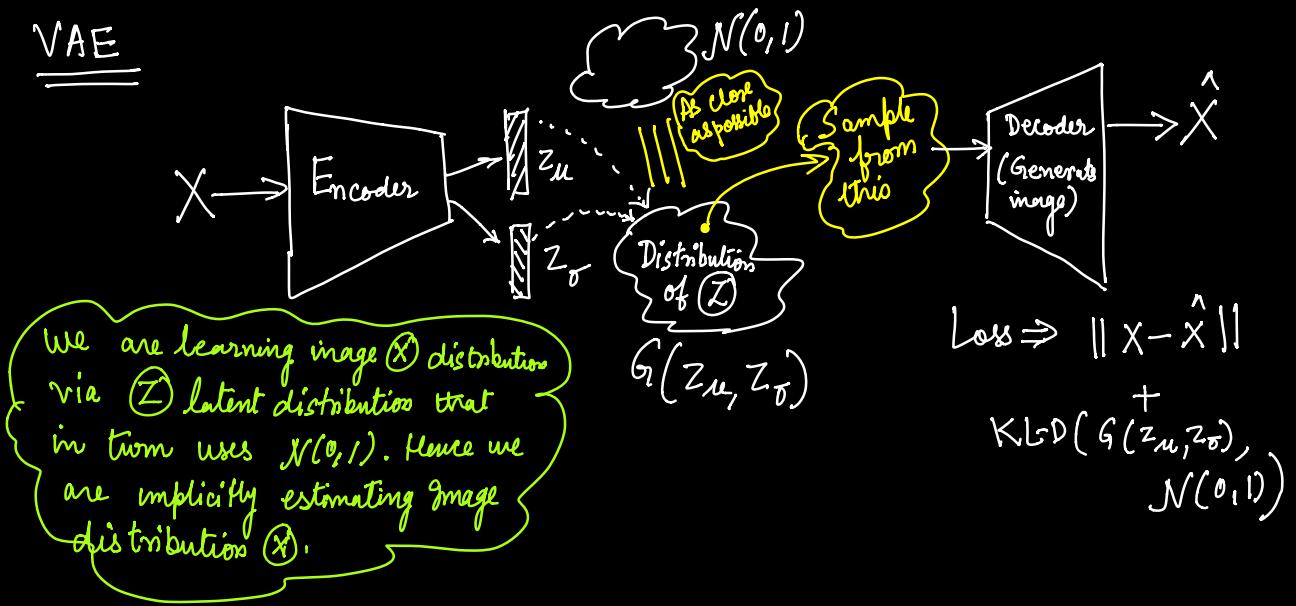
\* Instead of KL-D they utilizes  $JSD$ .

Let us redraw Basic A.E



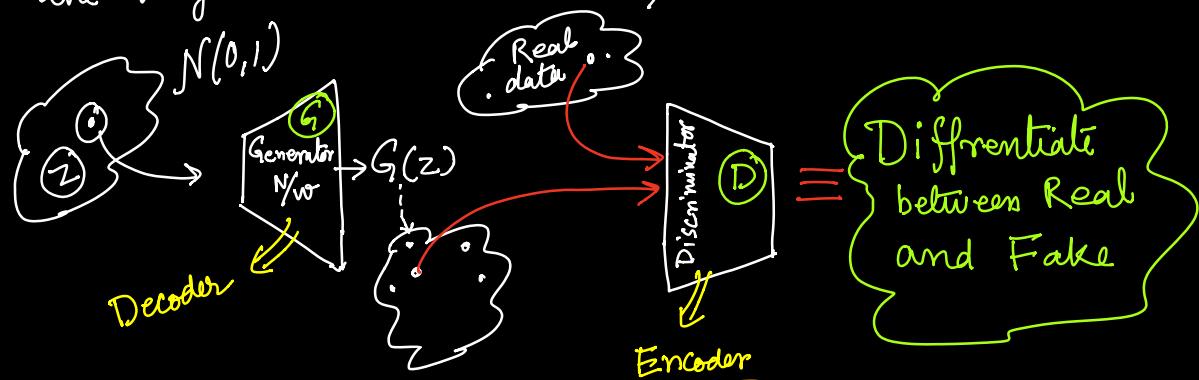
$$\text{Loss : } \| X - \hat{X} \|_2$$

VAE

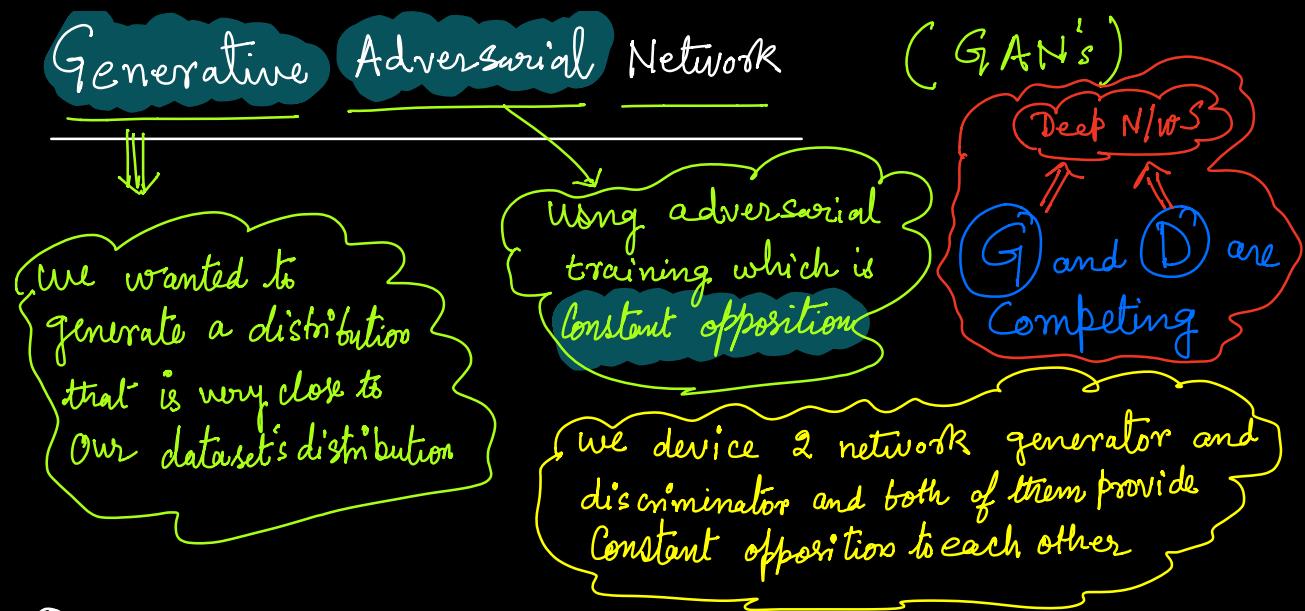


GAN [Generative Adversarial Network]

The sequence of Decoder (Generator) and Encoder (Discriminator) get changed. Their objective are also tweaked keeping the overall objective of image reconstruction being fixed. Instead of implicitly learning the image distributions (as in VAE), GAN's explicitly learns it.



Explicitly learning the image distribution  
Just trained as a classifier N/w

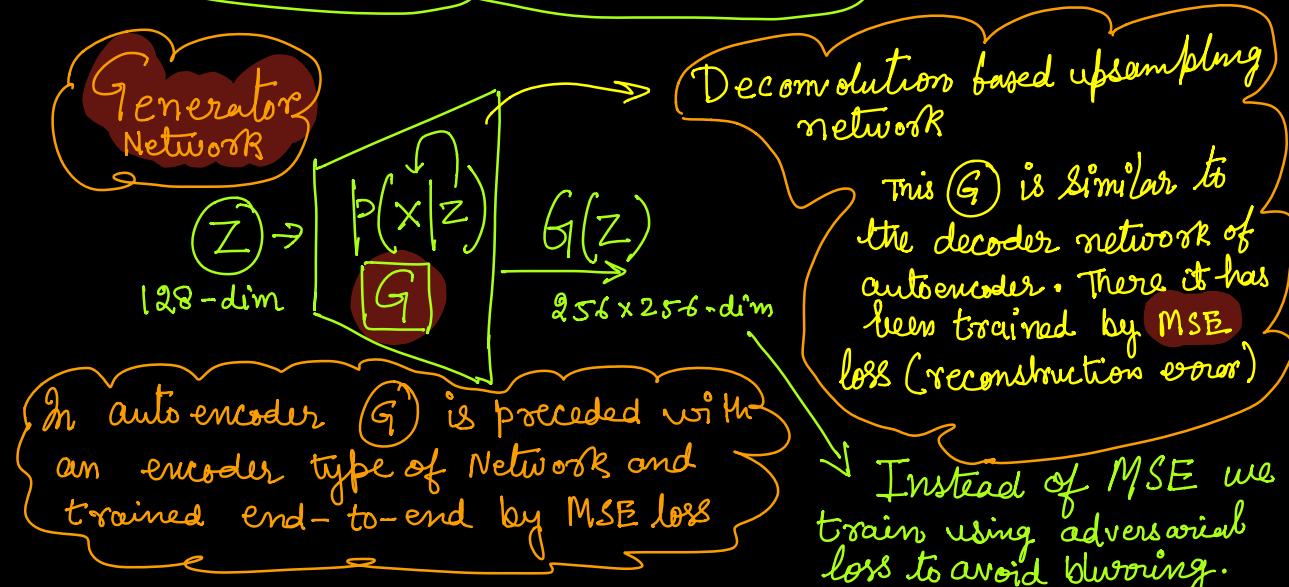


\* Adversarial training is essentially required as MSE introduces blurring effect.

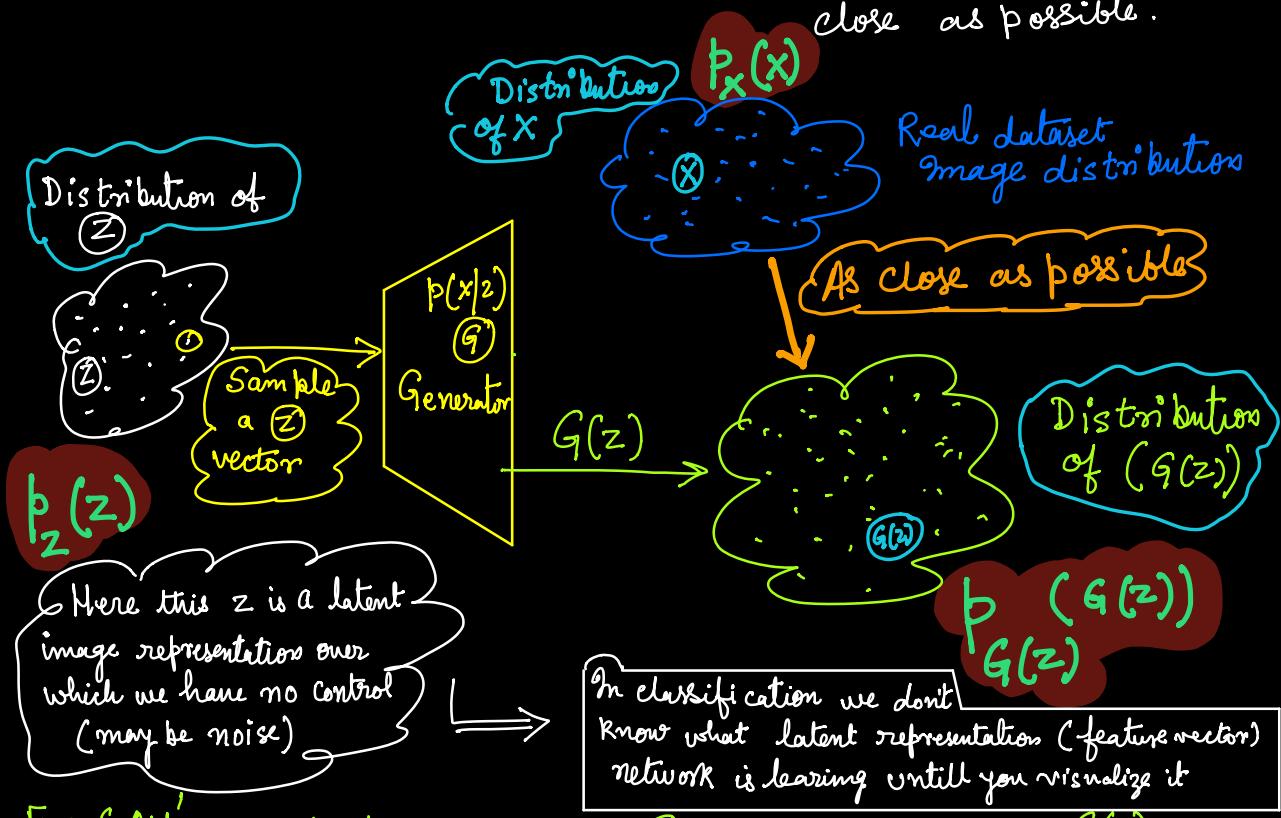
(training with an adversarial loss)

Generator network can be seen as a function that takes a random vector  $(z)$  as an input and generates  $G(z)$  an image. (say  $256 \times 256$ )

(say 128-dim) [Essentially a deconv-upsampling based NW]

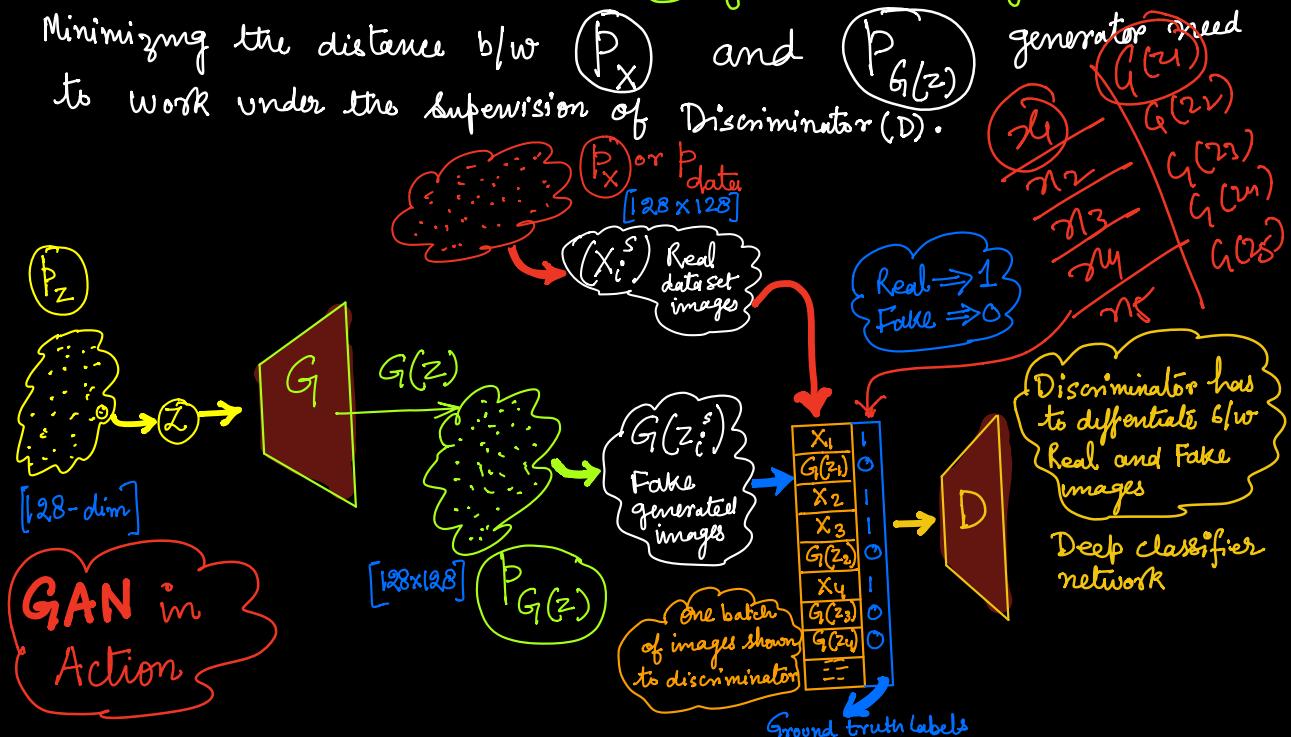


Objective of Generator network : Given  $\mathbb{Z}$  random vector generate a real image as close as possible.



For GANs, in order to know about  $\mathbb{Z}$  generate some image  $G(z)$ .

Minimizing the distance b/w  $p_x$  and  $p_{G(z)}$  to work under the supervision of Discriminator (D).



Discriminator network needs to figure out features that contributes to natural images and can help to discriminate real/fake images

While training, the feedback of discriminator network need to be backpropagated to generator

This will help generator( $G$ ) to figure out the missing features. Basically ( $G$ ) is trying to understand that

What features current generator is interpreting as real.

Success of Discriminator  $\textcircled{D}$

realized as

loss to the generator  $\textcircled{G}$  -  
(needs to be minimized)

Adversarial game :

Job of generator( $G$ ) is to fool discriminator.

Job of discriminator( $D$ ) is to Judge the generator.

The discriminator  $\textcircled{D}$  and generator  $\textcircled{G}$  networks are trained alternatively (one after other), while keep the other one fixed.

Eventually generator( $G$ ) started to make image real wrt  $\textcircled{D}$  trained over the dataset.

loss function : Binary Cross entropy for classification

$$L(\hat{y}, y) = y \log \hat{y} + (1-y) \log (1-\hat{y})$$

Predicted label  
Ground truth label

P	T.P	FN
N	FP	TN

Discriminator loss function :  $I/P \Rightarrow$  Either  $(X)$  OR  $G(z)$

[Classifier N/w]

$$\therefore L(D(x), 1) = \log [D(x)]$$

$D(x)$  minimization

$\log [D(x)]$  minimization

$O/P \Rightarrow D(x) \text{ OR } D(G(z))$

$y$  Ground Truth (label)  $\Rightarrow \begin{cases} 1 & (\text{Real}) \\ 0 & (\text{Fake}) \end{cases}$

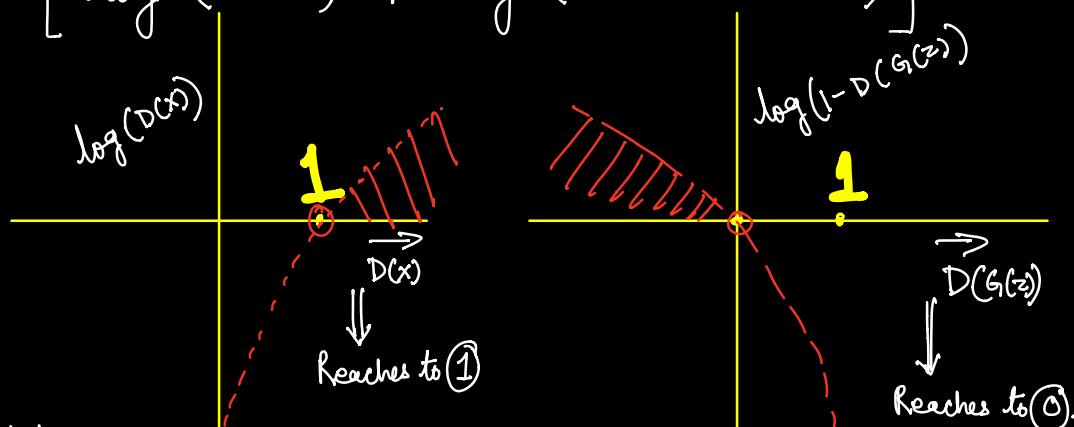
$$L(D(G(z)), 0) = \log [1 - D(G(z))]$$

$D$  need to minimize  $D(G(z))$   $\equiv$  maximize  $1 - D(G(z))$

maximize  $\log (1 - D(G(z)))$

Discriminator loss fn : (Single Sample)

$$\therefore \underset{D}{\text{Max}} \left[ \log (D(x)) + \log (1 - D(G(z))) \right]$$

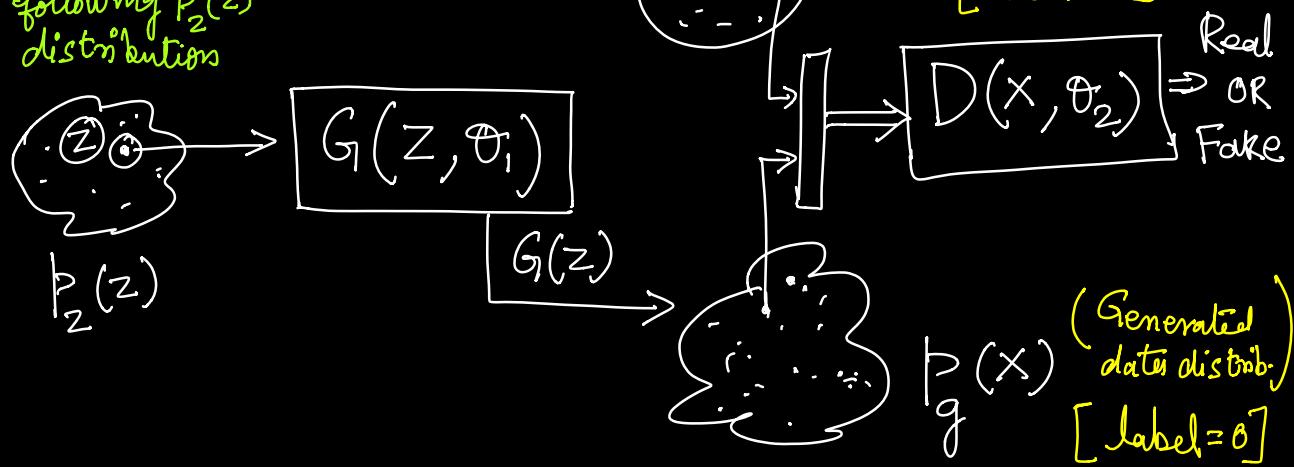


Over a batch

$$\underset{D}{\text{Max}} \left[ \underset{x \sim p(x) \text{ data}}{\mathbb{E}} (\log (D(x))) + \underset{z \sim p_z(z)}{\mathbb{E}} (\log (1 - D(G(z)))) \right]$$

Generator needs to fool Discriminator by generating image as real as possible hence:  $\underset{G}{\text{Min}} \underset{z \sim p_z(z)}{\mathbb{E}} (1 - D(G(z)))$

$\mathbb{Z}$  is a randomly distributed r.v following  $p_z(z)$  distribution



$$\begin{aligned} \text{Generator loss fn: } & I/P : G(\mathbb{Z}) \\ & (\hat{y}) O/P : D(G(\mathbb{Z})) \\ & (\hat{y}) G/T : 0 \end{aligned}$$

" From binary cross entropy (Single Sample)

$$L(D(G(\mathbb{Z})), 0) = \log(1 - D(G(\mathbb{Z})))$$

(over a batch)  $\Rightarrow \mathbb{E}_{\mathbb{Z} \sim p_z(z)} [\log(1 - D(G(\mathbb{Z})))]$ .

Since  $G$  is never going to see any real data, just for completeness one can write Generators complete loss fn as.

$$\underset{G}{\text{Min}} \left[ \mathbb{E}_{x \sim p_{data}(x)} [\log(D(x))] + \mathbb{E}_{\mathbb{Z} \sim p_z(z)} [\log(1 - D(G(\mathbb{Z})))] \right]$$

# Total Combined Loss Fn :

$$\min_G \max_D V(D, G) =$$

$$\min_G \max_D \left( \mathbb{E}_{\substack{x \sim p(x) \\ \text{Real data}}} [\log(D(x))] + \mathbb{E}_{\substack{z \sim p_z \\ \text{Random vector}}} [\log(1 - D(G(z)))] \right)$$

↓  
 Image generated by random vector

## Generative Adversarial Nets

10-June  
2014

Ian J. Goodfellow, Jean Pouget-Abadie\*, Mehdi Mirza, Bing Xu, David Warde-Farley,  
Sherjil Ozair† Aaron Courville, Yoshua Bengio†  
Département d'informatique et de recherche opérationnelle  
Université de Montréal  
Montréal, QC H3C 3J7

**Algorithm 1** Minibatch stochastic gradient descent training of generative adversarial nets. The number of steps to apply to the discriminator,  $k$ , is a hyperparameter. We used  $k = 1$ , the least expensive option, in our experiments.

```

for number of training iterations do
  for  $k$  steps do
    • Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
    • Sample minibatch of  $m$  examples  $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$  from data generating distribution  $p_{\text{data}}(\mathbf{x})$ .
    • Update the discriminator by ascending its stochastic gradient:
```

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[ \log D(\mathbf{x}^{(i)}) + \log (1 - D(G(\mathbf{z}^{(i)}))) \right].$$

**end for**

- Sample minibatch of  $m$  noise samples  $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\}$  from noise prior  $p_g(\mathbf{z})$ .
- Update the generator by descending its stochastic gradient:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log (1 - D(G(\mathbf{z}^{(i)}))).$$

**end for**

The gradient-based updates can use any standard gradient-based learning rule. We used momentum in our experiments.

Point-01 : First we will show that given any fixed  $G$ , the optimally achieved  $D_G^*$  can be defined as :

$$D_G^*(x) = \arg \max_{\mathcal{D}} V(D, G)$$

Best possible  $\mathcal{D}$

$$= \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

Point-02 : Secondly for any given (say best The optimal  $G^*$  can be defined as  $(D_G^*)$ )

$$G^* = \arg \min_G V(D_G^*, G)$$

fixed  $\mathcal{D}_G^*$  to get best generator.

It is been shown that this optimization have a unique solution achieved when

$$[P_g = P_{\text{data}}]$$

desirable form the problem itself

Also it has been shown that it is equivalent to minimize  $JSD(P_{\text{data}} || P_g)$ .

Proving point - 01 : For a given  $G$

Our objective is to Get the optimal  $D$ .

$$\therefore D_G^* = \arg \max_D V(D, G)$$

Best possible  $D$  for a given  $G$

$$\therefore D_G^* = \arg \max_D \left\{ \mathbb{E}_{x \sim p(x)_{\text{data}}} [\log(D(x))] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \right\}$$

$$\mathbb{E}_{p(x)} [x] = \int_X x p(x) dx$$

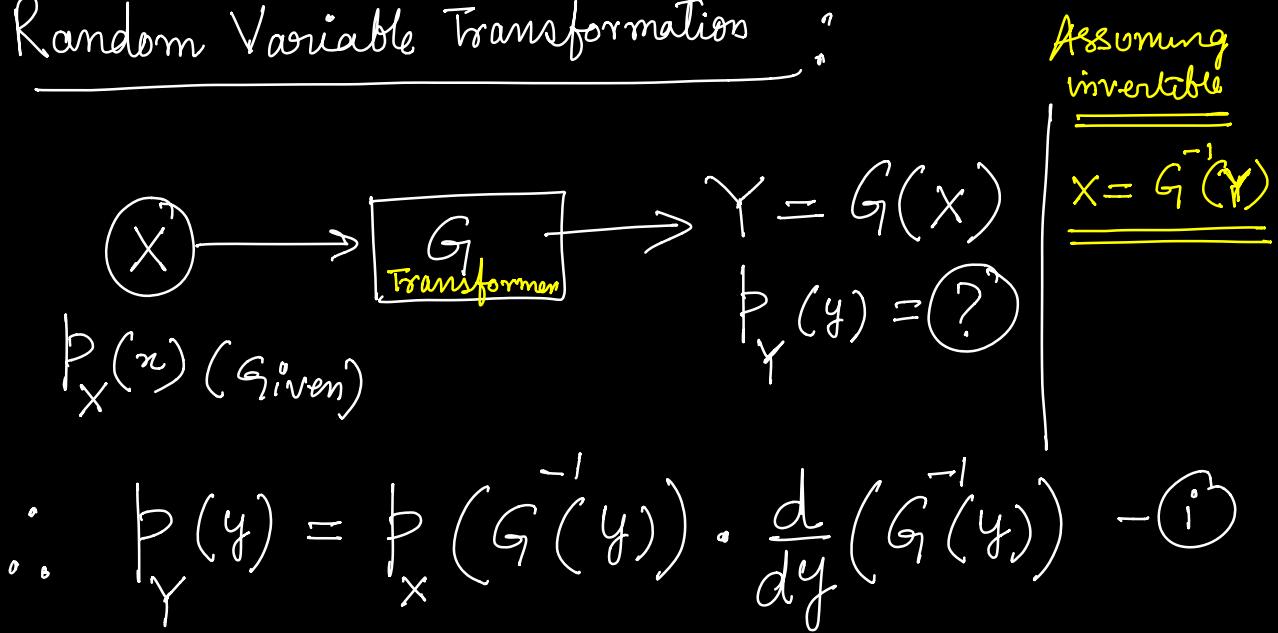
Expected value  
of  $RoV(x)$ ,  
following PDF  
 $\Rightarrow p(x)$

Hence  $V(G, D)$ :

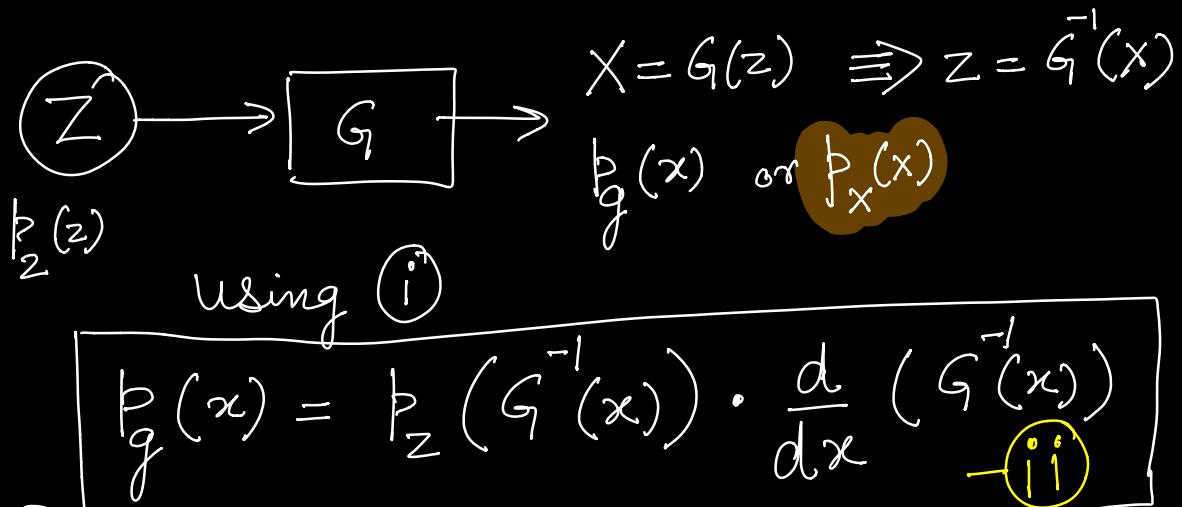
$$V(G, D) = \overbrace{\int_X p(x)_{\text{data}} \cdot \log(D(x)) dx}^A + \overbrace{\int_Z p_z(z) \cdot \log(1 - D(G(z))) dz}^B$$

Eq - 01

## Random Variable Transformation



Now Consider the case of GAN's.



(Now let us see what total loss again)

$V(G, D) :$

$$V(G, D) = \int_{\text{data}} P_G(x) \cdot \log(D(x)) dx \quad \boxed{\text{Eq-01}}$$

$$+ \int_z P_Z(z) \cdot \log(1 - D(G(z))) dz \quad \boxed{\text{Eq-02}}$$

Term (B):

$$\int_z p_z(z) \log(1 - D(G(z))) dz$$

*Changing  
of variable*

$$\text{Let } z = G^{-1}(x) \quad \boxed{x = G(z)}$$

$$\therefore dz = \frac{d(G^{-1}(x))}{dx} dx$$

$$\therefore \int_x p_z(G^{-1}(x)) \cdot \log(1 - D(x)) \cdot \frac{d(G^{-1}(x))}{dx} dx$$

$\Downarrow$   $p_g(x)$  (using (ii))

$$\boxed{\int_x p_g(x) \log(1 - D(x)) dx}$$

Hence finally  $V(G, D)$ :

$$\int_x p_{\text{data}}(x) \log(D(x)) dx + \int_x p_g(x) \log(1 - D(x)) dx$$

$$D_G^*(x) = \arg \max_D [V(D, G)]$$

optimal      fixed

*(only in  $x$ )*

Hence in order to minimize  $V(D, G)$   
we need to differentiate it wst  $[D(x)]$ .

$$\therefore \frac{d(V(D, G))}{d[D(x)]} = \frac{p_{\text{data}}(x)}{D(x)} - \frac{p_g(x)}{1 - D(x)} = 0$$

$$\therefore \boxed{D_G^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}} \quad \begin{array}{l} \text{Practically} \\ \text{Cannot} \\ \text{be Computed} \\ \text{but analytically,} \\ \text{proved.} \end{array}$$

Now the role of Generator is reverse of  $(D)$ .

Hence for this best an optimally chosen  $D_G^*$ , we need to obtain the optimal  $G^*$

by minimizing the loss fm  $V(D, G)$   
over all possible  $G$ 's for fixed  $D_G^*$

$$\therefore G^* = \underset{G}{\operatorname{argmin}} V(D_G^*, G)$$

$$G^* = \underset{G}{\arg \min} \left\{ \int_X p_{\text{data}}(x) \log(p_G^*(x)) dx + \int_X p_g(x) \log(1 - p_G^*(x)) dx \right\}$$

$$= \underset{G}{\arg \min} \left\{ \int_X \left( p_{\text{data}}(x) \cdot \log \left( \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right) + p_g(x) \cdot \log \left( \frac{p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right) \right) dx \right\}$$

$$+ \log(2) \cdot p_{\text{data}}(x) - \log(2) \cdot p_{\text{data}}(x)$$

$$+ \log(2) \cdot p_g(x) - \log(2) \cdot p_g(x)$$

they got combined

$$= \underset{G}{\arg \min} \left\{ \int_x \left( -\log_2 \left( p_{\text{data}}(x) + p_g(x) \right) \right) \right.$$

$$+ p_{\text{data}}(x) \cdot \log \frac{2 \cdot p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)}$$

$$p_g(x) \cdot \left. \log \frac{2 \cdot p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right) dx$$

$$\Rightarrow \int_x p_g(x) = \int_x p_{\text{data}}(x) = 1$$

$$\therefore G^* = \underset{G}{\arg \min} \left[ \int_x \left( -\log_2(2) \cdot (2) + p_{\text{data}}(x) \log \left( \frac{2 \cdot p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right) \right. \right. \\ \left. \left. + p_g(x) \log \left( \frac{2 \cdot p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right) \right) dx \right]$$

Hence

$$\begin{aligned} G^* &= \underset{G}{\operatorname{argmin}} \left[ -\log(4) + \int_x p_{\text{data}}(x) \cdot \log \left( \frac{2 \cdot p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \right) dx \right. \\ &\quad \left. + \int_x p_g(x) \cdot \log \left( \frac{2 \cdot p_g(x)}{p_{\text{data}}(x) + p_g(x)} \right) dx \right] \\ &= \underset{G}{\operatorname{argmin}} \left[ -\log(4) + 2 \cdot \text{JSD}\left(p_{\text{data}}(x), p_g(x)\right) \right] \end{aligned}$$

This shows that  $G^*$  optimization is equivalent to JSD minimization between  $p_{\text{data}}(x)$ ,  $p_g(x)$

→ This implies that minimum JSD can be equal to (zero) that too when  $p_g(x) = p_{\text{data}}(x)$

• After the above minimization Convergence,

$$\boxed{\text{JSD} \rightarrow 0} \quad \boxed{p_g(x) \rightarrow p_{\text{data}}(x)}$$

and  $\boxed{V(D, G) \rightarrow -\log(4).}$

This shows that :

The global minima of the Criterion for  $V(D_g^*, G)$  is achieved iff  $(P_g = P_{\text{data}})$  and the loss will be equal to  $[-\log(4)]$ .

Hence the goal of GAN's is to generate the optimal  $(G^*)$  which is effectively done by minimizing the  $JSD(P_{\text{data}}, P_g)$  and converges to  $\boxed{P_g = P_{\text{data}}}$ .

$\therefore$  The corresponding  $D_g^* = \frac{P_{\text{data}}}{P_{\text{data}} + P_g} = \frac{1}{2}$   
(loss) can be

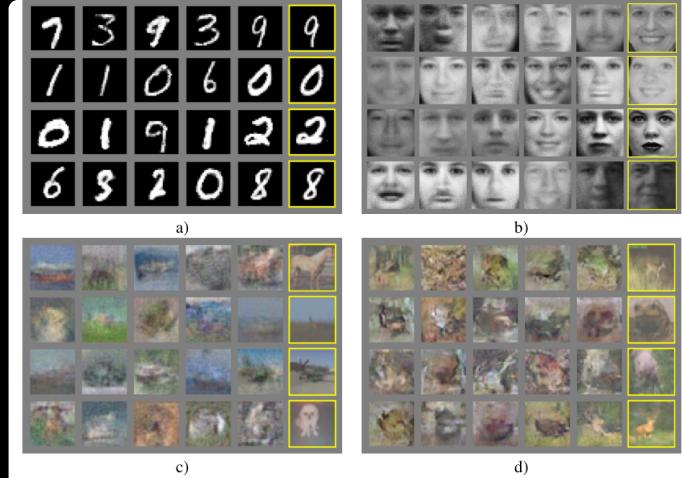


Figure 2: Visualization of samples from the model. Rightmost column shows the nearest training example of the neighboring sample, in order to demonstrate that the model has not memorized the training set. Samples are fair random draws, not cherry-picked. Unlike most other visualizations of deep generative models, these images show actual samples from the model distributions, not conditional means given samples of hidden units. Moreover, these samples are uncorrelated because the sampling process does not depend on Markov chain mixing. a) MNIST b) TFD c) CIFAR-10 (fully connected model) d) CIFAR-10 (convolutional discriminator and "deconvolutional" generator)

## Problems/Challenges/Limitations of the Current GAN models [AK/JH]

I

Diminishing Gradients

- \* Initially  $G$  is generating very poor samples and  $D$  will actually be trained fast. Training of  $G$  is also very slow.

$x$  &  $G(z)$  are very different and  $D(x) \approx D(G(z))$



$G$  is learning from  $D$ 's mistakes that getting rarer

Loss for D

$$\max_D \left\{ \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z)))] \right\}$$

Real  
Fake generated image  
Random vector

$D$  tries to make it close to 1

Loss for G

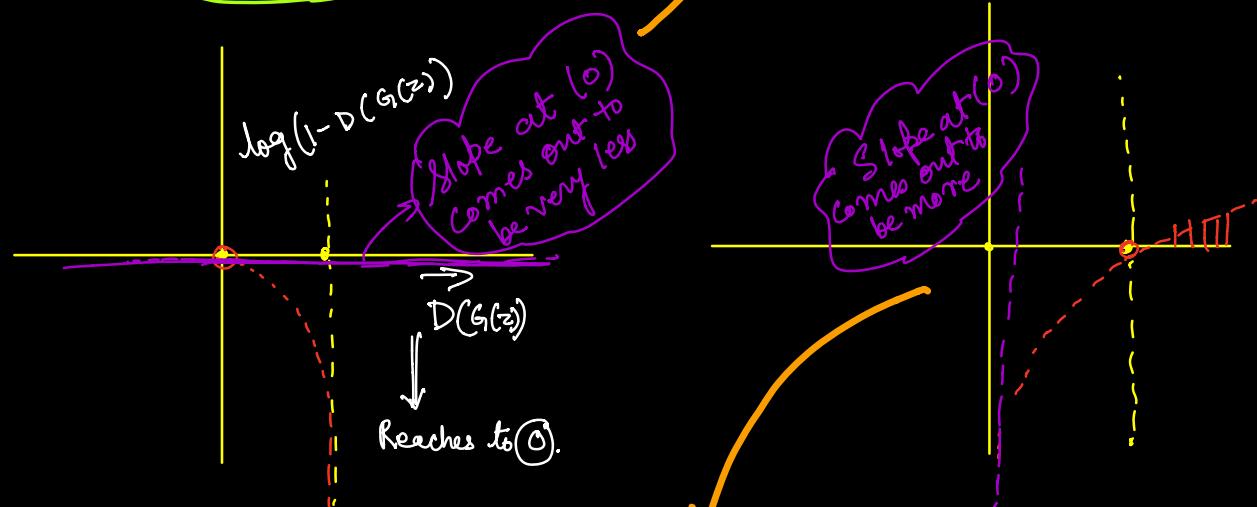
$$\min_G \left\{ \mathbb{E}_{z \sim p_z(z)} [\log (1 - D(G(z)))] \right\}$$

$G$  makes 0  
 $D$  makes 1

(A)

Early on the generator network training is very slow

Hence we need to tweak the loss so as to manage the derivative early on.



Equivalent to

$$\max_G \left\{ \mathbb{E}_{z \sim p_z(z)} [\log [D(G(z))]] \right\}$$

