

# Deep Unsupervised Learning

## L2 - Auto regressive models

Problems to solve: with likelihood based models.

- generating data
- compressing data
- anomaly detection
- learning a distri  $P$ , - computing  $P(x)$  for  $x$ 
  - sampling from  $P(x)$
- estimate distri of complex, high dim dat.
- computational and statistical efficiency.
  - Efficient training & model rep
  - Expressiveness and generalisation
  - Sampling quality, & speed
  - Comprehension rate & speed.

### Histograms

estimate part of underlying data distri from samples.

- histogram plots on occurrence of each date point on sampling.
- Inference  $\rightarrow$  lookup.
- Sampling  $\rightarrow$  sample from uniform distri, then move through the space from that point by probability mass that point, then pick it.
- failure in high dim.
- won't generalise for hot cases.

Scanned with CamScanner

### Parameterised Distri:

- fit a param distri over training data
- use fn approx.  
 $P_0(x) \approx P_{\text{data}}(x)$
- solve using MLE, so
- $P_0$  will be d. dim.
- probability distri over  $x$  should satisfy  $\sum_x P_0(x) = 1$   
 so we always have  $\sum_x P_0(x) = 1$ .

### Auto regressive models

Bayes net struct

chain rule

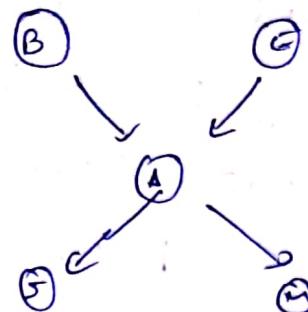
$$P(B, E, A, J, M)$$

$$= P(B) \cdot P(E|B) \cdot P(A|E, B)$$

$$P(J|A, E, B)$$

$$P(M|J, A, E, B)$$

$$p(M|A)$$



→ assumption

The sparsity from assumption leads to small tables for the conditionals. This makes it representable.

Auto regressive models parameterise their conditionals.

$$\log P_0(x) = \sum_i^d \log P_0(x_i | \text{parents}(x_i))$$

for an image, every pixel is an entry  
in  $x$ , each column channel is an entry.

$$\log P(x) \rightarrow \sum_i^d \log p(x_i | x_{1:i-1})$$

this is going to be parameterised

$(x_1, x_2) \in \mathcal{X}$

model:  $p(x_1, x_2) \rightarrow p(x_1) p(x_2 | x_1)$

$\downarrow$  histogram       $\downarrow$  parameterize using MLP,  
o/p on softmax

for high dims

- order of  $d$  params ( $O(d)$ )
- would have been  $O(\exp(d))$  in tabular

issues:

- for problems like text gen,  $d$  is very large
- No info sharing b/w conditionals, limited generalisation.

Soln:

- RNNs,
- Masking

RNN autoregressive models: (char-RNN)

$$\log p(x) = \sum_{i=1}^d \log p(x_i | x_{1:i-1})$$

↑  
character at  
i<sup>th</sup> position.

↑  
sequence of  
chars

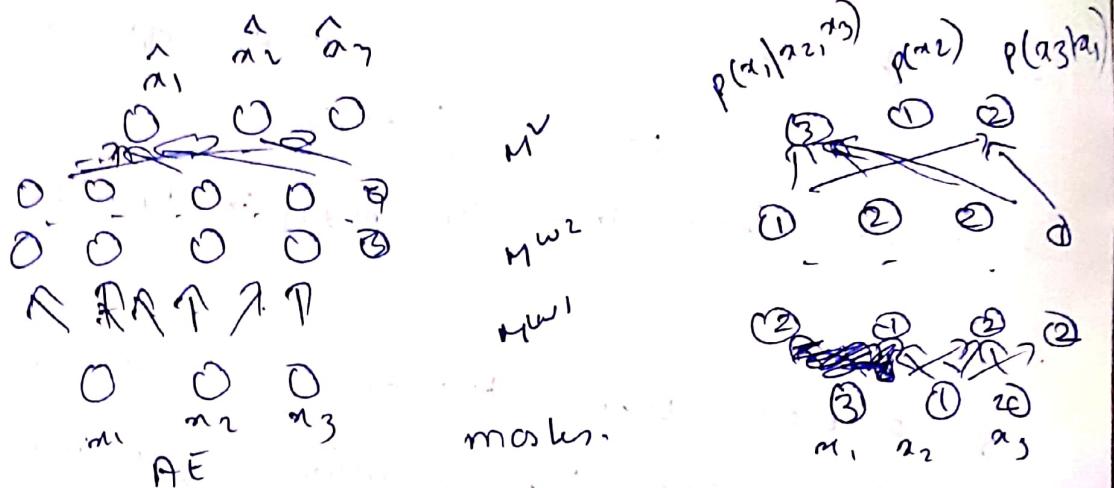
- can also generate images using RNN this way.

Masking based models

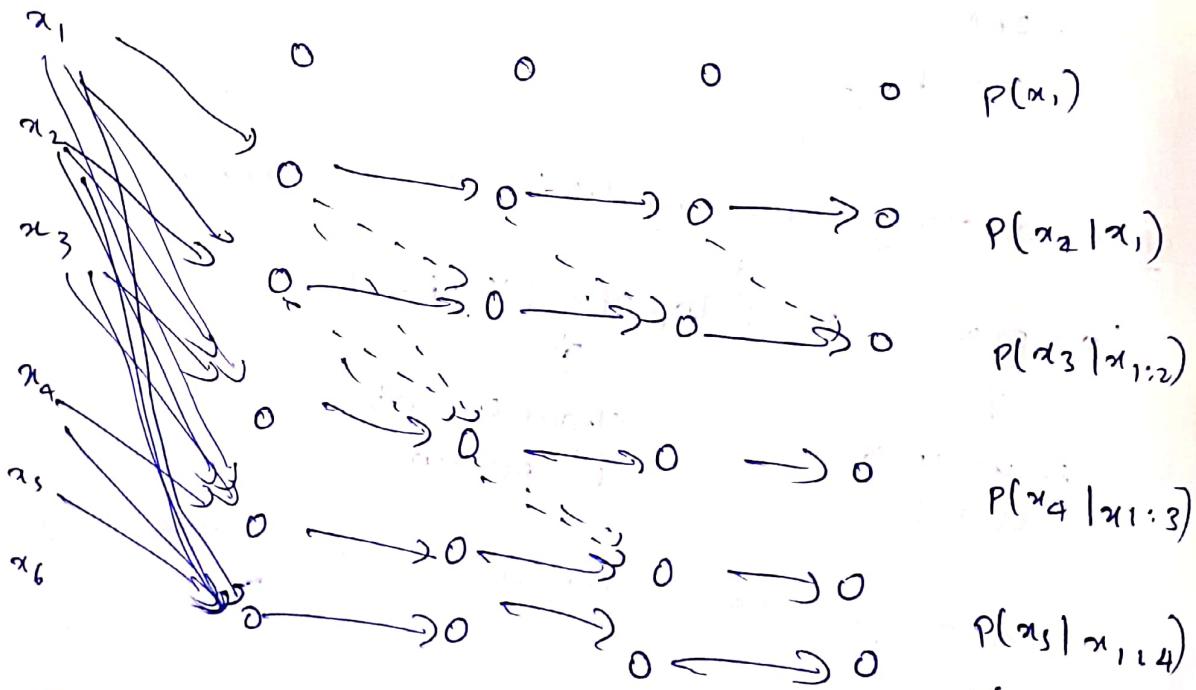
- parallelised computation of all conditionals
- Masked MLP
- Masked conv & self-attn.

Made autoencoder for distri estimation  
(MADE)

Autoencoder  $\times$  Masks  $\rightarrow$  MADE

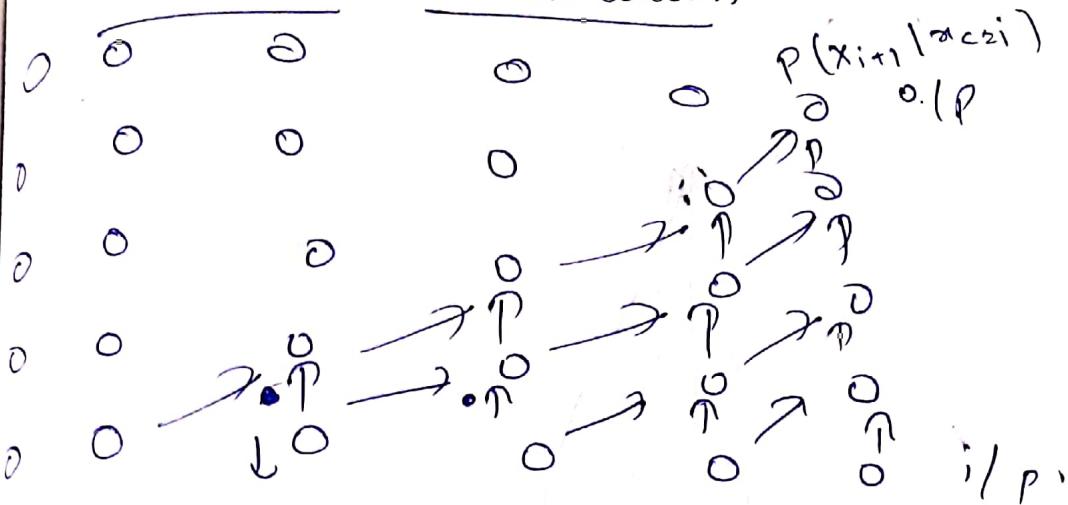


- P conditional probs can be obtained by removing a few edges.
- Done by defining which edges go where.



- all the dotted lines can be kept.
- Orderings matter.

### Masked 1D convolution



- these ops can be masked conv,  
Same operation everywhere.
- constant param count for variable length
- conv highly hardware optimised
- But limited receptive field

### WaveNet

- Improved receptive field, dilated conv, with exponential, dilation
- gated Res blocks, skip connection
- works on mnist with position encoding.

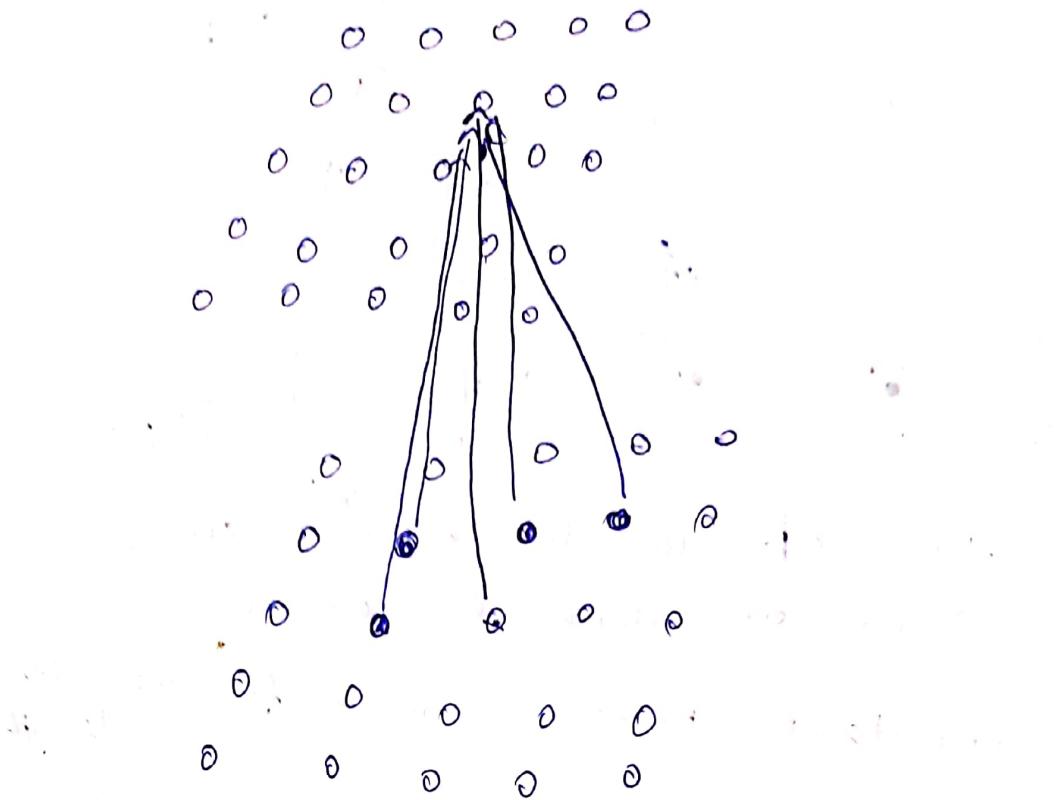
### Pixel CNN

- Masked 2D conv
  - change the filter, mask it directly.
- $x_1 \dots x_n$   $\rightarrow$  raster scan order.

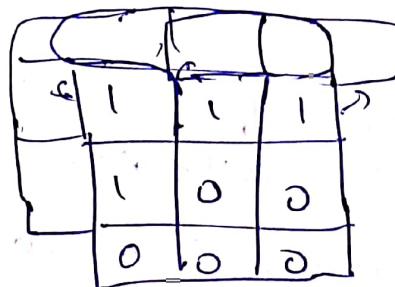
$x_i$

- mask  $\rightarrow$   $\begin{matrix} 1 & 1 & 1 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{matrix}$

- multiple layers



receptive field of 2. layer

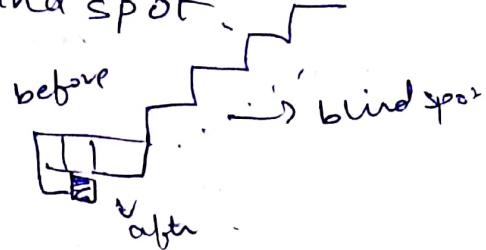


this will keep growing.

It grows such that the receptive field shifts to things before, nothing on itself on what comes after.

Keep doing pixel after pixel,

- But there is a blind spot.

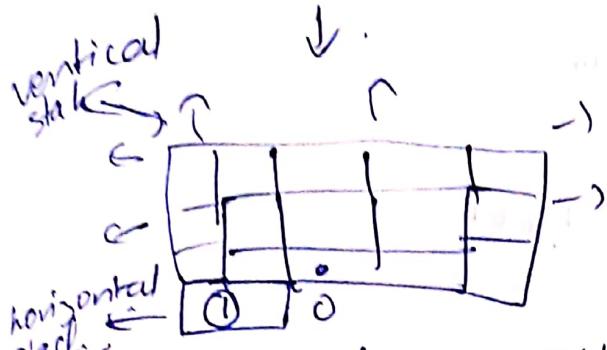


To avoid the blindspot, change kernel to



$o/p \rightarrow \text{softmax}$  over 256

receptive field



This is used on all neurons.  
But in last layer, the centred point  
is trained again to predict the next one.  
(would be cheating otherwise as it saw  
the point before)

But now there is a blindspot at ①  
a 1D filter is used for that.

gated Pixel CNN

$$y = \tanh(w_{k,f} * x) \odot \sigma(w_{k,g} * x),$$

Pixel CNN +

- nearby pixels are likely to co-occur.

$$p(x | \pi, \mu, s) = \sum_i^K \pi_i [\sigma((x + 0.5 - \mu_i) / s_i)]$$

$$\sigma((x - 0.5 - \mu_i) / s_i)]$$

$\mu \rightarrow \text{mean}$   
 $s \rightarrow \text{std}$

parametrising sigmoid  
cumulative disti

- The cliff b/w  $\sigma$  & sigmoid give diversity
- Capture long dependencies efficiently by downSampling (skip connections)

### Masked Attention

- Multi-headed self attn.
- good masking
- unlimited receptive field

### Attention

query, a hidden vector, typically generated by a previous layer; Another part of n/w will have key, value pairs. The query will do inner products with keys. Then weighted sum of values based on how good of match b/w key and query.

$$A(q, k, v) = \sum_i \frac{e^{q_k i}}{\sum_j e^{q_k j}} v_i$$

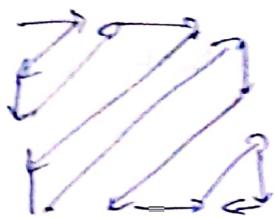


masking unnecessary connection

- = subtract a very -ve number

$$\left\{ \begin{array}{l} \underbrace{e^{q_k i} - \text{masked}(k, q)}_{\text{attn.}} \cdot 10^{10} \\ \underbrace{\sum_j e^{q_k j} - \text{masked}(k, q)}_{j} \cdot 0^{-10} \end{array} \right.$$

arbitrary ordering.



Zigzag ordering.

- lot of param sharing still happening

PixelSNAIL:

Masked Attn + Conv.

Class conditional PixelCNN

- One-hot encoding of label
- force label image to be generated

Disadvantages of masked auto-regressive models

- sampling each pixel: 1 forward pass
- it takes time to generate  $16 \times 32 \times 32$  on GPU

How to overcome?

- breaking auto-regressive pattern
- make groups and do parallelization
- Speedup by caching activations

How to get latent rep of Pixel CNN?

- use fisher score

$$J(x; \theta) = P_\theta \log P_\theta(x)$$

### L3- flow models

- mixture of gaussians (Refer DeepBayes)
- Practical Param of flows.
  - cumulative density fun
  - eg: mix of gaussians
- nn each layer is a flow, then  
sequencing of layer = flow
  - Can't use ReLU as we can't be inverted.
  - sigmoid can be used.
  - tanh not a good idea as depends on when to fit data

#### 2D flows

$$x_1 \rightarrow z_1 = b_\theta(x_1)$$

$$x_2 \rightarrow z_2 = b_\theta(x_1, x_2)$$

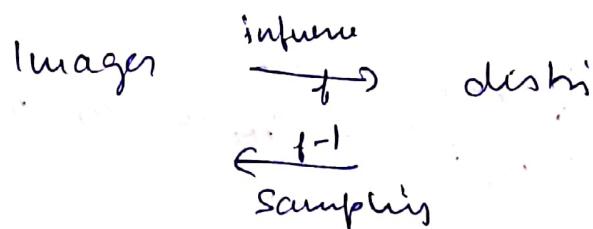
training obj:

$$\max_{\theta} \sum \log P_{z_1} \log(b_\theta(x_1)) + \log \left| \frac{dz_1}{dx_1} \right|$$

$$+ \log P_{z_2} \log(b_\theta(x_1, x_2)) + \log \left| \frac{dz_2}{dx_2} \right|$$

#### N-D flows

for high d data



#### N-D using Autoregressive flows.

$$x_1 \sim p_\theta(x_1)$$

$$x_1 \rightarrow f_\theta^{-1}(z_1)$$

- $x_2 \sim p_\theta(x_2|x_1)$   $x_2 = f_\theta^{-1}(z_2; x_1)$
- sampling is an invertible mapping
- lot of sampling for large image.

### Invert AR flow

$$x \rightarrow z$$

$$z_1 = f_\theta^{-1}(x_1) \quad x_1 = f_\theta(z_1)$$

$$z_2 = f_\theta^{-1}(x_2; z_1)$$

- IAF Sampling is faster

- Both AF & IAF is P due to same no of variables
- (soln. param sharing)
- again train with MLE

### NICE / Real NVP

splitting variables in half

$$z_{1:d/2} \in \mathbb{R}^{d/2}$$

$$z_{d/2:d} \in \mathbb{R}^{d-d/2} = \text{sol}(x_1:d/2) + t_\theta(x_1:d/2)$$

so  $\Sigma$  can be WNs with no restriction

Papers to look at

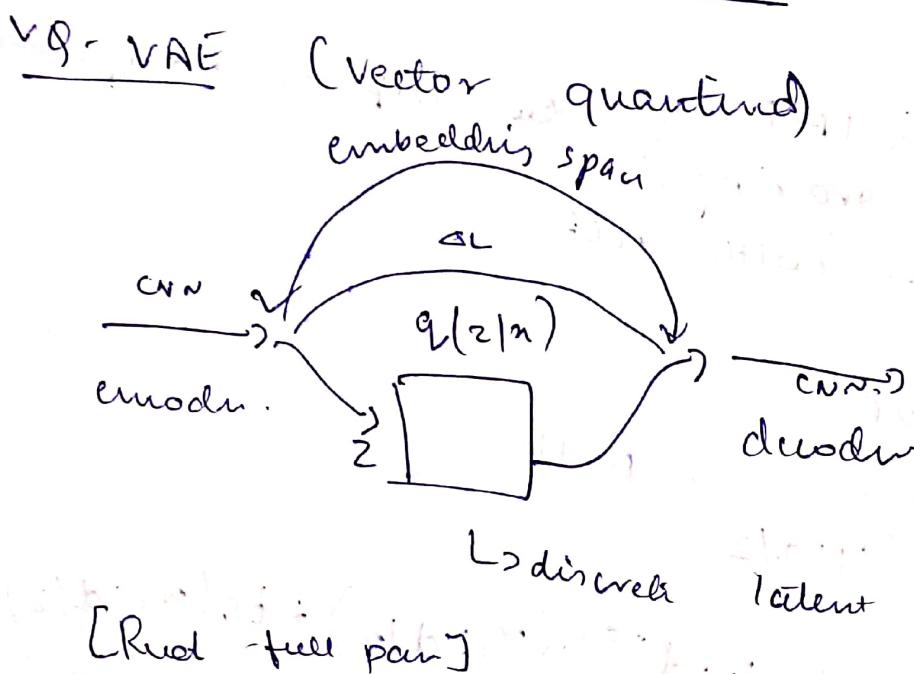
glow, flow++, bfjord.

## Latent Variable Models (

other)

- AP models are slow to sample as all pixels are oriented dependent
- make some part of obs. span independent on some LV
- Training Eu Opti  $\rightarrow$  Raybayers

additional work on VAEs.



### VQ-VAE 2.0

- two levels of VAE on encoder decoder
- hierarchical latent space

## GANs

original gan.

- loss

$$\min_G \max_D \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p_{\text{z}}} [\log(1 - D(h(z)))]$$

↓                      ↗  
clarify.            generate

Inception score for evaluation

- train take pre trained Inception vs a IM.

$$P(y|x)$$

- Marginal label distri  $P(y) = \sum_p P(y|p) P_g(p)$

$$IS = \exp(\mathbb{E}_{x \sim p_g} [D_{KL}[P(y|x) || P(y)]])$$

$$= \exp(H(y) - H(y|x))$$

- The more realistic the image, the higher the inception score.

Frechet Inception distance

- distance w.r.t mean of a feature vector of pre trained & a new data point. It looks at covariance over generated data & over training data.

$$d^2((\mu_c, \Sigma_c), (\mu_w, \Sigma_w)) = \| \mu - \mu_w \|^2_2 + \text{Tr}(\Sigma_c + \Sigma_w - 2(\Sigma_c \Sigma_w)^{-1})$$

- In GANs, there isn't an optimal metric to minimize at first.

## Some theory behind NNAs

Bayes Optimal Discriminators

what is the optimal discriminator given gen  $\mathbf{G}_1$  tree distri?

$$V(\mathbf{G}_1, D) = \mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p(z)} [\log (1 - D(\mathbf{G}_1(z)))]$$

$$= \int_x P_{\text{data}}(x) \log D(x) dx + \int_z p(z) \log (1 - D(\mathbf{G}_1(z)))$$

$$= \int_x P_{\text{data}}(x) \log D(x) dx + \int_z P_g(z) \log (1 - D(z)) dz$$

$$= \int_x [P_{\text{data}}(x) \log D(x) + P_g(x) \log (1 - D(x))] dx$$

$$\nabla_y [\alpha \log y + b \log (1-y)] = 0 \Rightarrow y^* = \frac{a}{a+b} \quad \forall [a,b] \in \mathbb{R}^2$$

$$\Rightarrow D^*(x) = \frac{P_{\text{data}}(x)}{(P_{\text{data}}(x) + P_g(x))} \rightarrow \text{optimal discr}$$

$$\therefore V(\mathbf{G}_1, D) = \mathbb{E}_{x \sim P_{\text{data}}} [\log D^*(x)] + \mathbb{E}_{x \sim P_g} [\log (1 - D^*(x))]$$

$$= \mathbb{E}_{x \sim P_{\text{data}}} \left[ \log \frac{P_{\text{data}}(x)}{P_{\text{data}}(x) + P_g(x)} \right] + \mathbb{E}_{x \sim P_g} \left[ \log \frac{P_g(x)}{P_g(x) + P_{\text{data}}(x)} \right]$$

$$= -\log(a) + KL \left( P_{\text{data}} \parallel \left( \frac{P_{\text{data}} + P_g}{2} \right) \right) +$$

$$KL \left( P_g \parallel \left( \frac{P_{\text{data}} + P_g}{2} \right) \right)$$

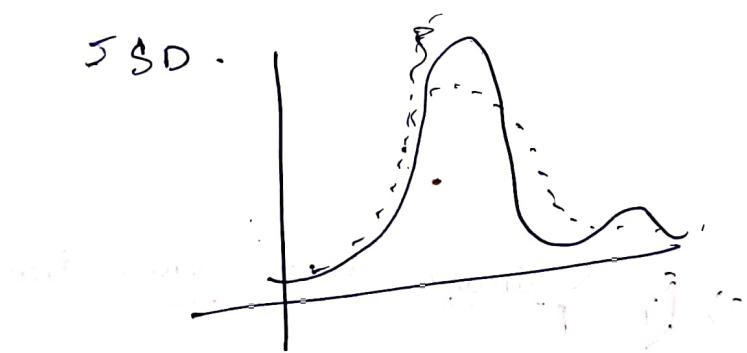
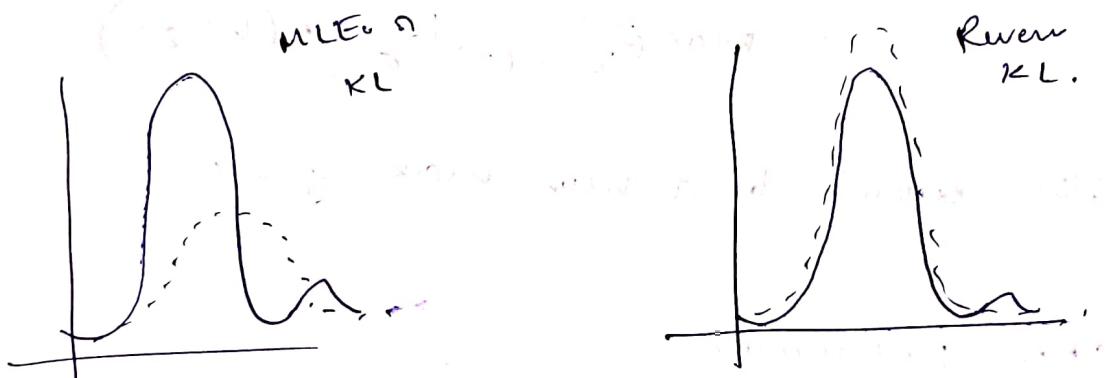
$$JS D = \frac{1}{2} \sum P_{\text{data}} \ln \frac{P_{\text{data}}}{P_g}$$

Assuming the discr. fully optimizes, then the gen. tries to optimize a type of KL div called SSD.

MLP tries to out-fit or find best scoring distri.

Rewire KL gives highest score on focusing on one mode.

SSD is something in b/w MLE & RKL.



Discriminator saturation

- If optimal  $D$  is reached, then the gen. can't be trained as  $D$  generates close to 0 grad.

① So alternative optimisation is done.

② or Non saturating formulation,

$$L^{(n)} \approx -L^D$$

In this case

$$L^{(G)} = -L^D = \min_G E_{z \sim p(z)} \log (1 - D(G(z)))$$

$$L^{(D)} = -E_{x \sim P_{\text{data}}} [\log D(x)] - E_{z \sim p(z)} [\log (1 - D(G(z)))]$$

↑  
non zero sum  
so want grad.

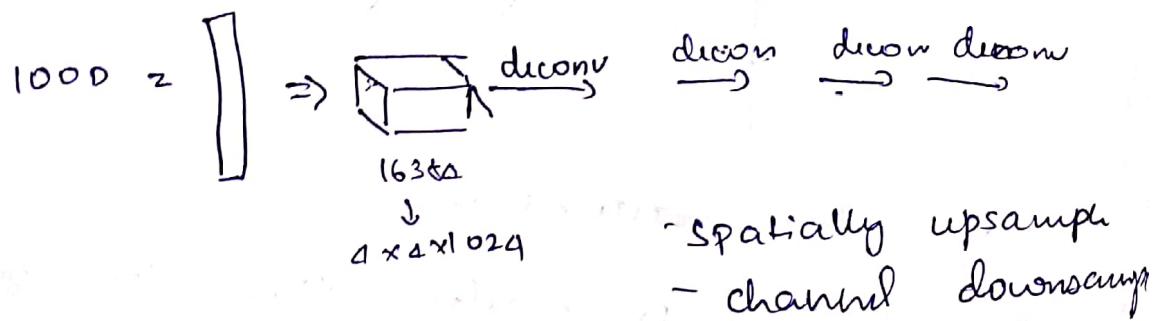
$$L^{(G)} = -E_{z \sim p(z)} \log (D(G(z)))$$

$$\max_G E_{z \sim p(z)} \log (D(G(z)))$$

This would be a min max game.

## GAN progression

### DGAN



- Most decons are normalized.

They are used to tie  $G(0) \rightarrow \mathbb{R}^3$ .

$$\left( \frac{n}{127.0} \right)$$

tanh at final op.

Because sampling time, have to go back to rgb

- for dis
- remove max pooling, mean pooling
- Standard conv., & avg pooling at end.
- Leaky ReLU, ReLU for gen.  
(grad can become sparse)
- tanh gen, sig dis.
- BS used <sup>for prevent mode collapse</sup>
- BS not applied at start of G & D<sup>if D</sup>
- Small lr, small momen, bcs 128
- . vector arithmetic with DCGAN
- . (disent) DC feature map TL did better, than  
Supervised on CIFAR.
- still unstable training
- Brittle architecture / hyperparam.

Improved training of GANs (Salimans et al.)

- feature matching
- intermediate feature f in dis.
- pass real data forward, collect f. on dis, same on gen.
- $\| \mathbb{E}_{\mathbf{x} \sim P_{\text{data}}} f(\mathbf{x}) - \mathbb{E}_{\mathbf{z} \sim p(z)} f(g(z)) \|_2^2$
- force feature of gen to match real data.

- minibatch discr
  - Mode collapse  $\rightarrow$  all noise in distri is mapped to single point by gen  $\Sigma$  discr is not aware this is happening as it looks at samples independently
  - need to make discr aware that gen needs to produce diverse samples
  - let  $x_i$  be a sample in minibatch & it produces a feature  $f(x_i)$
  - A tensor is used to project  $f(x_i)$  to a matrix
 
$$T \in \mathbb{R}^{A \times B \times C} \quad M_i \in \mathbb{R}^{B \times C}$$

$$\downarrow$$
 consider  $x_i \Sigma x_j$ ,  
 $c_b(x_i, x_j) = \exp(-\text{norm}(o(x_i) - o(x_j)))$
  - Do this for every sample in minibatch, add all the distances. This gives a metric of how diff is a particular  $f(x_i)$
  - $o(x) = [o(x_1), \dots, o(x_B)]$   
 $o(x) \in \mathbb{R}^{n \times B}$ 
    - to next layer of dis,  
 $[f(x), o(x)]$

## Historical Averaging

- inspired from Nash Equi.
- On a complicated 2 player game, with the issue of saddle point optim, it will basically be circling around in orbits.
- params in gen & discr kept ~cloud have penalty to keep closer to historical avg. Keep backups of params

## One-sided label smoothing

- mid in discr loss
  - with optim D,  $D(\alpha) = \frac{(1-\alpha)P_{\text{data}}(x) + \alpha P_{\text{model}}(x)}{P_{\text{data}}(x) + P_{\text{model}}(x)}$
- this is removed because for mode collapsing, it'll help appearing, so it should not be included for gen.

## Virtual Batch Norm

- use a sub-bn to compute norm stats

$$g_i = [r_1, \dots, r_m] \quad \downarrow \mu, \sigma$$

$$x_i = [x^1, \dots, x^m] \quad \frac{x - \mu}{\sigma}$$

same ✓ for every such forward batch.

- But model could overfit if 1 mini batch is used
- To avoid this, append sample with ref bn  $\rightarrow [x^1, r^1, \dots, r^m]$

- Introduced semi-supervised for gan
- idea is that if dis takes some img & predict, in addition, can predict category
- Then a class loss in dis
- Inception score

Kantorovich distance

given a joint disti b/w  $P_r$  &  $P_g$ ,

$$W(P_r, P_g) = \inf_{\pi \in \Pi(P_r, P_g)} E_{(x, y) \sim \pi} [ \|x - y\| ]$$

dual version  
both are intractable

expectation over samples from joint disti  
if calc. for various joint disti over all possible, then it is EMD

Kantorovich-Rubinstein Duality:

$$W(P_r, P_g) = \sup_{\|f\|_L \leq 1} E_{x \sim P_r}[f(x)] - E_{x \sim P_g}[f(x)]$$

supremum over all Lipschitz fns.

$W_{GAN}$

- what is EMD and in  $W_{GAN}$
- how is it tractable?

given any family  $\omega$

$$\max_{\omega \in \mathcal{W}} \underbrace{\mathbb{E}_{x \sim P_r} [f_\omega(x)] - \mathbb{E}_{x \sim P_g} [f_\omega(x)]}_{\leq \sup_{\|f_\omega\| \leq K} \mathbb{E}_{x \sim P_r} [f_\omega(x)] - \mathbb{E}_{x \sim P_g} [f_\omega(x)]}$$

$$= K \cdot W(P_r, P_g)$$

→ considered approx Leibnitz, lower bound is optimized on the EMD

- In the paper, they clip the weights to do this.

$$\sum_i^m f_\omega(x_i) - \frac{1}{m} \sum_i^m f_\omega(g_\theta(z^{(i)}))$$

this is clipped to make approx 2-1.

- No clipping in this case, but a critic gen tries to max critic while discriminating mini critic.

GRAD-UP

gradient penalty for Leip..

- read paper
- (no weight clipping here)
- no batch norm in discr (followed in all games now)

negatives

- training slower (due to loss)

- discr on grad penalty

- heuristic could be unstable,

## Progressive GAN - Nvidia

- high quality samples
- progressively grows size of images  
for gen & dis, via trained layers of lower levels

## SN GAN (spectral normalization)

SN  $\rightarrow$  (dugbaus)

Basic idea

weight  $w = \frac{w}{\sigma(w)} \rightarrow$  target singular value

do this for every weight matrix

- Hinge loss used, dis became something like a SVM classi

$$V_D(\hat{G}, D) = \mathbb{E}_{x \sim q_{\text{data}}(x)} [\min(0, -1 + D(x))] +$$

$$\mathbb{E}_{z \sim p(z)} [\min(0, -1 - D(\hat{G}(z)))]$$

$$V_G(h, \delta) = -\mathbb{E}_{z \sim p(z)} [\delta(h(z))]$$

## Implementation

- power implementation  
for square matrix  
(Dugbaus)

## self Attn GAN (SAGAN)

- Read paper.

BigGAN

↪ orthogonal reg

$$\beta \| w^T w - I \|$$

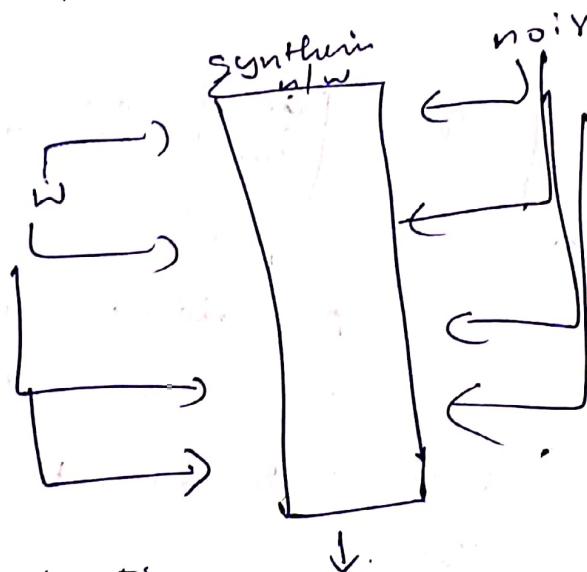
(and don't penalize  
diagonal terms)

- truncation trick
  - clip sampling test trim

StyleGAN

$$z \xrightarrow{f_{\text{cs}}}$$

thin as  $c \rightarrow$   
style bedn.

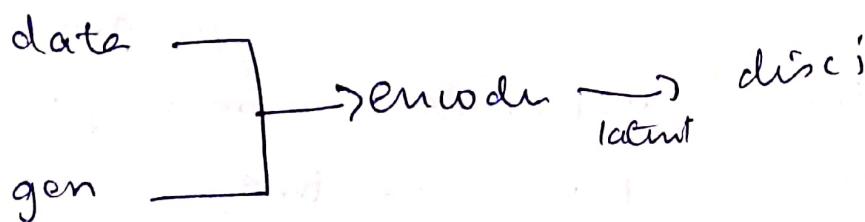


## Instance Normalisation

$$B \times H \times W \times C$$

- take avg stats over  $H \times W$

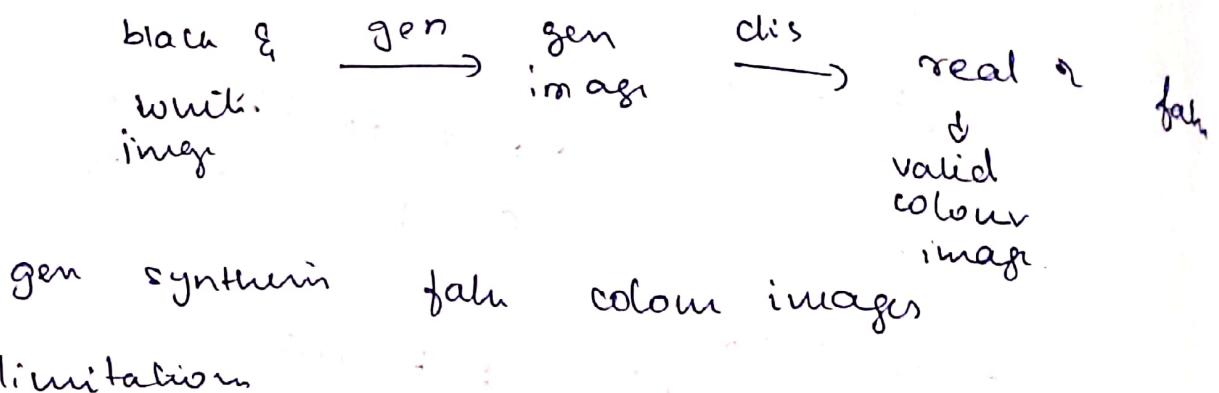
## Information Bottleneck



here discr. classifies  
in a lower dim

## Conditional GAN | pix2pix

- modifying an image  
(black & white to colour)



gen did not do the task it

was assigned. Bad objective.

- we both ① & ② to classify discrimination loss

$$\arg \min_{G,D} \mathbb{E}_{x,y} [\log D(x, G(y)) + \log (1 - D(x, y))]$$

$$G^* = \arg \min_G \max_D \mathcal{L}_{GAN}(G, D)$$

$$+ \lambda G_L(x_i)$$

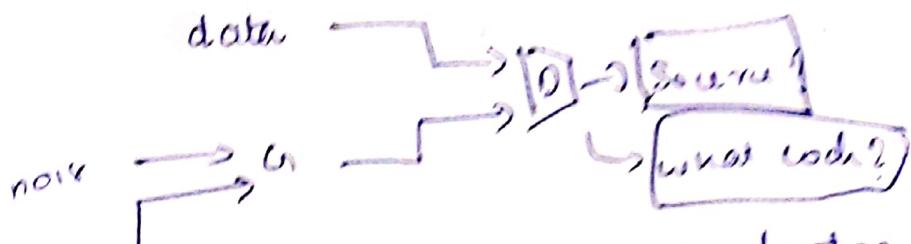
→ reconstruction loss

- shrinking discrimination - making it to focus on patches of image

# GRANS and representations

## InfoGAN

- info., type, rotation, width



code  
(high level  
aspects) → Independent factors should  
maximally explain variations  
in generated image

$$\text{max } I(c; z) = H(z) - H(z|c)$$

$$H(z|c) = H(c) - H(c|z)$$

$\downarrow$   $z = G(z, c)$   
variational lower bound to  
solve it as MI is intractable

$$I(c, G(z, c)) \approx H(c) - H(c|G(z, c))$$

$$= H(c) + E_{z \sim G} (z, c) [E_{c' \sim P} (c'|z)]$$

$$[\log P(c'|z)]$$

$$= H(c) + E_{z \sim G(z, c)} [E_{c' \sim P(c|z)} [\log Q(c'|z)]]$$

$$+ \underbrace{D_{KL}(P(\cdot|z) || Q(\cdot|z))}_{\geq 0}$$

$$\geq H(c) + E_{z \sim G(z, c)} [E_{c' \sim P(c|z)} [\log Q(c'|z)]]$$

$\uparrow$   
MLE

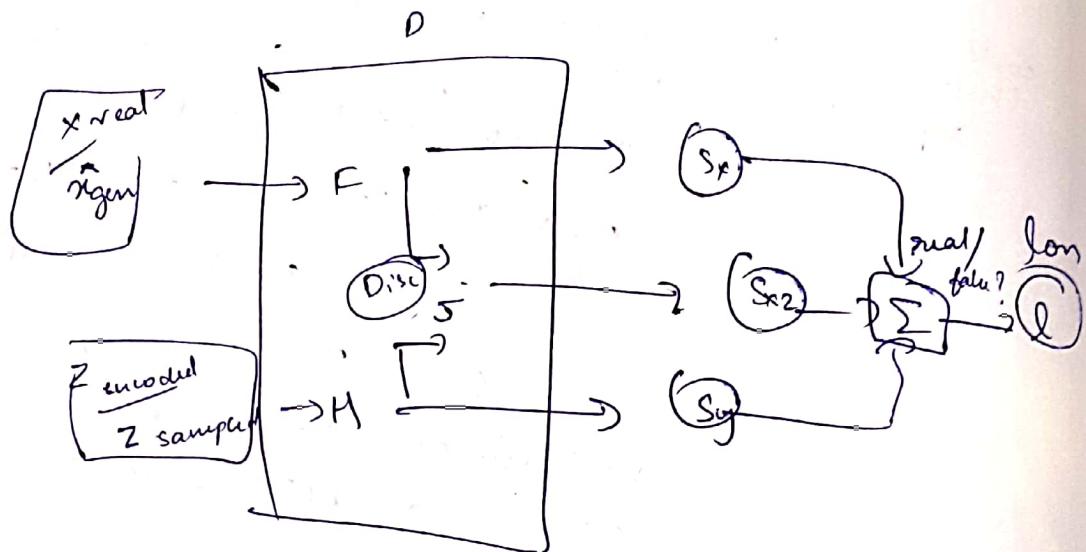
- In BigGAN,

if the gen noise  $g_z$  is changed

$\Rightarrow$  concat ( $[N(0, 1)]^{100}$ , Uniform ( $[0, 1]$ ))

optimal for different classes.

- BigGAN



Take find Encoder for Rep learning

- nearest neighbors in latent space given similar image

GANs as energy models

align energy  $E(x)$  to every  $x$

$$P(x) = \frac{1}{Z} e^{-E(x)}$$

$$Z = \sum_{x \in X} e^{-E(x)}$$

lower the energy  $\Rightarrow$  all possible world state (normalizing)

-  $x$  can be really large

for an RGB  $28 \times 28 \times 3$  img,

$$256 \times (28 \times 28 \times 3) = 256 \text{ GB}$$

low energy for real looking image & high  
for those that don't look real.

$$\max_{\theta} \mathbb{E}_{x \sim P_{\text{data}}} [\log P_{\theta}(x)]$$

$$= \max_{\theta} \mathbb{E}_{x \sim P_{\text{data}}} \log \left( \frac{e^{-E_{\theta}(x)}}{Z(\theta)} \right)$$

$$= \max_{\theta} \mathbb{E}_{x \sim P_{\text{data}}} [-E_{\theta}(x)] - \log Z(\theta)$$

(minimizing energy)      (hand to computer)

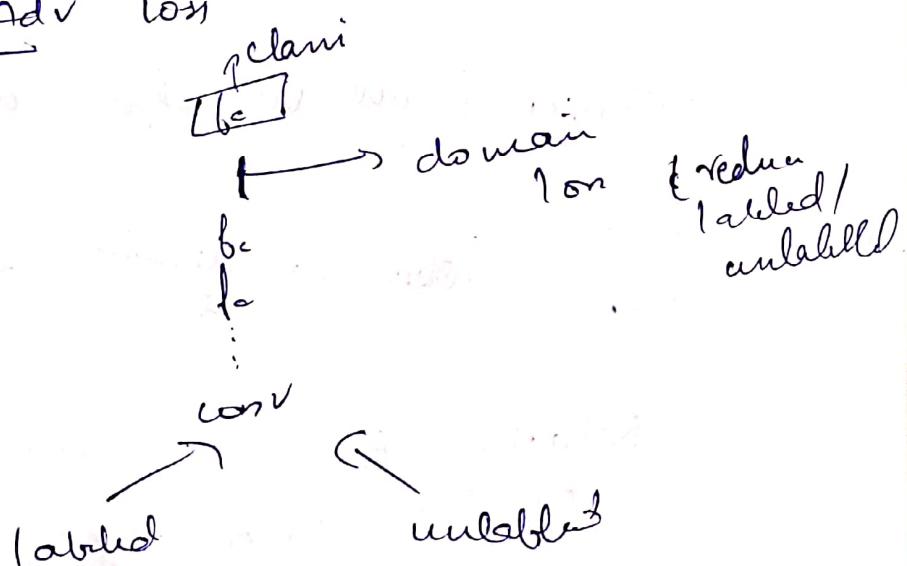
variational lower bound for  $Z$

$$- \log Z \leq \min_{\varphi} \mathbb{E}_{x \sim P_{\text{data}}} [E_{\theta}(x)] - H(q_{\varphi})$$

$$\geq \max_{\theta} \mathbb{E}_{x \sim P_{\text{data}}} [-E_{\theta}(x)] + \underbrace{\min_{\varphi} \mathbb{E}_{x \sim P_{\text{data}}} [E_{\theta}(x)] - H(q_{\varphi})}_{\text{gan obj}}$$

- not easy to compute  $H(q_{\varphi})$   
but can be computed for softmax

TL with Adv loss



- can also do this for facemask.

## Self-Supervised

### Cognitive principles

- Reconstruct from corrupted
- Visual common sense task
- Contrastive learning

### Reconstruct from corrupted

- denoising AE
  - Additive, isotropic gaussian noise
  - masking noise
  - Salt & pepper noise
  - stacked denoising AE  
(iterate on features)
- predicting missing piece
  - mask a piece, reconstruct
  - Joint loss
- Predicting one view from another
-  → crop → denois  $\underset{AE}{\rightarrow}$  grayscale
- Relative position of image patches
  - predict relative position

- either both spatially & columnar patches shouldn't overlap
  - ↳ to prevent chromatic aberration
- Solving jigsaw puzzles
  - take patches from random crop, shuffle and predict proper img
- Rotation
  - '0, 180' worked best on CIFAR

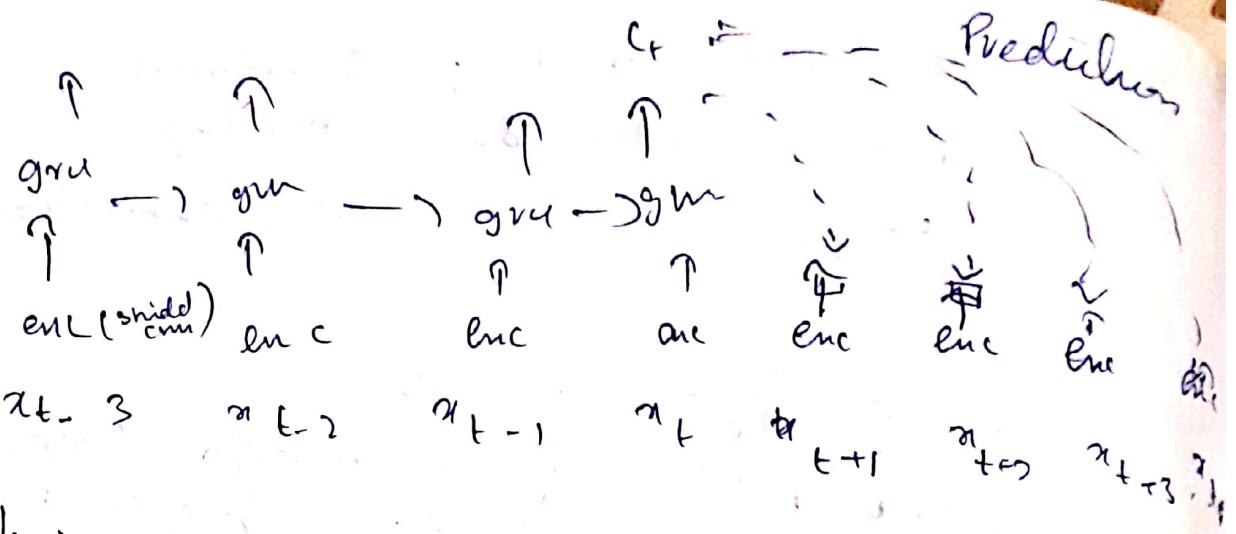
## Contrastive learning

- word2vec
  - ↳ make word matrix
  - apply SVD
    - sparsity
    - computation
    - Infrequent word
    - noise from frequent
  - ↳ n-gram models

## Contrastive Predictive Coding

$\max_{\theta} \mathbb{E}_{z \sim p(z)} \mathbb{E}_{c \sim p(c|z)} \mathbb{E}_{\text{encodes } z \text{ in } c}$ 
  
 Jensen's Inequality

$$\max_{\theta} \mathbb{E}_{z \sim p(z)} \mathbb{E}_{c \sim p(c|z)} \text{MI}(c; z)$$



## Instance Discrimination

- MoCo
- SimCLR (need full papers)

### MoCo

query  $\rightarrow$  encoder  $\rightarrow q$

Similarity  $\rightarrow$  C.L.

key  $\rightarrow$  encoder  $\rightarrow k_0, k_1, k_2, \dots$

momentum encoder  $\rightarrow$

historically averaged version of encoder

$\rightarrow$   $q, k$

## Summary of unsupervised learning models

### Autoencoders

- MADE
- PixelRNN / CNN
- WaveNet
- Video Pixel Networks
- GPT 2

Success due to:

- large batch sizes, more compute, ~~less~~  
- better stability, training
- wider, deeper
- conditional models on auxiliary variables
- Fewer assumptions
- Architectures
  - masked conv
  - dilated conv
  - transformers
- loss on cross entropy

Futute:

- model parallelism, larger models
- same model for multiple modalities
- Fair sampling
- Architecture innovations

Negatives:

- no single layer of rep.
- sampling low for deployment
- No interpolation

## Flow models

- NICE
- WLOW
- Flow++

## LV models

- PixelVAE
- Scratch RNN
- World models

## Advantages

- Compressed Bottleneck Rep
- Approx density estimates
- Disentangled reps

## Dis:

- Blurry
- Factorised gaussian posterior or doesn't assumption may limit

## Futur:

- Better modules (cross entropy loss and weakly autoregressive)
- More powerful posteriors

## CiGANs

- original
- DCGAN
- Wgan
- StyGAN
- Biggan

Future:

- move fine grained
- Video generation
- L-1 constraints (new approach)
- Arch (Upsampling, downs.)
- Obj fun

Negatives:

- engineering tricks
- mode dropping
- Evaluation metrics

GANs or Density models?

more ~ same level of hard engineering details

- Blurry samples vs mode collapse  
(comprehension vs sample quality)
- GANs work well with less computation

For training density models:

- if only can obt density estimator, not Sampling, choose A.R.
- if both rep. & sampling, use VAEs

When gaus?

- good samples
- Large high quality images

## Contractive learning

- Dictionary look up task
- Predictive Coding & Instance Divergence
  - ↳ end to end
  - ↳ momentum
- CPCv2

Molto  
SimCLR

## Futur:

- gap not close b/w self-s & s
- down streak task TL ~~not~~ gains not that significant.

## Modeling futur in latent space

- internal world models that run on abstract space

## Semi Supervised

D<sub>U</sub> is an p(x)

D<sub>S</sub>: (x, y) ~ p(x, y)

run this to make p(x, y) better

## Core Concepts

- Confidence vs Entropy

\* Entropy minimization

- taking label date & making sure that classifier training on label date has minimal entropy on unlabeled date.

## 1 Pseudo Labelling

- make classifier to predict labels on unlabeled data, take confidence predictions, convert them to extra ground truth. (so called self training)

## 2 Virtual Adversarial Training

- Label consistency
  - Make sure that augmentations of sample have the same class
  - Pi Model, Temporal Ensemble, mean teach
- Regularisation
  - weight decay
  - Dropout
  - Data Aug
  - Unsupervised Data Aug (UDA)
  - MixMatch
- Co-training / Self-training / Pseudo labelling.  
(Noisy student)

Entropy minim



pi Model

