

BLG561E FALL 2021

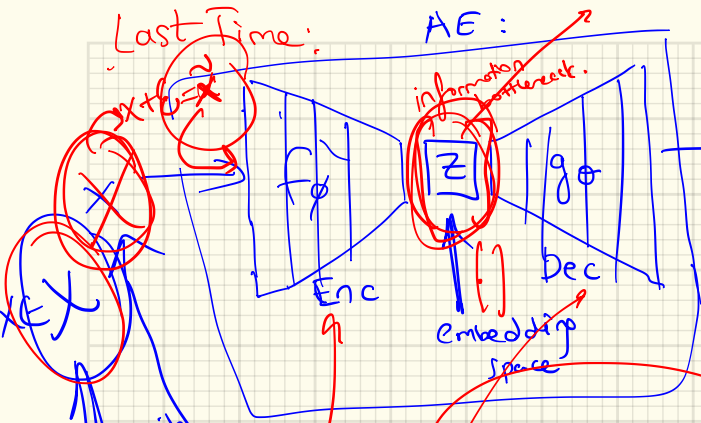
Deep Learning

07.12.2021

Görde ÜNAL

Last Time:

AE:



① Autoencoder hypothesis

$$\hat{x} = h_{\theta, \phi}(x) = g_{\theta}(f_{\phi}(x))$$

$$L_{rec} = \|h_{\theta, \phi}(x) - x\|_2^2$$

DAE Inference



$$L_{rec} = \|\hat{x} - x\|_2^2$$

② Loss

Limit capacities of NNs
Enc
Dec.

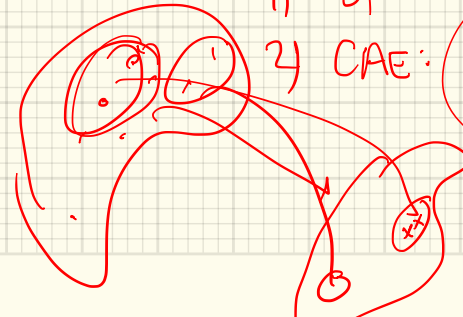
+ 2 Regularizers.

1) Sparse AE $\Rightarrow \|z\|_1$

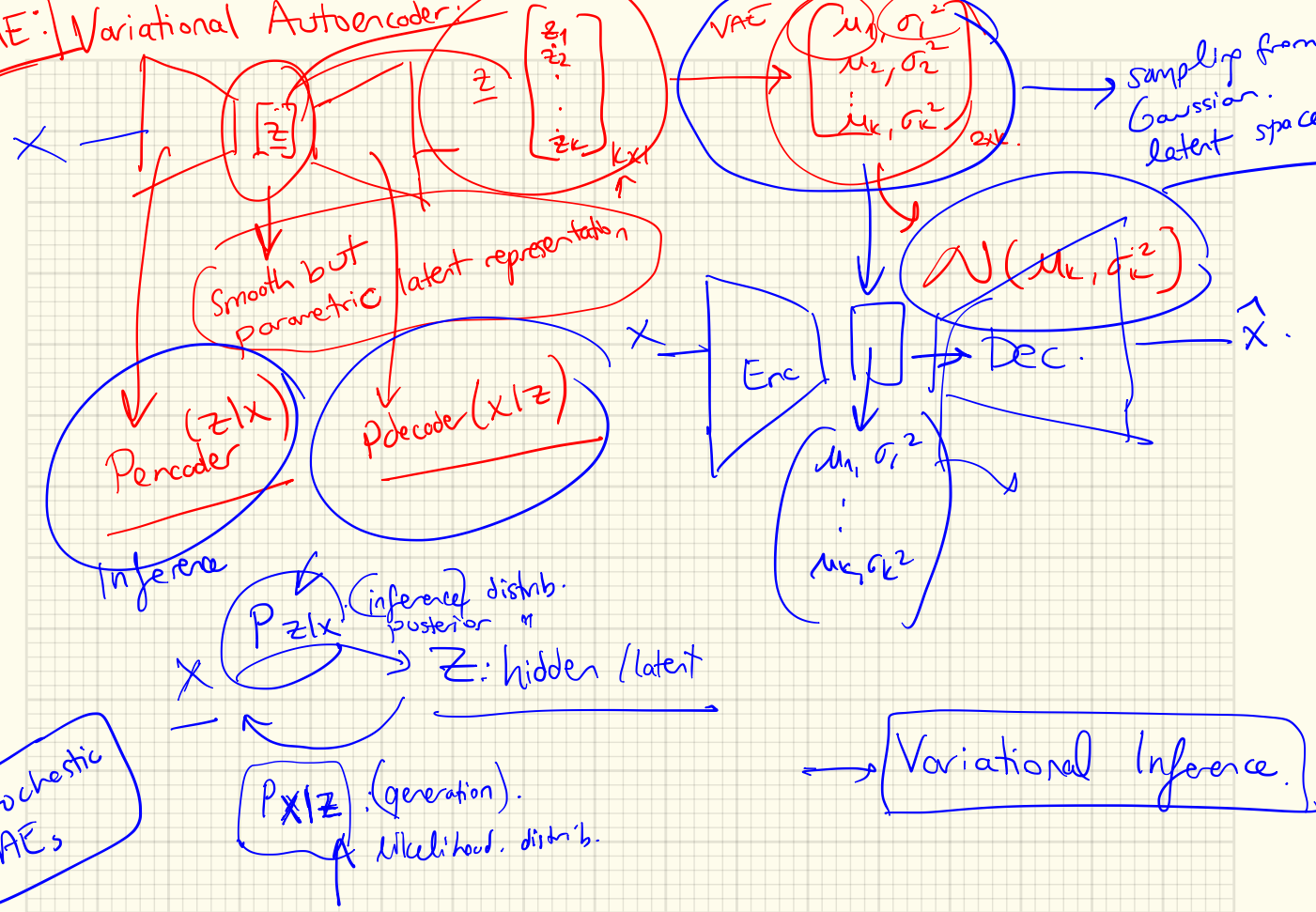
2) CAE: $\left\| \frac{\partial z}{\partial x} \right\|_F^2$ small.

3) Denoising AE

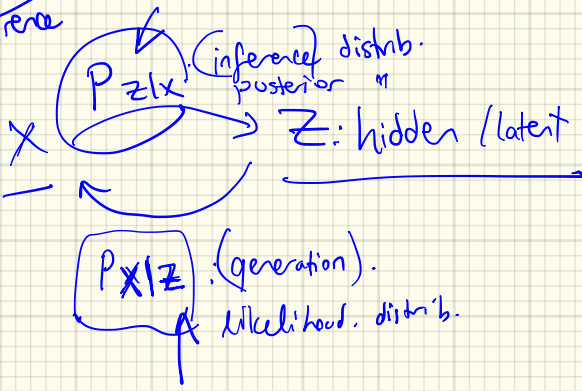
$\frac{\partial f}{\partial x} \rightarrow$ Jacobian matrix



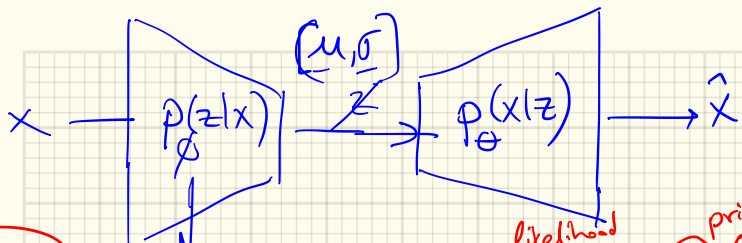
VAE: Variational Autoencoder



Stochastic AEs

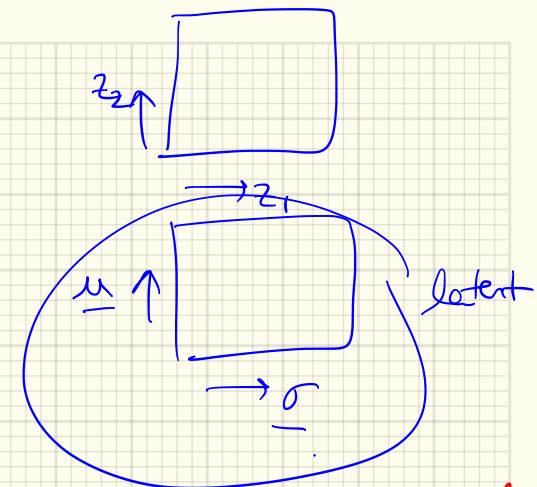


Variational Inference



posterior distrib

$$p(z|x) = \frac{p(x|z)}{p(x)} = \frac{\underbrace{p(x|z)}_{\text{likelihood}} \underbrace{p(z)}_{\text{prior}}}{\underbrace{p(x)}_{\text{Evidence}}}$$



$$p(x) = \int p(x|z) p(z) dz$$

↑ ↑ z unobserved for each input x.

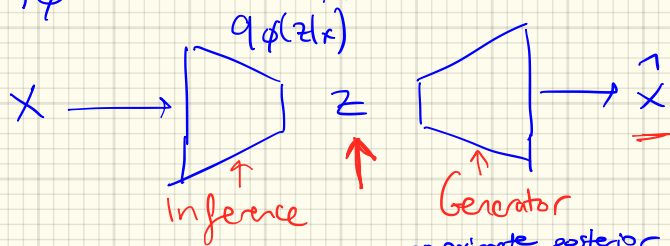
huge # dimensions

→ Intractable integral.

Intractable both $p(z|x)$ & $p(x)$

→ Variational Inference tries to approximate $p(z|x)$ by an approximate posterior distribution.

→ $q_\phi(z|x)$: approximate posterior distribution:



$$q_\phi(z|x) \stackrel{\text{approximate}}{\sim} \stackrel{\text{true}}{p}(z|x)$$

$$KL(p_1 || p_2) = \sum p_1 \log \frac{p_1}{p_2}$$

★ VAE wants to make q as "close" as possible to the true $p(z|x)$

$q_\phi(z|x)$: tractable distribution

posterior distrib.
(intractable)

eg. a Gaussian distribution

$$\phi = (\mu_i, \Sigma_i)$$

$$\begin{bmatrix} \sigma_1^2 & & 0 \\ & \sigma_2^2 & \\ 0 & & \sigma_k^2 \end{bmatrix}$$

: uncorrelated across dimensions of the hidden code.

Q. How do I measure closeness of 2 probability distributions?

A. Resort to Information-theoretic Divergence measures

like KL or Jensen-Shannon Divergence. ⇒

\Rightarrow Minimize $\underbrace{KL(q \parallel p)}_{\substack{\text{Kullback-Leibler} \\ \text{div. between } q \text{ \& } p.}} \triangleq \sum q \log \frac{q}{p} = - \sum q \log \frac{p}{q}$

$\min KL(q_p(z|x) \parallel p(z|x)) \leftarrow$

$KL = - \sum_z q(z|x) \log \frac{p(z|x)}{q(z|x)} = - \sum_z q(z|x) \log \frac{p(z|x)}{q(z|x)} \cdot \frac{1}{p(x)}$

$= - \sum_z q(z|x) \left[\log \frac{p(x,z)}{q(z|x)} - \log p(x) \right]$

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

Any pdf $p(x)$
 Satisfies
 i) $p(x) \geq 0$
 ii) $\sum_x p(x) = 1$

$\rightarrow \log p(x) = \underbrace{KL(q(z|x) \parallel p(z|x))}_{\geq 0} + \sum_z q(z|x) \log \frac{p(x,z)}{q(z|x)}$

const.

KL div. properties ≥ 0

$\rightarrow \log p(x) \geq \mathcal{L}_{VB}$

\mathcal{L}_{VB} : called the Variational Lower Bound for $\log p(x)$.

want to $\max \log p(x)$ $\log p(x) \geq \mathcal{L}_{VB}$ \rightarrow instead $\equiv \max \mathcal{L}_{VB}$: maximize the variational lower bound

$$\rightarrow \mathcal{L}_{VB} = \sum_z q(z|x) \log \frac{p(x|z) p(z)}{p(z|x)}$$

$$\mathcal{L}_{VB} = \underbrace{\sum_z q(z|x) \log p(x|z)}_{\text{Data Term}} + \underbrace{\sum_z q(z|x) \log \frac{p(z)}{q(z|x)}}_{\text{Regularizer Term}}$$

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

model distribution $\sim \mathcal{N}(0, 1)$

$$\max \mathcal{L}_{VB} = E_{q(z|x)} [\log p(x|z)] - KL(q(z|x) \parallel p(z))$$

$\log p(x|z)$
log likelihood of x given z

Regularizer Term: \sim minimizes the KL approx posterior $q(z|x)$ & $p(z)$

eg. Data Term: $\mathcal{L}_{data} + \lambda \mathcal{L}_{reg.}$
Reconstruction Loss

→ or minimize \mathcal{L}_{VB} : ϕ θ

Given $\{x^i\}_{i=1}^m$
min

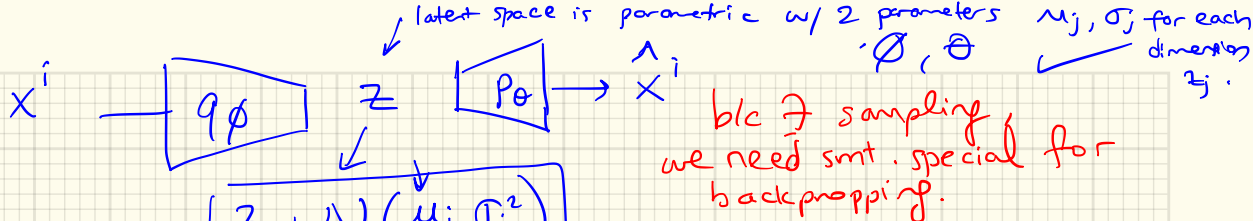
$$\mathcal{L}(\theta, \phi) = \sum_{i=1}^m \underbrace{\mathcal{L}_{\text{rec}}(\text{data})}_{-E_{z \sim q_{\phi}(z|x^i)} [\log p_{\theta}(x^i|z)]} + d \cdot \underbrace{\text{KL}(q_{\phi}(z|x^i) || p(z))}_{\mathcal{N}(0,1)}$$

$$+ d \sum_j \text{KL}(q_j(z|x^i) || \mathcal{N}(0,1))$$

j (for every dimension of the latent space)

$$\mathcal{N}(\mu_j, \sigma_j^2)$$

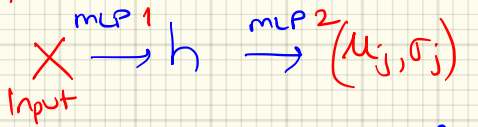
VAE paper : (Kingma, Welling, 2014)



We need

Reparameterization Trick:

Say for example, we have a 2-layer MLP Encoder:



$h = \tanh(\underline{w}^1 X + \underline{b}^1)$

activation at 1st layer:

$\mu_j = \underline{w}^2 h + b^2$

$\log \sigma_j^2 = \underline{w}^3 h + b^3 \rightarrow \sigma_j^2 = \exp(\log \sigma_j^2)$

ensures σ^2 is +ve.

now, we have (μ_j, σ_j^2) from the encoder ϕ

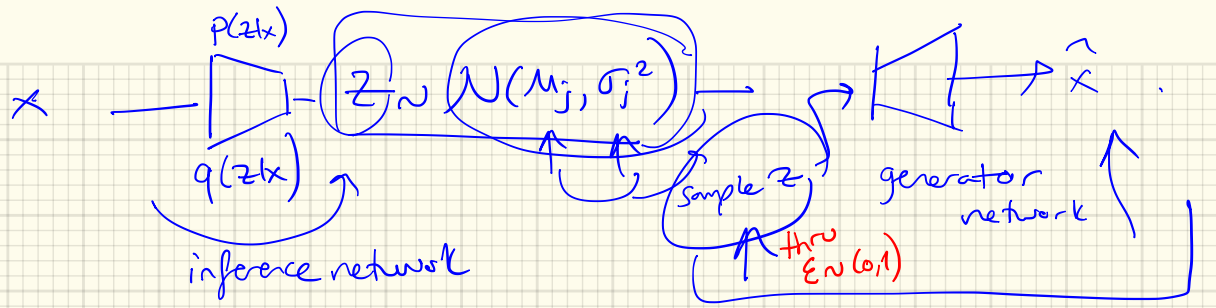
$\mathcal{N}(\mu_j, \sigma_j^2) \rightarrow \underline{z}_j$: sampled by

to keep the variance > 0

Helps us to be able to differentiate w.r.t. μ_j 's & σ_j 's during backprop.

$z_j = \mu_j + \sigma_j \cdot \epsilon$, $\epsilon \sim \mathcal{N}(0, 1)$

ϵ : sampled from standard normal.



Summarize VAE →

Loss terms is based on Variational Lower Bound:

- 1) Data Likelihood \approx Reconst loss
- 2) Regularizer: KL loss $\xrightarrow{\text{btw}}$ approximate inference distrib & a Gaussian distrib.

Generative Adversarial Networks (GANs) / Implicit Generative Models

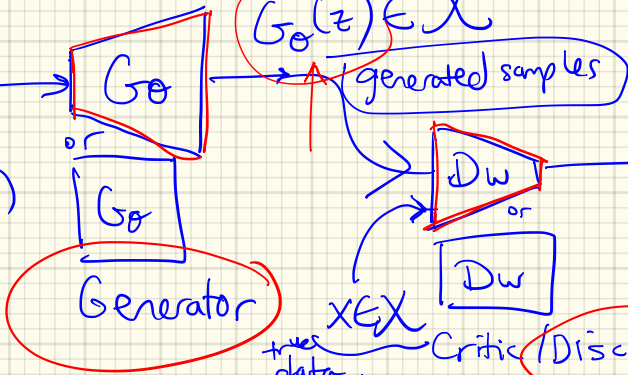
Problem: To learn a generative model of a distribution of data points

$P_{\text{data}}(x)$ implicitly. $x \in X$

Given a set of data sample x from X , learn a data distribution in the form of a deep NN $G_{\theta}(\cdot)$, w/ parameters θ .

G takes as input a Noise sample

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



pushes the GAN towards realistic looking outputs.

GAN: creates a ZERO-sum Game between the Generator & Discriminator

Generator tries to maximize error of the D network \rightarrow zero-sum game.

Vanilla GAN: $D_w(x) \in [0, 1)$ w/ BCE loss.
 (Goodfellow 2014)

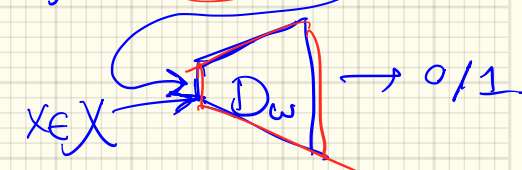
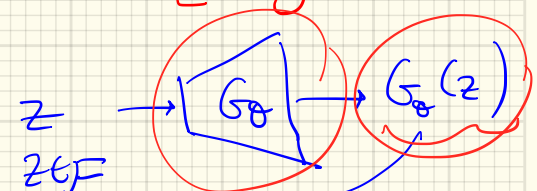
Optimization problem of the vanilla GAN

$$\min_{\theta} \max_w E[\log D_w(x)] + E[\log(1 - D_w(G_{\theta}(z)))]$$

$\times \sim \text{update}^{(i)}$ \uparrow $\text{z} \sim \mathcal{N}$ \uparrow

minimize $\log(1-D)$
 for $D \approx 0$
 very flat area } Vanishing gradient

$\mathcal{L}(\theta, w)$
 $E[-\log D_w(G_{\theta}(z))]$



steeper,
 so larger
 gradients for
 small D.

$$\left. \begin{array}{l} \text{Discriminator: } \max L^D = L(\theta, w) \\ \text{Generator: } \min L^G = -L^D \end{array} \right\} \underbrace{L^D + L^G = 0}_{\text{zero-sum game.}}$$

GAN Loss used in practice

$$\min \log(1 - D(G(z))) \rightarrow \min -\log D(G(z))$$

Loss

$$\min_{\theta} \max_w E_x [\log p_w(x)] - E_z [\log D(G(z))]$$

(Goodfellow 2014) : Generator minimizes the ^(JS) Jensen-Shannon divergence between the true distribution & the generator distribution

WGAN :



Wasserstein GAN : [Arjovsky et al]

Instead JS divergence, they optimize between the true distrib & generator distrib.

WGAN leads to the loss function

WGAN loss:

2

min

θ

↑

max

w

$$E_{x \sim p_{data}} [D_w(x)] - E_{z \sim F} [D_w(G_\theta(z))]$$

$L(G, w)$

Note:

Generator only operates on this term.

WGAN

E_{μ}
Earth-mover distance

(Wasserstein - 1) metric.



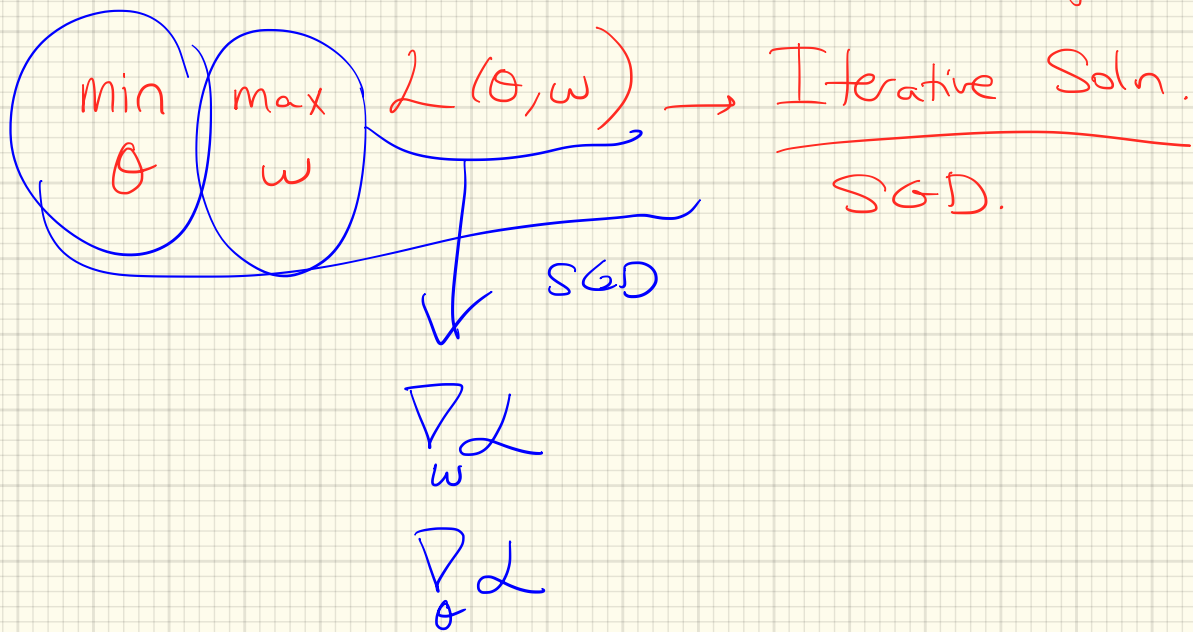
3

Optimization

? \equiv Training of GANs?

→ Is $L(\theta, w)$ a $\underbrace{\text{convex}}_{\text{in } \theta}$ - $\underbrace{\text{concave}}_{\text{in } w}$ function? No!

Optimization: Via SGD or its variants Train parameters θ & w in an alternate fashion.



for $t \in \{1, \dots, T-1\}$; with $\nabla_{w,t} \mathcal{L}(\theta_t, w_t)$ ^{fixed}

$\nabla_{\theta,t} \mathcal{L}(\theta_t, w_t)$ ^{fixed.}

Gradient Ascent for w :

$$w_{t+1} = w_t + \eta \nabla_{w,t} \mathcal{L}$$

Gradient Descent for θ :

$$\theta_{t+1} = \theta_t - \eta \nabla_{\theta,t} \mathcal{L}$$

Q. Are there any guarantees for an optimal solution to this problem?

(Other than in a convex-concave setting) No guarantees for an optimal solution.

★ GANs are difficult to train \rightarrow K even unstable!

LS-GAN:
(Least-Squares GAN) Rather than cross-entropy loss in the vanilla GAN;

Loss functions:

$$D: \min_w E_{x \sim p_{\text{data}}} \left[\underbrace{(D_w(x) - 1)}_{\substack{\uparrow \\ \text{true data}}} \right]^2 + E_{z \sim p(z)} \left[\underbrace{(D_w(G_\theta(z)) - 1)}_{\substack{\uparrow \\ \text{Generated}}} \right]^2$$

$\underbrace{z \sim p(z)}_{\mathcal{N}(0,1)}$

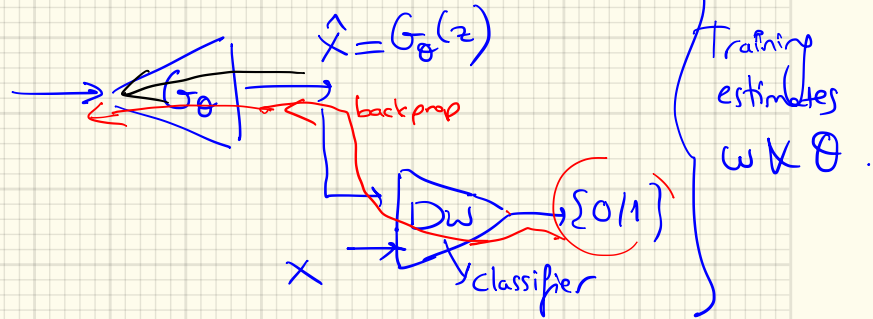
$$G: \min_\theta E_{z \sim p(z)} \left[(D_w(G_\theta(z)) - 1)^2 \right]$$

LSGAN:

Note: A variant of GAN \rightarrow \exists tons of them

GAN Training Stage:

$$z \in \mathcal{N}(0,1)$$



Now $w, \& \theta$ are learned:

GAN Inference Stage:

want to Generate new Samples:

~~No D_w network~~

Sample a $z \sim \mathcal{N}(0,1)$

