

Deep Learning for Computer Vision

Combining VAEs and GANs

Vineeth N Balasubramanian

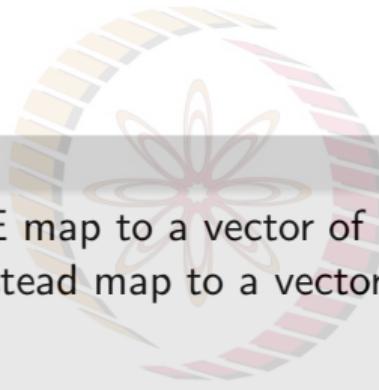
Department of Computer Science and Engineering
Indian Institute of Technology, Hyderabad



Review: Questions

Questions

- Why does the encoder of a VAE map to a vector of means and a vector of standard deviations? Why does it not instead map to a vector of means and a covariance matrix?

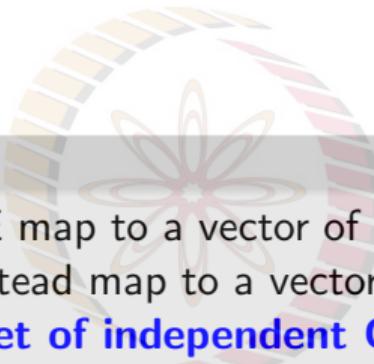


NPTEL

Review: Questions

Questions

- Why does the encoder of a VAE map to a vector of means and a vector of standard deviations? Why does it not instead map to a vector of means and a covariance matrix?
We are explicitly learning a set of independent Gaussians, which makes the learning easier - and of course, it works!

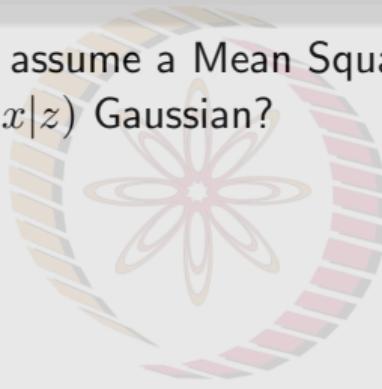


NPTEL

Review: Questions

Questions

- What about the decoder? If we assume a Mean Squared Error for the reconstruction loss, what is the covariance of the $p(x|z)$ Gaussian?



NPTEL

Review: Questions

Questions

- What about the decoder? If we assume a Mean Squared Error for the reconstruction loss, what is the covariance of the $p(x|z)$ Gaussian?

Equivalent to modeling $p(x|z)$ as Gaussian with identity covariance; in this case, decoder output is mean $\mu(t)$ and, therefore, for an example x_i , you get the following reconstruction loss:

$$\begin{aligned}-\log(p(x_i|t_i)) &= -\log \left(\frac{1}{\sqrt{(2\pi)^k |I|}} \exp \left(-\frac{1}{2} (x_i - \mu(t_i))^T I (x_i - \mu(t_i)) \right) \right) \\ &= \frac{1}{2} \|x_i - \mu(t_i)\|^2 + \text{const.}\end{aligned}$$

This is MSE!

VAE vs GAN

VAE

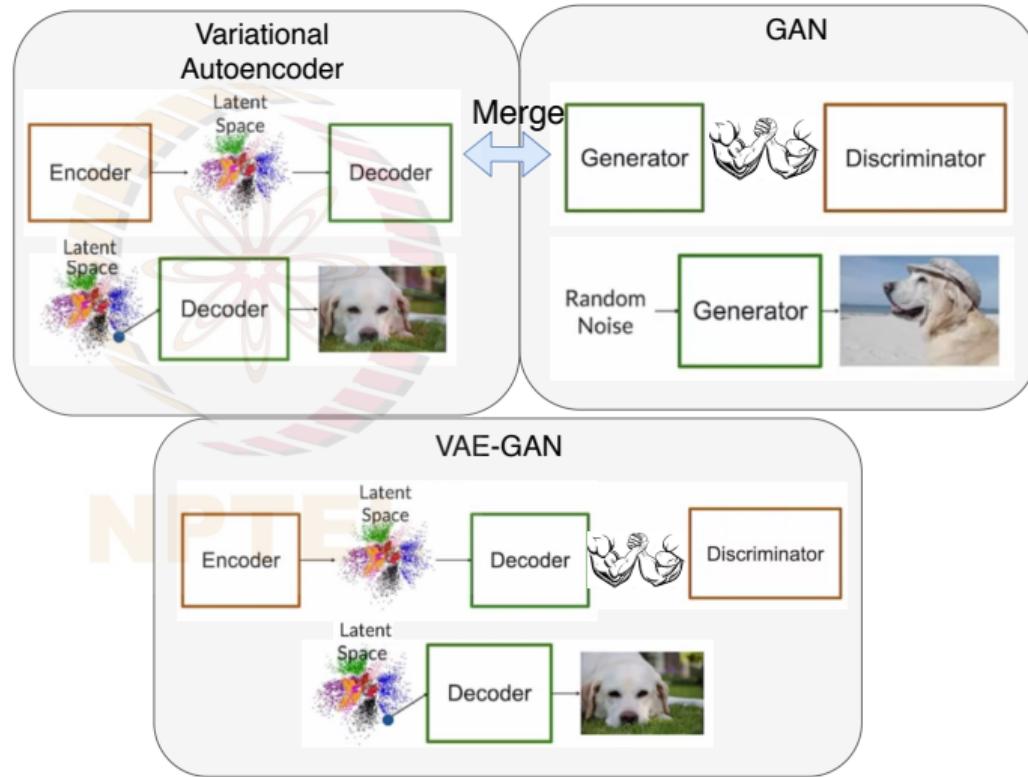
- **+ve** Learns an inference machine by mapping data to a latent space with distribution of choice enabling **fast/efficient inference**
- **-ve** Tends to distribute probability mass diffusely over data space resulting in blurred/**low quality image** samples

GAN

- **+ve:** Bypass inference and learn generative model that produces high quality samples without sacrificing sampling speed
- **-ve:** **Lacks** an effective **inference mechanism** preventing from reasoning about data at an abstract level

Can we combine VAEs and GANs?

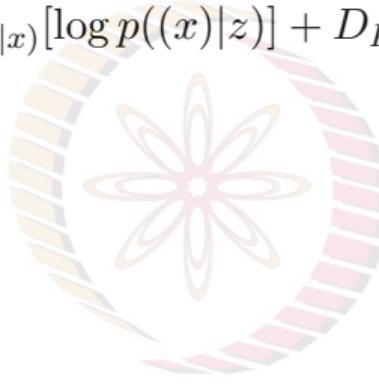
- **Solution:** Bridge the gap between VAEs and GANs, learn models that generate high-quality samples along with an effective inference network



VAE Limitations

Let's recall the VAE objective:

$$\mathcal{L}_{VAE} = -\mathbb{E}_{q(z|x)}[\log p((x)|z)] + D_{KL}(q(z|x)||p(z))$$

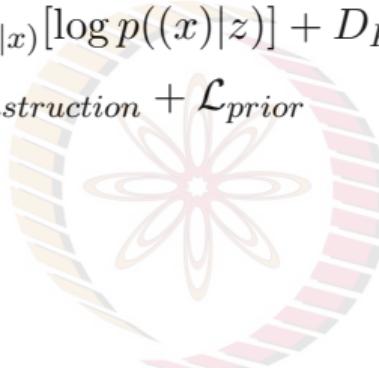


NPTEL

VAE Limitations

Let's recall the VAE objective:

$$\begin{aligned}\mathcal{L}_{VAE} &= -\mathbb{E}_{q(z|x)}[\log p((x)|z)] + D_{KL}(q(z|x)||p(z)) \\ &= \mathcal{L}_{reconstruction} + \mathcal{L}_{prior}\end{aligned}$$

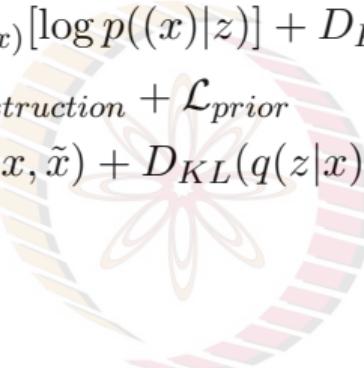


NPTEL

VAE Limitations

Let's recall the VAE objective:

$$\begin{aligned}\mathcal{L}_{VAE} &= -\mathbb{E}_{q(z|x)}[\log p((x)|z)] + D_{KL}(q(z|x)||p(z)) \\ &= \mathcal{L}_{reconstruction} + \mathcal{L}_{prior} \\ &= MSE(x, \tilde{x}) + D_{KL}(q(z|x)||p(z))\end{aligned}$$



NPTEL

VAE Limitations

Let's recall the VAE objective:

$$\begin{aligned}\mathcal{L}_{VAE} &= -\mathbb{E}_{q(z|x)}[\log p((x)|z)] + D_{KL}(q(z|x)||p(z)) \\ &= \mathcal{L}_{reconstruction} + \mathcal{L}_{prior} \\ &= MSE(x, \tilde{x}) + D_{KL}(q(z|x)||p(z))\end{aligned}$$

Mean-Squared Error

NPTEL

VAE Limitations

Let's recall the VAE objective:

$$\begin{aligned}\mathcal{L}_{VAE} &= -\mathbb{E}_{q(z|x)}[\log p((x)|z)] + D_{KL}(q(z|x)||p(z)) \\ &= \mathcal{L}_{reconstruction} + \mathcal{L}_{prior} \\ &= MSE(x, \tilde{x}) + D_{KL}(q(z|x)||p(z))\end{aligned}$$

Mean-Squared Error

- Assumes signal fidelity is independent of temporal/spatial relationships → does not hold for images

NPTEL

VAE Limitations

Let's recall the VAE objective:

$$\begin{aligned}\mathcal{L}_{VAE} &= -\mathbb{E}_{q(z|x)}[\log p((x)|z)] + D_{KL}(q(z|x)||p(z)) \\ &= \mathcal{L}_{reconstruction} + \mathcal{L}_{prior} \\ &= MSE(x, \tilde{x}) + D_{KL}(q(z|x)||p(z))\end{aligned}$$

Mean-Squared Error

- Assumes signal fidelity is independent of temporal/spatial relationships → does not hold for images
- Element-wise metric unable to model human perception of image fidelity and quality → low image quality

VAE Limitations

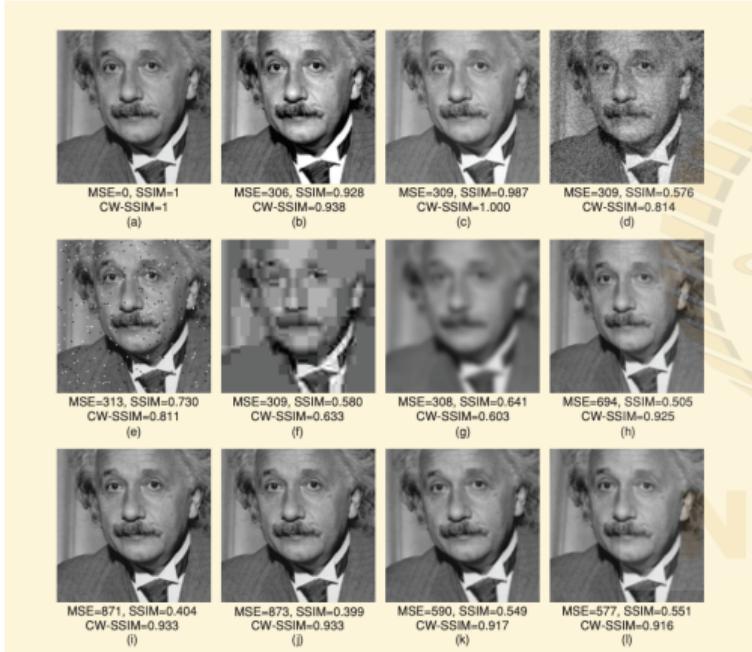
Let's recall the VAE objective:

$$\begin{aligned}\mathcal{L}_{VAE} &= -\mathbb{E}_{q(z|x)}[\log p((x)|z)] + D_{KL}(q(z|x)||p(z)) \\ &= \mathcal{L}_{reconstruction} + \mathcal{L}_{prior} \\ &= MSE(x, \tilde{x}) + D_{KL}(q(z|x)||p(z))\end{aligned}$$

Mean-Squared Error

- Assumes signal fidelity is independent of temporal/spatial relationships → does not hold for images
- Element-wise metric unable to model human perception of image fidelity and quality → low image quality
- Pixel-based loss metric does not respect semantic-preserving transforms, e.g. scaling/translation

MSE and Image Fidelity¹



[FIG2] Comparison of image fidelity measures for "Einstein" image altered with different types of distortions. (a) Reference image. (b) Mean contrast stretch. (c) Luminance shift. (d) Gaussian noise contamination. (e) Impulsive noise contamination. (f) JPEG compression. (g) Blurring. (h) Spatial scaling (zooming out). (i) Spatial shift (to the right). (j) Spatial shift (to the left). (k) Rotation (counter-clockwise). (l) Rotation (clockwise).

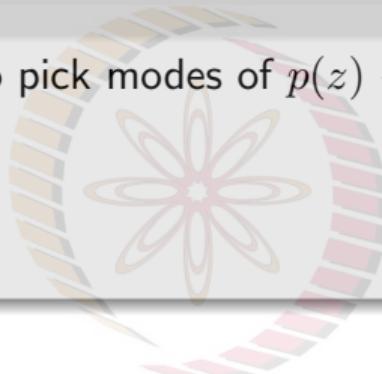
- MSE values of all distorted images - images (b)–(g) relative to (a) - are nearly identical, even though images are perceptually different
- Images with small geometrical modifications - images (h)–(i) - present large MSE values relative to (a), yet show negligible change in perceived quality

¹Wang and Bovik, Mean Squared Error: Love It or Leave It? IEEE Signal Processing Magazine, 2009

VAE Limitations

KL Divergence

- Focused on encouraging $q(z)$ to pick modes of $p(z) \rightarrow$ unable to match $q(z)$ to whole distribution of $p(z)$ well

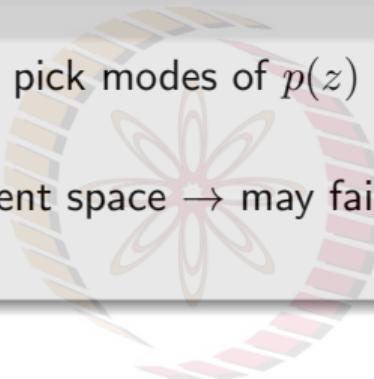


NPTEL

VAE Limitations

KL Divergence

- Focused on encouraging $q(z)$ to pick modes of $p(z) \rightarrow$ unable to match $q(z)$ to whole distribution of $p(z)$ well
- 'Spaces' or 'holes' in learned latent space \rightarrow may fail to capture data manifold

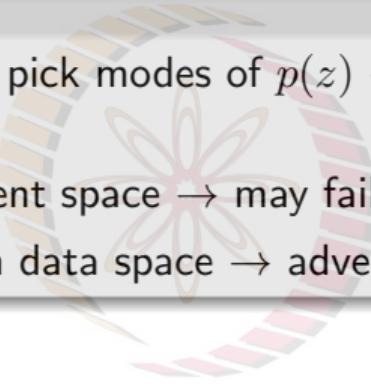


NPTEL

VAE Limitations

KL Divergence

- Focused on encouraging $q(z)$ to pick modes of $p(z) \rightarrow$ unable to match $q(z)$ to whole distribution of $p(z)$ well
- 'Spaces' or 'holes' in learned latent space \rightarrow may fail to capture data manifold
- May miss several local regions in data space \rightarrow adversely effect generalization capability



NPTEL

VAE Limitations

KL Divergence

- Focused on encouraging $q(z)$ to pick modes of $p(z) \rightarrow$ unable to match $q(z)$ to whole distribution of $p(z)$ well
- 'Spaces' or 'holes' in learned latent space \rightarrow may fail to capture data manifold
- May miss several local regions in data space \rightarrow adversely effect generalization capability

Form of Prior

- Requires access to exact functional form of prior

VAE Limitations

KL Divergence

- Focused on encouraging $q(z)$ to pick modes of $p(z) \rightarrow$ unable to match $q(z)$ to whole distribution of $p(z)$ well
- 'Spaces' or 'holes' in learned latent space \rightarrow may fail to capture data manifold
- May miss several local regions in data space \rightarrow adversely effect generalization capability

Form of Prior

- Requires access to exact functional form of prior
- Difficult to optimize; not computable in closed form always for all priors \rightarrow limits choice of priors that can be used

How to address?

VAE Limitations

KL Divergence

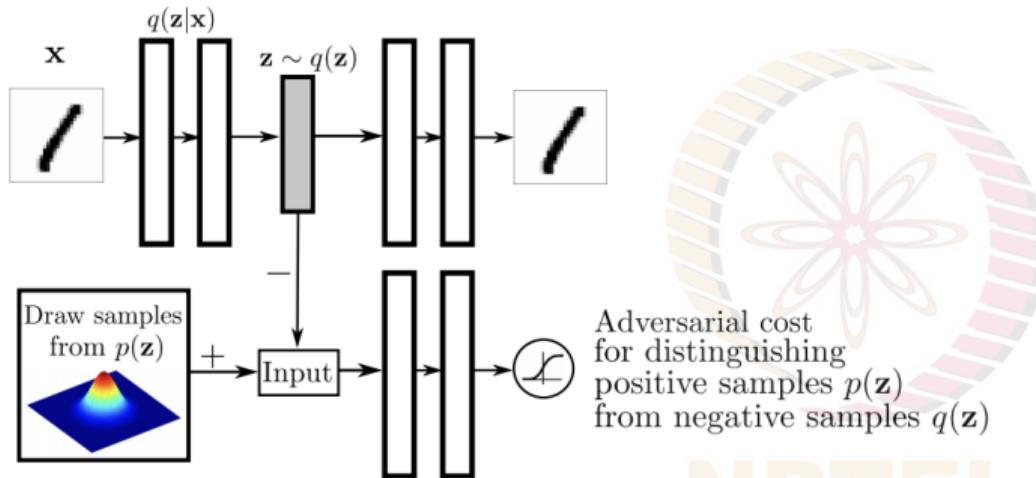
- Focused on encouraging $q(z)$ to pick modes of $p(z) \rightarrow$ unable to match $q(z)$ to whole distribution of $p(z)$ well
- 'Spaces' or 'holes' in learned latent space \rightarrow may fail to capture data manifold
- May miss several local regions in data space \rightarrow adversely effect generalization capability

Form of Prior

- Requires access to exact functional form of prior
- Difficult to optimize; not computable in closed form always for all priors \rightarrow limits choice of priors that can be used

How to address? Integrate with GANs to use its positives to overcome limitations!

Adversarial Autoencoder (AAE)²



(Top) Standard VAE

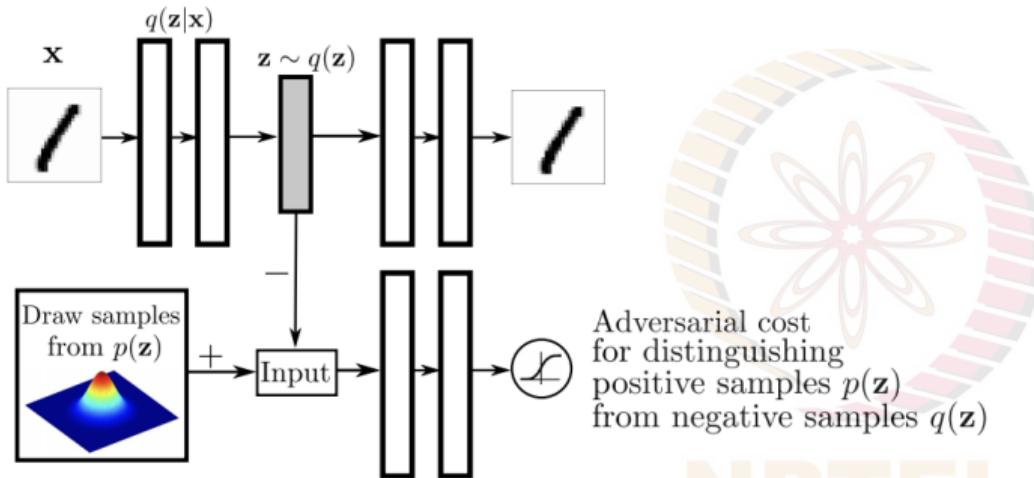
(Bottom) Second network trained to discriminatively predict whether a sample arises from hidden code of autoencoder or input training distribution

- Aims to match aggregated posterior, $q(z)$, to an arbitrary prior, $p(z)$ via adversarial objective-based training

NPTEL

²Makhzani et al, Adversarial Autoencoders, ICLRW 2016

Adversarial Autoencoder (AAE)²



(Top) Standard VAE

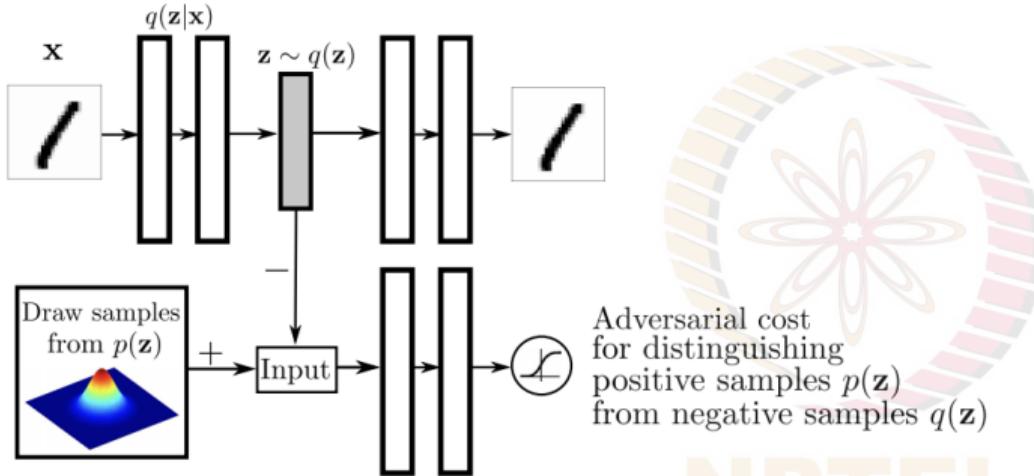
(Bottom) Second network trained to discriminatively predict whether a sample arises from hidden code of autoencoder or input training distribution

- Aims to match aggregated posterior, $q(z)$, to an arbitrary prior, $p(z)$ via adversarial objective-based training
- Renders continuous learned latent space → captures data manifold well

NPTEL

²Makhzani et al, Adversarial Autoencoders, ICLRW 2016

Adversarial Autoencoder (AAE)²



(Top) Standard VAE

(Bottom) Second network trained to discriminatively predict whether a sample arises from hidden code of autoencoder or input training distribution

- Aims to match aggregated posterior, $q(z)$, to an arbitrary prior, $p(z)$ via adversarial objective-based training
- Renders continuous learned latent space → captures data manifold well
- Encoder converts data distribution to prior distribution, while decoder learns a deep generative model that maps imposed prior to data distribution

²Makhzani et al, Adversarial Autoencoders, ICLRW 2016

Training AAE

Objective:

$$\mathcal{L} = \underbrace{\mathbb{E}_x [\mathbb{E}_{q(z|x)} [-\log p(x|z)]]}_{\text{Reconstruction Error}} + \underbrace{\mathbb{E}_x [\text{KL}(q(z|x)||p(z))]}_{\text{KL Regularizer}}$$

↓
Replaced by adversarial loss in AAE

NPTEL

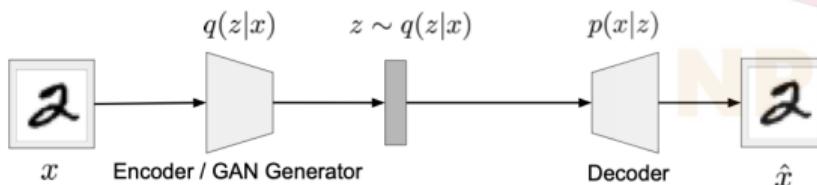
Training AAE

Objective:

$$\mathcal{L} = \underbrace{\mathbb{E}_x [\mathbb{E}_{q(z|x)} [-\log p(x|z)]]}_{\text{Reconstruction Error}} + \underbrace{\mathbb{E}_x [\text{KL}(q(z|x)||p(z))]}_{\text{KL Regularizer}}$$

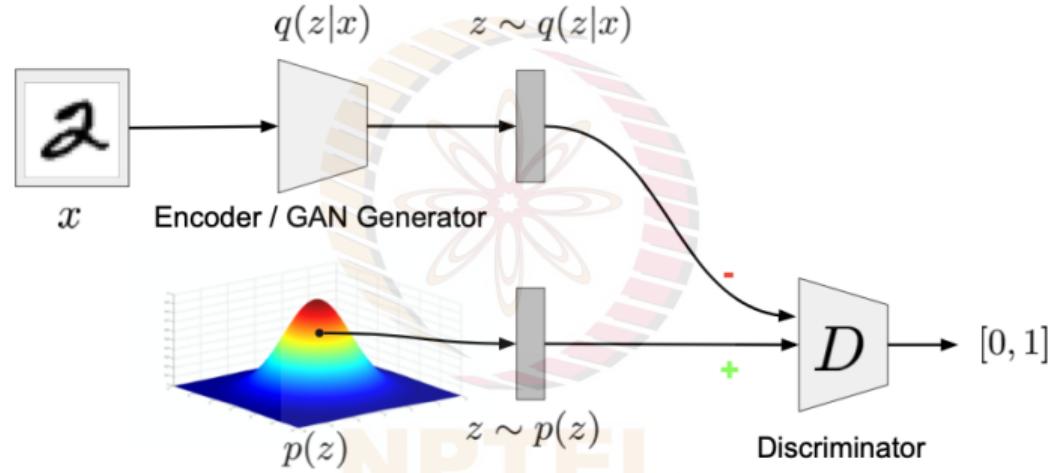
Replaced by adversarial loss in AAE

Reconstruction Phase:



- Introduce latent variable z with simple prior $p(z)$ (e.g. Gaussian)
- Sample $z \sim p(z)$, pass it through Generator $\hat{x} = G(z)$; where $\hat{x} \sim p_G$
- Introduce mechanism to ensure $p_G \approx p_{data}$

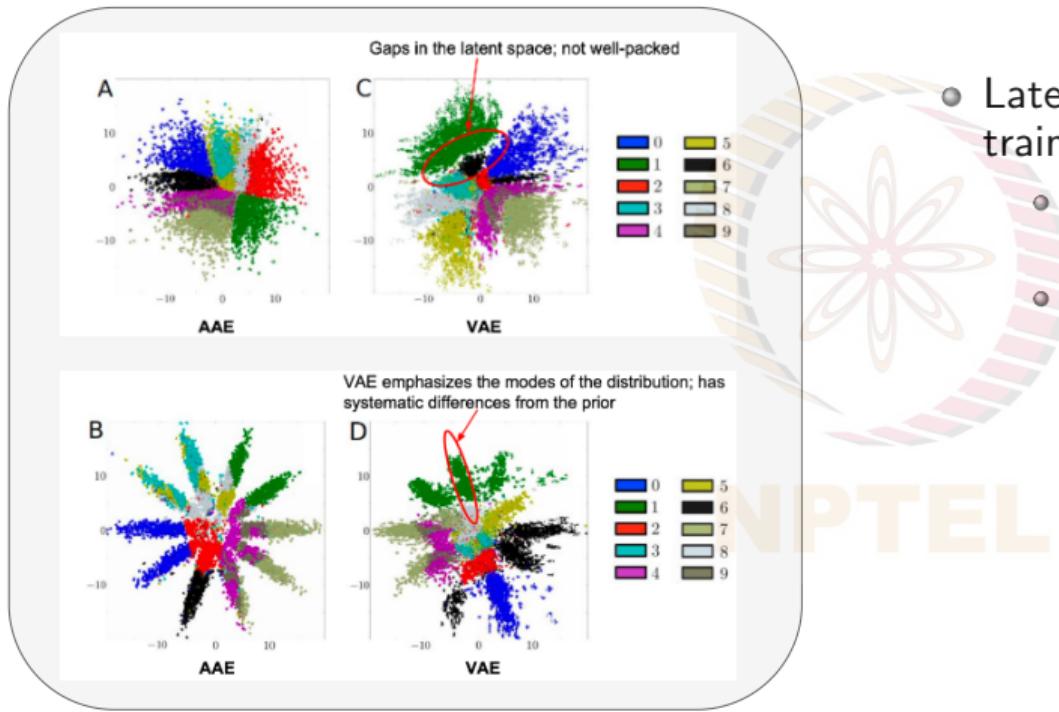
Training AAE



Regularization Phase

- Aims to match aggregated posterior, $q(z)$, to an arbitrary prior, $p(z)$ via adversarial objective-based training

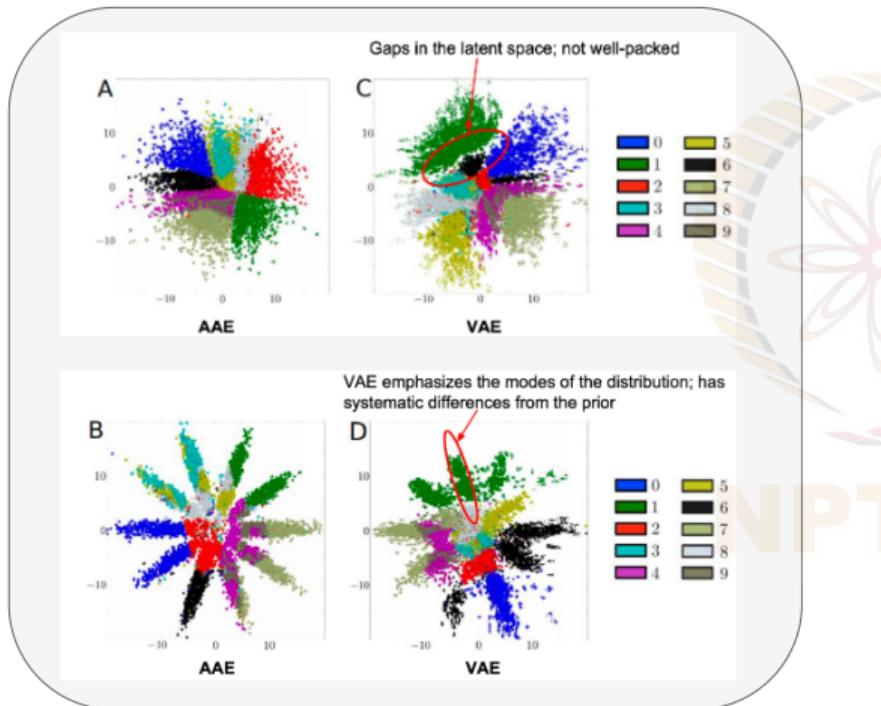
Learned Latent Space: AAE vs VAE³



- Latent space of test data using models trained on MNIST:
 - Top:** Spherical 2-D Gaussian prior distribution
 - Bottom:** Mixture of 10 2-D Gaussian

³Makhzani et al, Adversarial Autoencoders, ICLRW 2016

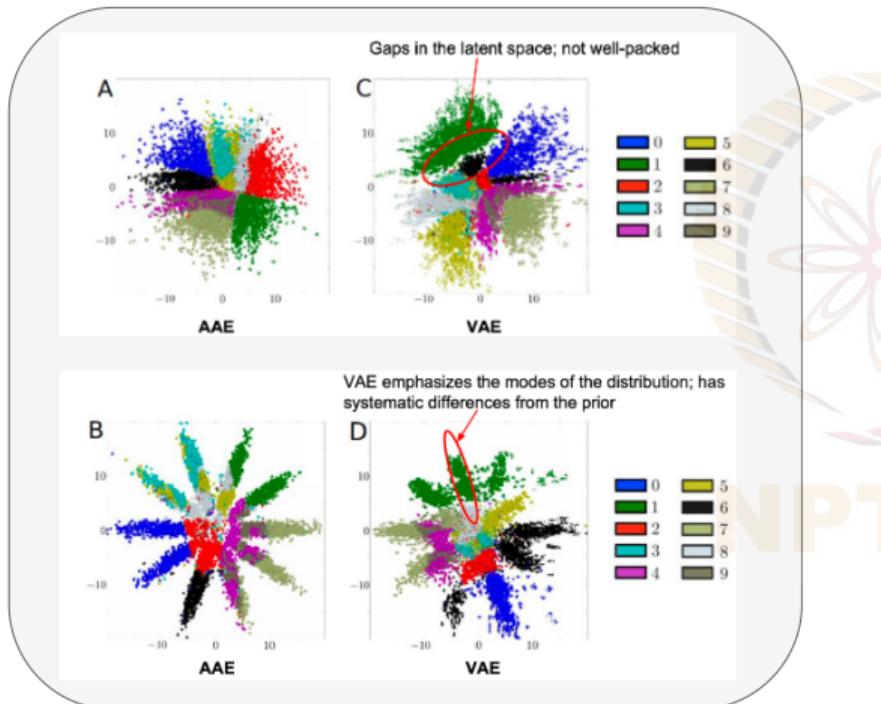
Learned Latent Space: AAE vs VAE³



- Latent space of test data using models trained on MNIST:
 - Top:** Spherical 2-D Gaussian prior distribution
 - Bottom:** Mixture of 10 2-D Gaussian
- Adversarial training to impose prior (AAE) renders a more continuous latent space relative to KL divergence-based (VAE) distribution alignment (**Top**)

³Makhzani et al, Adversarial Autoencoders, ICLRW 2016

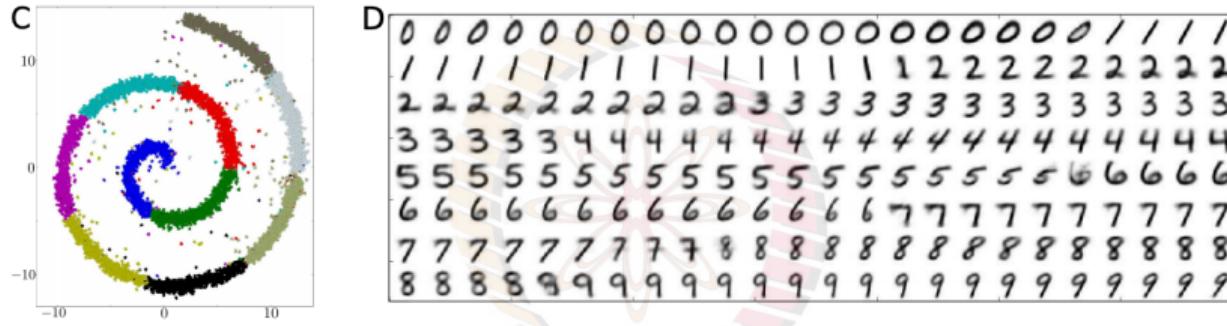
Learned Latent Space: AAE vs VAE³



- Latent space of test data using models trained on MNIST:
 - Top:** Spherical 2-D Gaussian prior distribution
 - Bottom:** Mixture of 10 2-D Gaussian
- Adversarial training to impose prior (AAE) renders a more continuous latent space relative to KL divergence-based (VAE) distribution alignment (**Top**)
- AAE imposes multi-modal distributions better than VAE (**Bottom**)

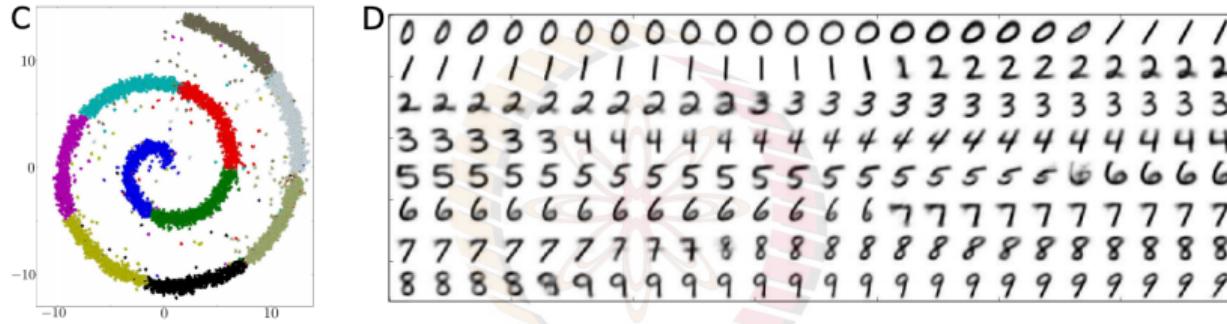
³Makhzani et al, Adversarial Autoencoders, ICLRW 2016

Imposing Complex Priors



- **Left:** Latent space of AAE trained on MNIST dataset with Swiss roll distribution as prior $p(z)$
- **Right:** Samples generated by walking along the main Swiss roll axis

Imposing Complex Priors



- **Left:** Latent space of AAE trained on MNIST dataset with Swiss roll distribution as prior $p(z)$
- **Right:** Samples generated by walking along the main Swiss roll axis
- For AAE, only require sampling from prior distribution in order to induce $q(z)$ to match $p(z)$; exact functional form of prior is not required

VAE-GAN⁴

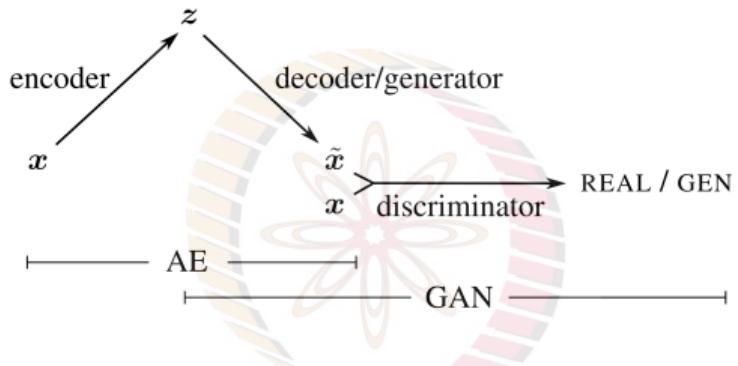


Figure 1. Overview of our network. We combine a VAE with a GAN by collapsing the decoder and the generator into one.

- Replace element-wise MSE in pixel space with feature-wise metric between discriminator's hidden representations

⁴Larsen et al, Autoencoding beyond Pixels using a Learned Similarity Metric, ICML 2016

VAE-GAN⁴

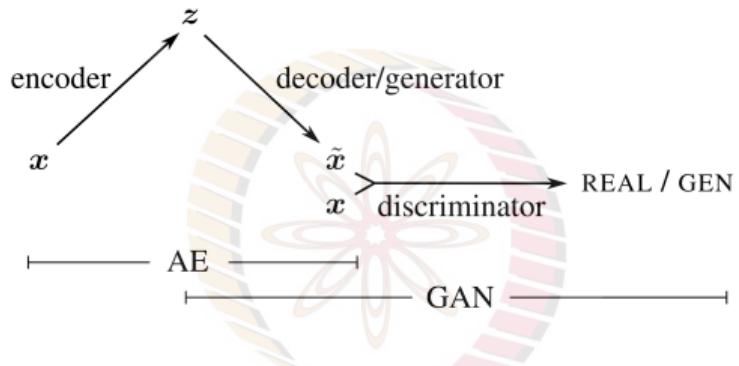


Figure 1. Overview of our network. We combine a VAE with a GAN by collapsing the decoder and the generator into one.

- Replace element-wise MSE in pixel space with feature-wise metric between discriminator's hidden representations
- Combines advantage of GAN as high-quality generative model and VAE as a method that produces an encoding of data into latent space z

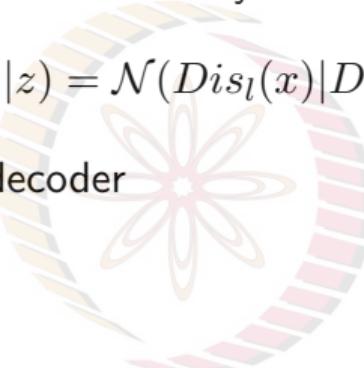
⁴Larsen et al, Autoencoding beyond Pixels using a Learned Similarity Metric, ICML 2016

Loss Formulation

Let $Dis_l(x)$ denote hidden representation of l^{th} layer of discriminator; then:

$$p(Dis_l(x)|z) = \mathcal{N}(Dis_l(x)|Dis_l(\tilde{x})|, I)$$

where $\tilde{x} \sim Dec(z)$ is a sample from decoder



Loss Formulation

Let $Dis_l(x)$ denote hidden representation of l^{th} layer of discriminator; then:

$$p(Dis_l(x)|z) = \mathcal{N}(Dis_l(x)|Dis_l(\tilde{x})|, I)$$

where $\tilde{x} \sim Dec(z)$ is a sample from decoder

$$\mathcal{L}_{recon-content}^{Dis_l} = -\mathbb{E}_{q(z|x)}[\log p(Dis_l(x)|z)]$$

NPTEL

Loss Formulation

Let $Dis_l(x)$ denote hidden representation of l^{th} layer of discriminator; then:

$$p(Dis_l(x)|z) = \mathcal{N}(Dis_l(x)|Dis_l(\tilde{x})|, I)$$

where $\tilde{x} \sim Dec(z)$ is a sample from decoder

$$\mathcal{L}_{recon-content}^{Dis_l} = -\mathbb{E}_{q(z|x)}[\log p(Dis_l(x)|z)]$$

$$\mathcal{L}_{recon-style}^{GAN} = \log(Dis(x)) + \log(1 - Dis(Dec(Dec(z))))$$

Loss Formulation

Let $Dis_l(x)$ denote hidden representation of l^{th} layer of discriminator; then:

$$p(Dis_l(x)|z) = \mathcal{N}(Dis_l(x)|Dis_l(\tilde{x})|, I)$$

where $\tilde{x} \sim Dec(z)$ is a sample from decoder

$$\mathcal{L}_{recon-content}^{Dis_l} = -\mathbb{E}_{q(z|x)}[\log p(Dis_l(x)|z)]$$

$$\mathcal{L}_{recon-style}^{GAN} = \log(Dis(x)) + \log(1 - Dis(Dec(Dec(z))))$$

$$\mathcal{L}_{prior} = D_{KL}(q(z|x)||p(z))$$

Loss Formulation

Let $Dis_l(x)$ denote hidden representation of l^{th} layer of discriminator; then:

$$p(Dis_l(x)|z) = \mathcal{N}(Dis_l(x)|Dis_l(\tilde{x})|, I)$$

where $\tilde{x} \sim Dec(z)$ is a sample from decoder

$$\mathcal{L}_{recon-content}^{Dis_l} = -\mathbb{E}_{q(z|x)}[\log p(Dis_l(x)|z)]$$

$$\mathcal{L}_{recon-style}^{GAN} = \log(Dis(x)) + \log(1 - Dis(Dec(Dec(z))))$$

$$\mathcal{L}_{prior} = D_{KL}(q(z|x)||p(z))$$

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{recon-content}^{Dis_l} + \mathcal{L}_{recon-style}^{GAN} + \mathcal{L}_{prior}$$

Training Algorithm

Algorithm 1 Training the VAE/GAN model

$\theta_{\text{Enc}}, \theta_{\text{Dec}}, \theta_{\text{Dis}} \leftarrow$ initialize network parameters

repeat

$\mathbf{X} \leftarrow$ random mini-batch from dataset

$\mathbf{Z} \leftarrow \text{Enc}(\mathbf{X})$

$\mathcal{L}_{\text{prior}} \leftarrow D_{\text{KL}}(q(\mathbf{Z}|\mathbf{X}) \| p(\mathbf{Z}))$

$\tilde{\mathbf{X}} \leftarrow \text{Dec}(\mathbf{Z})$

$\mathcal{L}_{\text{llike}}^{\text{Dis}_l} \leftarrow -\mathbb{E}_{q(\mathbf{Z}|\mathbf{X})} [p(\text{Dis}_l(\mathbf{X})|\mathbf{Z})]$

$\mathbf{Z}_p \leftarrow$ samples from prior $\mathcal{N}(\mathbf{0}, \mathbf{I})$

$\mathbf{X}_p \leftarrow \text{Dec}(\mathbf{Z}_p)$

$\mathcal{L}_{\text{GAN}} \leftarrow \log(\text{Dis}(\mathbf{X})) + \log(1 - \text{Dis}(\tilde{\mathbf{X}}))$
+ $\log(1 - \text{Dis}(\mathbf{X}_p))$

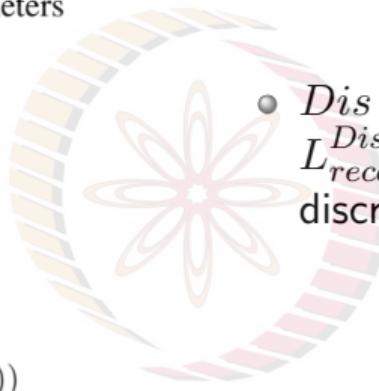
// Update parameters according to gradients

$\theta_{\text{Enc}} \xleftarrow{+} -\nabla_{\theta_{\text{Enc}}} (\mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{llike}}^{\text{Dis}_l})$

$\theta_{\text{Dec}} \xleftarrow{+} -\nabla_{\theta_{\text{Dec}}} (\gamma \mathcal{L}_{\text{llike}}^{\text{Dis}_l} - \mathcal{L}_{\text{GAN}})$

$\theta_{\text{Dis}} \xleftarrow{+} -\nabla_{\theta_{\text{Dis}}} \mathcal{L}_{\text{GAN}}$

until deadline



NPTEL

- *Dis* should not try to minimize $L_{\text{recon-content}}^{\text{Dis}_l}$ as this would collapse the discriminator to 0

Training Algorithm

Algorithm 1 Training the VAE/GAN model

$\theta_{\text{Enc}}, \theta_{\text{Dec}}, \theta_{\text{Dis}} \leftarrow$ initialize network parameters

repeat

$\mathbf{X} \leftarrow$ random mini-batch from dataset

$\mathbf{Z} \leftarrow \text{Enc}(\mathbf{X})$

$\mathcal{L}_{\text{prior}} \leftarrow D_{\text{KL}}(q(\mathbf{Z}|\mathbf{X}) \| p(\mathbf{Z}))$

$\tilde{\mathbf{X}} \leftarrow \text{Dec}(\mathbf{Z})$

$\mathcal{L}_{\text{llike}}^{\text{Dis}_l} \leftarrow -\mathbb{E}_{q(\mathbf{Z}|\mathbf{X})} [p(\text{Dis}_l(\mathbf{X})|\mathbf{Z})]$

$\mathbf{Z}_p \leftarrow$ samples from prior $\mathcal{N}(\mathbf{0}, \mathbf{I})$

$\mathbf{X}_p \leftarrow \text{Dec}(\mathbf{Z}_p)$

$\mathcal{L}_{\text{GAN}} \leftarrow \log(\text{Dis}(\mathbf{X})) + \log(1 - \text{Dis}(\tilde{\mathbf{X}}))$
+ $\log(1 - \text{Dis}(\mathbf{X}_p))$

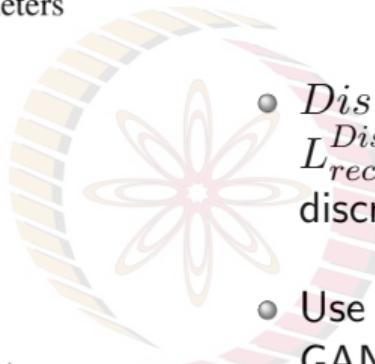
// Update parameters according to gradients

$\theta_{\text{Enc}} \xleftarrow{+} -\nabla_{\theta_{\text{Enc}}} (\mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{llike}}^{\text{Dis}_l})$

$\theta_{\text{Dec}} \xleftarrow{+} -\nabla_{\theta_{\text{Dec}}} (\gamma \mathcal{L}_{\text{llike}}^{\text{Dis}_l} - \mathcal{L}_{\text{GAN}})$

$\theta_{\text{Dis}} \xleftarrow{+} -\nabla_{\theta_{\text{Dis}}} \mathcal{L}_{\text{GAN}}$

until deadline



NPTEL

- Dis should not try to minimize $L_{\text{recon-content}}^{\text{Dis}_l}$ as this would collapse the discriminator to 0
- Use samples \mathbf{X}_p from prior directly in GAN loss in addition to $\tilde{\mathbf{X}}$

Training Algorithm

Algorithm 1 Training the VAE/GAN model

$\theta_{\text{Enc}}, \theta_{\text{Dec}}, \theta_{\text{Dis}} \leftarrow$ initialize network parameters

repeat

$\mathbf{X} \leftarrow$ random mini-batch from dataset

$\mathbf{Z} \leftarrow \text{Enc}(\mathbf{X})$

$\mathcal{L}_{\text{prior}} \leftarrow D_{\text{KL}}(q(\mathbf{Z}|\mathbf{X}) \| p(\mathbf{Z}))$

$\tilde{\mathbf{X}} \leftarrow \text{Dec}(\mathbf{Z})$

$\mathcal{L}_{\text{llike}}^{\text{Dis}_l} \leftarrow -\mathbb{E}_{q(\mathbf{Z}|\mathbf{X})} [p(\text{Dis}_l(\mathbf{X})|\mathbf{Z})]$

$\mathbf{Z}_p \leftarrow$ samples from prior $\mathcal{N}(\mathbf{0}, \mathbf{I})$

$\mathbf{X}_p \leftarrow \text{Dec}(\mathbf{Z}_p)$

$\mathcal{L}_{\text{GAN}} \leftarrow \log(\text{Dis}(\mathbf{X})) + \log(1 - \text{Dis}(\tilde{\mathbf{X}}))$
+ $\log(1 - \text{Dis}(\mathbf{X}_p))$

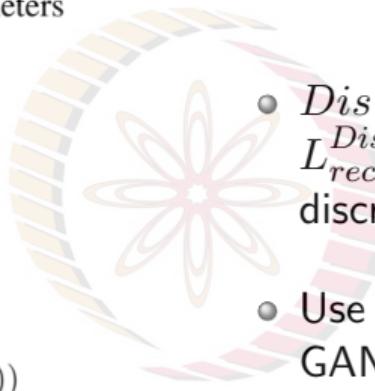
// Update parameters according to gradients

$\theta_{\text{Enc}} \xleftarrow{+} -\nabla_{\theta_{\text{Enc}}} (\mathcal{L}_{\text{prior}} + \mathcal{L}_{\text{llike}}^{\text{Dis}_l})$

$\theta_{\text{Dec}} \xleftarrow{+} -\nabla_{\theta_{\text{Dec}}} (\gamma \mathcal{L}_{\text{llike}}^{\text{Dis}_l} - \mathcal{L}_{\text{GAN}})$

$\theta_{\text{Dis}} \xleftarrow{+} -\nabla_{\theta_{\text{Dis}}} \mathcal{L}_{\text{GAN}}$

until deadline



NPTEL

- Dis should not try to minimize $L_{\text{recon-content}}^{\text{Dis}_l}$ as this would collapse the discriminator to 0
- Use samples \mathbf{X}_p from prior directly in GAN loss in addition to $\tilde{\mathbf{X}}$
- Weight $L_{\text{recon-content}}^{\text{Dis}_l}$ and $L_{\text{recon-style}}^{\text{GAN}}$ to weight ability to reconstruct vs fooling discriminator

Comparing Generated Samples

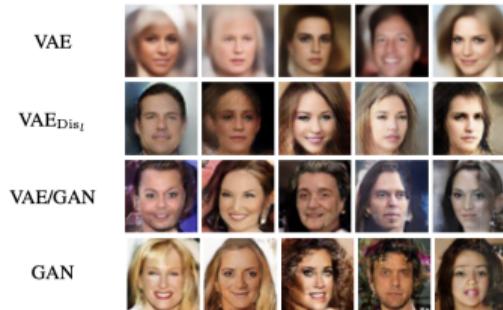


Figure 3. Samples from different generative models.



Figure 4. Reconstructions from different autoencoders.

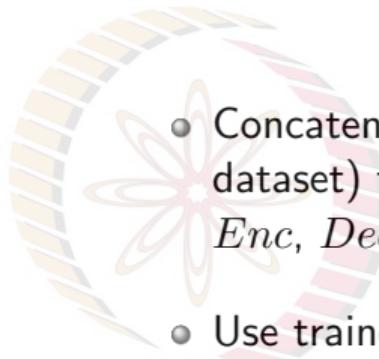
- Draw samples from $p(z)$ and propagate these through decoder to generate new images

- VAE draws frontal part well, but off-center gets blurry

- VAE_{Dis}^l produces sharper images even off-center

Conditional Generation

| | | |
|---------|---|--|
| Query |  | Prominent attributes: White, Fully Visible Forehead, Mouth Closed, Male, Curly Hair, Eyes Open, Pale Skin, Frowning, Pointy Nose, Teeth Not Visible, No Eyewear. |
| VAE |  | |
| GAN |  | |
| VAE/GAN |  | |
| Query |  | Prominent attributes: White, Male, Curly Hair, Frowning, Eyes Open, Pointy Nose, Flash, Posed Photo, Eyeglasses, Narrow Eyes, Teeth Not Visible, Senior, Receding Hairline. |
| VAE |  | |
| GAN |  | |
| VAE/GAN |  | |



- Concatenate face attribute vector (LFW image dataset) to vector representation of input in *Enc*, *Dec* and *Dis* while training
- Use trained model to generate faces conditioned on held-out test attributes
- Compared to an ordinary VAE, the VAE/GAN model yields significantly better images visually

Homework

Readings

- A wizard's guide to Adversarial Autoencoders: Part 1, Autoencoder
- A wizard's guide to Adversarial Autoencoders: Part 2, Exploring latent space with Adversarial Autoencoders
- What The Heck Are VAE-GANs?
- (YouTube video) VAE-GAN Explained

