

$$\textcircled{1} D_{KL}(q(x) || p(x)) = - \int q(x) \log \left[\frac{p(x)}{q(x)} \right] dx \geq 0$$

$$= - \int q(x) \log p(x) + \int q(x) \log q(x) \quad - \text{eqn } \textcircled{1}$$

$$q(x) = N(\mu, \sigma^2) \quad p(x) = N(0, 1)$$

$$D_{KL} = \int q(x) \log(p(x)) dx + \int q(x) \log(q(x)) dx - \textcircled{1}$$

2nd term in eqn \textcircled{1}

$$= \int q(x) \log(q(x)) dx$$

$$= \text{Expectation} [\log N(\mu, \sigma^2)]$$

$$\Rightarrow E \left\{ \log \left[(2\pi\sigma^2)^{-\frac{1}{2}} \exp \left(-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} \right) \right] \right\}.$$

$$= E \left[-\frac{1}{2} \log(2\pi\sigma^2) \right] + E \left[-\frac{1}{2\sigma^2} (x-\mu)^2 \right]$$

\downarrow
const.

$$= -\frac{1}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} [E(x-\mu)^2]$$

$$E(x-\mu)^2 = \sigma^2$$

$$\int q(x) \log q(x) = -\frac{1}{2} [1 + \log(2\pi\sigma^2)] \quad - \text{2}$$

1st term in equation ①.

$$\begin{aligned} & \int g(x) \log p(x) dx \quad \sigma_0=1, \mu_0=0 \\ = & \int \frac{1}{\sqrt{2\pi}\sigma} e^{-(\frac{x-\mu}{\sigma})^2} \log \left[\frac{1}{\sqrt{2\pi}\sigma_0} e^{-(\frac{x-\mu_0}{\sigma_0})^2} \right] dx \\ = & \frac{1}{\sqrt{2\pi}\sigma} \int e^{-(\frac{x-\mu}{\sigma})^2} \left[\log \left[\frac{1}{\sqrt{2\pi}\sigma_0^2} e^{-(\frac{x-\mu_0}{\sigma_0})^2} \right] \right] dx \\ = & \int_{-\infty}^{\infty} -\frac{1}{2} \log [2\pi\sigma_0^2] g(x) dx + \frac{g(x)}{2} \int \left[-\left(\frac{x-\mu_0}{\sigma_0} \right)^2 \right] dx \\ = & -\frac{1}{2} \log [2\pi\sigma_0^2] - \frac{g(x)_n}{\sigma_0^2} \int [x^2 - 2\mu_0 x + \mu_0^2] dx \\ = & -\frac{1}{2} \log [2\pi\sigma_0^2] - \frac{1}{2\sigma_0^2} \left[\int g(x)x^2 dx - 2\mu_0 \int g(x)x dx + \int g(x)\mu_0^2 dx \right] \\ & \boxed{\int g(x)x^2 dx = \frac{\sigma^2}{\mu^2} + \mu^2; \int g(x)x dx = \mu.} \\ = & -\frac{1}{2} \log [2\pi\sigma_0^2] - \frac{1}{2\sigma_0^2} \left[\sigma^2 - 2\mu_0 \mu + \mu^2 \right] \\ = & -\frac{1}{2} \left[\log(2\pi\sigma_0^2) + \frac{1}{\sigma_0^2} [\sigma^2 + (\mu - \mu_0)^2] \right] \\ & \sigma_0=1, \mu_0=0 \\ = & -\frac{1}{2} \left[\log(2\pi) + \frac{1}{\sigma^2} [\sigma^2 + \mu^2] \right] - ③ \end{aligned}$$

Putting ②③ back in ①
we get

$$-\frac{1}{2} \left[\log(2\pi) + (\sigma^2 + \mu^2) - 1 - \log(2000) \right]$$

$$= -\frac{1}{2} \left[\log \frac{2\pi\sigma^2}{2000} + 1 - \sigma^2 - \mu^2 \right]$$

$$= -\frac{1}{2} \left[\log(\sigma^2) + 1 - \sigma^2 - \mu^2 \right]$$

b) When α is too large, the KL divergence loss will give higher gradients to the network. So the KL-D Loss will try to force the α -neuron, to condition to encoder sampler. So all the data sampled will have similar latent vector. So all the reconstructed data will be similar to each other.

As the decoder after training is deterministic, if the sampled latent vectors are similar the reconstruction will yield similar data.

∴ All reconstructed data will be similar.

c)

① PCA, as a linear operation, that transforms a higher dimensional data to a low dimensional data. while VAE, does the same but with non-linearities between layers.

② PCA always reduce data. In VAEs, you can project data to even a higher dimension. This is more like a feature transformation.

③ In PCA there is no parameters learned, but VAE's can be trained to find the correct weights to encode the data.

Problem 2:-

(a) False.

We should update discriminator $k \geq 1$ times when we optimize generator once. If we do opposite as mentioned in the statement, generator becomes really better than discriminator, so every bad quality fake sample gets passed through the discriminator... So GAN's should always have a balance between discriminator and generator which will be broken when we train according to the statement.

(b) True

$D(G(z))$ initially will be close to 0, because the generator will be very bad initially and discriminator will be good, so all the fake samples generated from generator $G(z)$, will be close to 0, and will be classified as fake = 0, by the discriminator.

(c) Blue, non-saturating cost.

Initially as explained above the value of $D(G(z))$ will be close to 0, The figure shows that orange has a very low slope close to 0, initially while blue has a higher slope. So gradients will be higher for non-saturating cost.

function blue. As gradients are sufficient the generator's optimization will be faster. For the case of saturating orange cost, the optimizers will suffer from vanishing gradients.

(d) False

If $D(G(z)) \approx 1$, it means all the fake samples are getting classified as real by the discriminator. This can occur when discriminator is very strong or generator is very powerful. If discriminator is weak, then even the bad data samples generated will be passed through. As we want a balance between generator and discriminator. Ideally both of them should be around 0.5. Generator & discriminator should be both powerful.

Programming

Problem 1:

VAE Training Parameter Details

1. Encoder structure

Four Linear layers of the following input and output size

- nn.Linear(200, 128)
- nn.Linear(128, 64)
- nn.Linear(64, 32)
- nn.Linear(32, 16)

Activation layer: ReLU

Batchnorm1d between linear layers

number of trainable parameters: 24768

2. Decoder structure

Four Linear layers of the following input and output size

- nn.Linear(16, 32)
- nn.Linear(32, 64)
- nn.Linear(64, 128)
- nn.Linear(128, 200)

Activation layer: ReLU, at the end of decoder tanh is used

Batchnorm1d between linear layers

number of trainable parameters: 24324

3. Learning rate: 0.001

4. Batch size: 512

5. Number of epochs: 2000

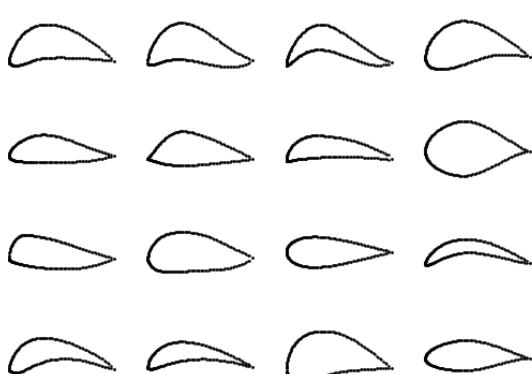
6. Alpha for kl divergence loss weightage = 0.005

7. Loss function = MSELoss() + alpha*KL_divergence, {reduction used is sum}

8. Optimizer = Adam()

9. Learning rate scheduler = StepLR(step_size=100, gamma=0.4)

Results



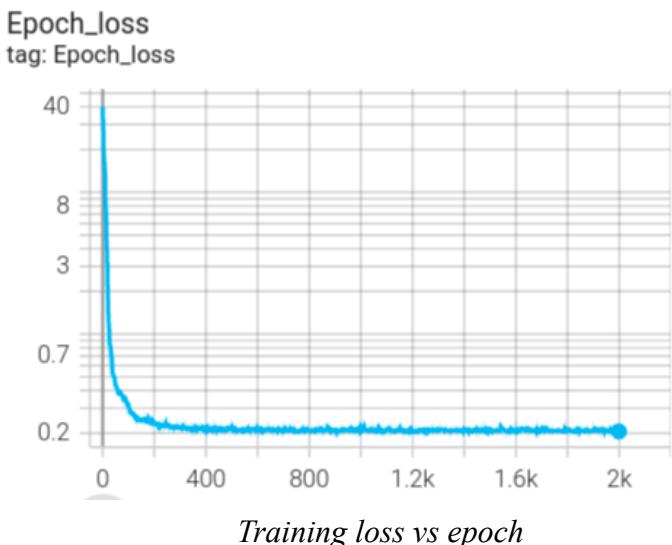
Real airfoils



Reconstructed airfoils



Synthesized airfoils



Training loss vs epoch

GAN

Problem 1:

GAN Training Parameter Details

10. Generator structure

Four Linear layers of the following input and output size

- Linear(in_features=16, out_features=32, bias=True)
- Linear(in_features=32, out_features=64, bias=True)
- Linear(in_features=64, out_features=128, bias=True)
- Linear(in_features=128, out_features=200, bias=True)

Activation layer: ReLU

number of trainable parameters: 36776

11. Discriminator structure

Four Linear layers of the following input and output size

- Linear(in_features=200, out_features=128, bias=True)
- Linear(in_features=128, out_features=64, bias=True)
- Linear(in_features=64, out_features=32, bias=True)
- Linear(in_features=32, out_features=1, bias=True)

Activation layer: ReLU,

number of trainable parameters: 36097

12. Gen Learning rate: 0.00015

13. Disc learning rate: 0.0002

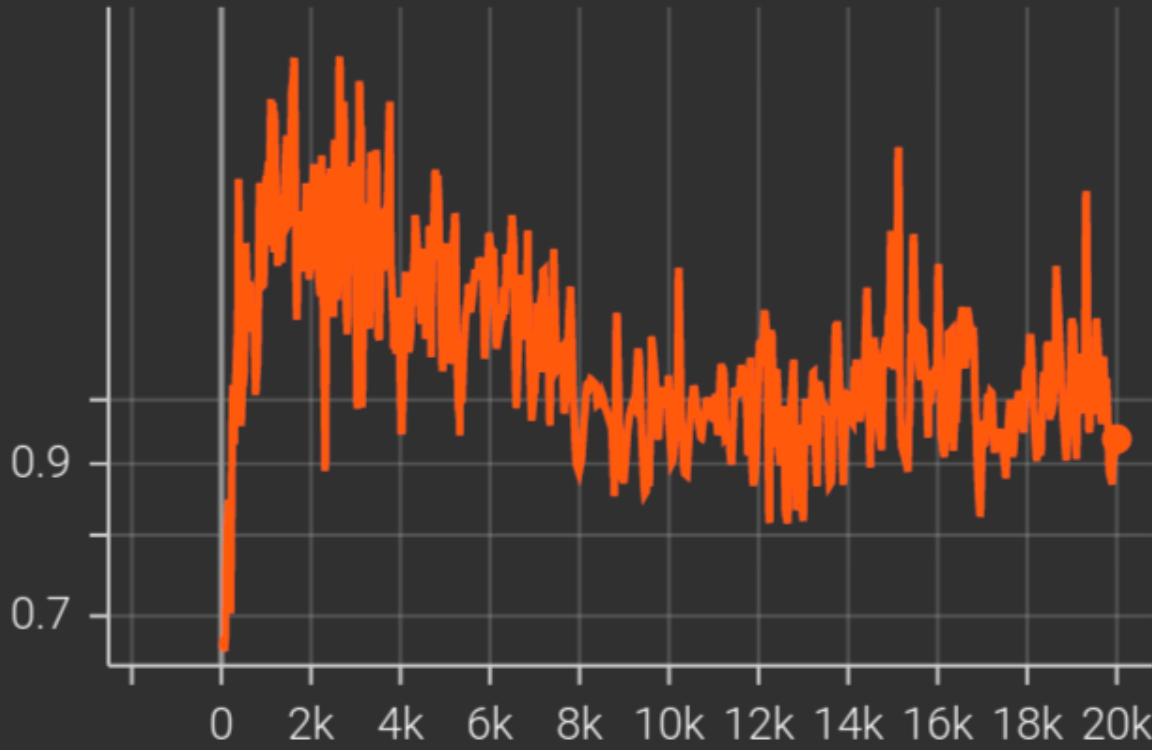
14. Batch size: 512

15. Number of epochs: 5000

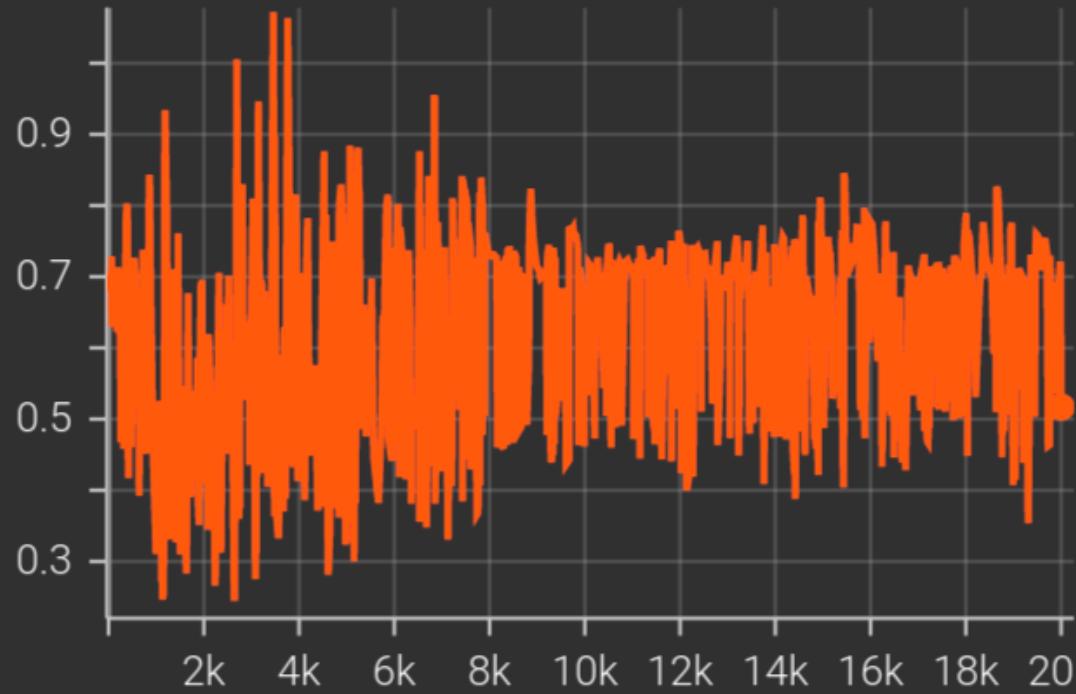
16. Loss function = BCEwithlogitsloss()

17. Optimizer = Adam() {for both gen and dis}

Gen loss
tag: Gen loss

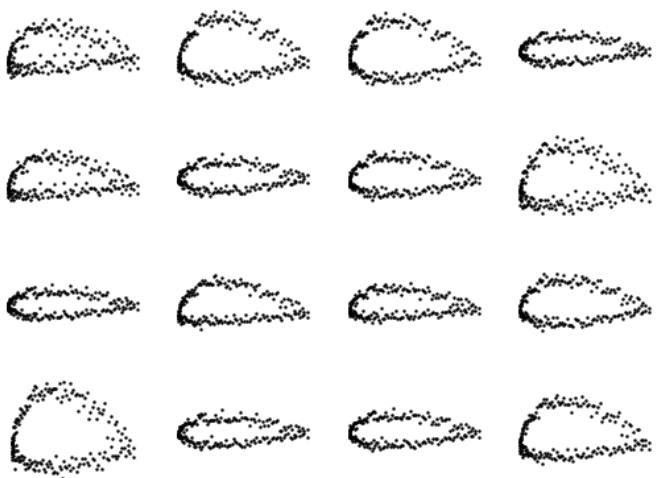


Disc loss
tag: Disc loss



Loss vs number of iterations

Generated data from noise



Observation: The samples from gan are very noisy. The reason might be unlike VAE we do not have a reconstruction component in the loss, which compares each point generated individually in my code MSELoss. Also, VAE directly optimizes the ELBO function of MLE, while GAN is just classifying generated and real data. For gan, all we have is a classification loss which suggests if the data is fake or real. That being said the data from gans appear more diverse.