

**Problem 1. VAE (6%)**

- 1) The architecture of our VAE network is shown in fig1\_1. batch size=32, optimizer uses Adam with learning rate  $lr=2e-4$ , betas=(0.5,0.999), KL lambda= $1e-5$ , Dim of z selected at 512.

```

VAE(
  (encoder): Sequential(
    (0): LeakyReLUConv2d(
      (model): Sequential(
        (0): Conv2d(3, 32, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
        (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True)
        (2): ReLU(inplace)
      )
    )
    (1): LeakyReLUConv2d(
      (model): Sequential(
        (0): Conv2d(32, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
        (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True)
        (2): ReLU(inplace)
      )
    )
    (2): LeakyReLUConv2d(
      (model): Sequential(
        (0): Conv2d(64, 128, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
        (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True)
        (2): ReLU(inplace)
      )
    )
    (3): LeakyReLUConv2d(
      (model): Sequential(
        (0): Conv2d(128, 256, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
        (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True)
        (2): ReLU(inplace)
      )
    )
    (4): Flatten()
  )
  (decoder): Sequential(
    (0): Linear(in_features=512, out_features=4096, bias=True)
    (1): unFlatten()
    (2): LeakyReLUConvTranspose2d(
      (model): Sequential(
        (0): ConvTranspose2d(256, 128, kernel_size=(6, 6), stride=(2, 2), bias=False)
        (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True)
        (2): LeakyReLU(0.01, inplace)
      )
    )
    (3): LeakyReLUConvTranspose2d(
      (model): Sequential(
        (0): ConvTranspose2d(128, 64, kernel_size=(6, 6), stride=(2, 2), bias=False)
        (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True)
        (2): LeakyReLU(0.01, inplace)
      )
    )
    (4): LeakyReLUConvTranspose2d(
      (model): Sequential(
        (0): ConvTranspose2d(64, 32, kernel_size=(6, 6), stride=(2, 2), bias=False)
        (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True)
        (2): LeakyReLU(0.01, inplace)
      )
    )
    (5): ConvTranspose2d(32, 3, kernel_size=(5, 5), stride=(1, 1))
    (6): Tanh()
  )
  (muParametrize): Linear(in_features=4096, out_features=512, bias=True)
  (logvarParametrize): Linear(in_features=4096, out_features=512, bias=True)
)

```

fig1\_1

- 2) Learning curve of training our model is shown in fig1\_2, the left curve is KLD against training steps and the right curve is the MSE against training steps.

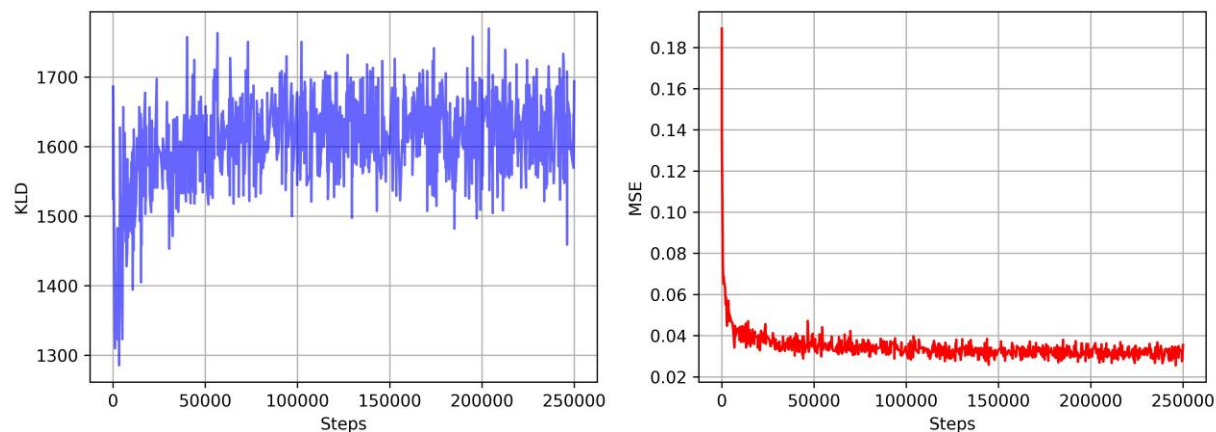


fig1\_2

- 3) The selected test images are shown in the first row of fig1\_3 and the reconstructed results are shown in the second row of fig1\_3. Besides, the MSE of entire test images is approximately **0.0291**.



fig1\_3

- 4) 32 random generated images of our model are plotted in fig1\_4.



fig1\_4

- 5) We visualize the latent space in fig1\_5 (map it from high dimensional to 2D using tSNE from sklearn), where the blue dots are faces have the *bangs* attribute, the red dots are faces don't have this attribute.

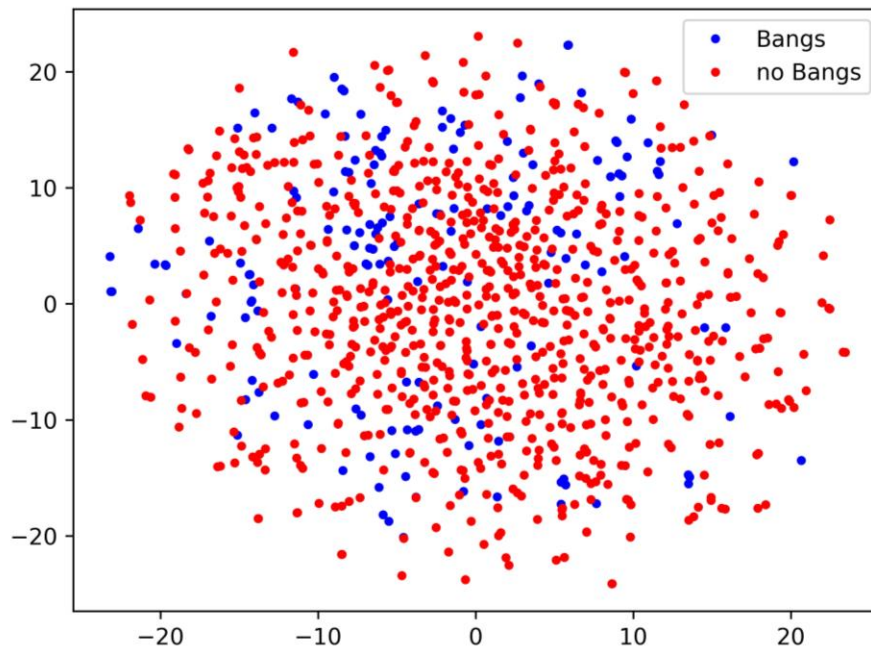


fig1\_5

6) What I learned from implement VAE:

- a. We learned how to construct the AE/VAE model and tune hyper-parameters. We start from Auto Encoder model, which is relatively easy to realize, then modified it to VAE.
- b. We understand the difference between VAE and AE. Using AE, we have no idea about how AE access the latent space. But using VAE, we can not only reconstruct image from image but also from noise, which is advanced than AE. It helps us to understand.
- c. At the begining, we didn't realize to normalize the image to -1~1, the results is bad. After fixed this bug, we realized the significance of data normalization.
- d. A bug that we use rand but not normal to generate noise, which can affect the results.
- e. At the begining, we set the batch size to 1000, with the hope that the training speed becomes quickly. However, the results becomes very bad.
- f. Besides, we also learned how to visulize the latent space using tSNE.
- g. KL lambda finds a trade-off between the reconstruction quality and the generated image quality.

## Problem 2. GAN (5%)

- 1) The architecture of our GAN network (DCGAN) is shown in fig2\_1, where the red box is the generator and the blue box is the discriminator. For generator, the optimizer uses Adam with learning rate G\_lr=2e-4, betas=(0.5,0.999). For discriminator, the optimizer uses Adam with learning rate D\_lr=5e-5, betas=(0.5,0.999). batch\_size=128 and z\_dim=100. The criterion we empolyed BCE loss:

$$loss(o, t) = -1/n \sum_i (t[i] * \log(o[i]) + (1 - t[i]) * \log(1 - o[i]))$$

In the process of training, we train the discriminator first, then the generator.

```

generator(
  (gen): Sequential(
    (0): ConvTranspose2d(100, 512, kernel_size=(4, 4), stride=(1, 1), bias=False)
    (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True)
    (2): LeakyReLU(0.1, inplace)
    (3): ConvTranspose2d(512, 256, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (4): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True)
    (5): LeakyReLU(0.1, inplace)
    (6): ConvTranspose2d(256, 128, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (7): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True)
    (8): LeakyReLU(0.1, inplace)
    (9): ConvTranspose2d(128, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (10): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True)
    (11): LeakyReLU(0.1, inplace)
    (12): ConvTranspose2d(64, 3, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (13): Tanh()
  )
)

discriminator(
  (dis): Sequential(
    (0): Conv2d(3, 64, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (1): LeakyReLU(0.2, inplace)
    (2): Conv2d(64, 128, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (3): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True)
    (4): LeakyReLU(0.2, inplace)
    (5): Conv2d(128, 256, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (6): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True)
    (7): LeakyReLU(0.2, inplace)
    (8): Conv2d(256, 512, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (9): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True)
    (10): LeakyReLU(0.2, inplace)
    (11): Conv2d(512, 1, kernel_size=(4, 4), stride=(1, 1), bias=False)
    (12): Sigmoid()
  )
)

```

fig2\_1

- 2) Learning curve is shown in fig2\_2. (left) shows the discriminator loss against steps. We can observe that the discriminator loss converges to approximately 0.5, which indicates that the discriminator cannot distinguish the fake images and the real images. Besides, the generator loss is increasing and converges after about 70k steps. (right) shows the discriminator accuracy against steps. We can observe that all the two accuracy curves converge after about 20k. The blue and the red curve results from feeding the real and generated images into the discriminator, respectively. We can found that the discriminator is trained well.

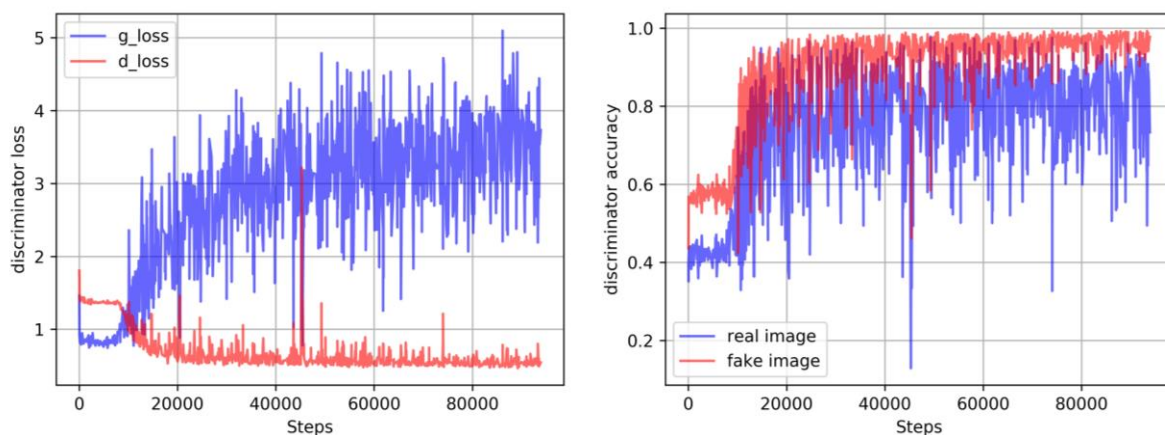


fig2\_2

- 3) 32 random generated images of our model is plotted in fig2\_3.





fig2\_3

- 4) After this homework, I agree that GAN is quite hard to train. We tried a lot of tricks to balance GAN. a) add noise to the output/real image. b) add noise to the label (i.e., soft label). c) flip the image, but the results are not so good. d) Train D more. e) Maximum  $D(G(z))$  instead of minimum  $1-D(G(z))$  to avoid vanish gradient issue. We also learned how to make a GAN experiment reproducible, To avoid the fast convergence of discriminator network, generator network is updated twice for each D network update, which differs from original paper. Discriminator Loss 0 is a failure mode which indicate the discriminator is too strong, therefore we need consider how to balance the training process. With the increase of Epoch, the training results could be much more worse. As shown in fig2\_4, the images boxed up using yellow box are the bad results.
- 5) The differences between images generated by VAE and GAN are that: a) VAE generated images are much more smooth than GAN, because VAE using MSE as the reconstruction loss which will result in the loss of high frequency signals; b) Compare with VAE, GAN generated images are much more sharp (contain a lot of high frequency signals) but they have more artifacts (as shown in fig2\_4), which we don't want.

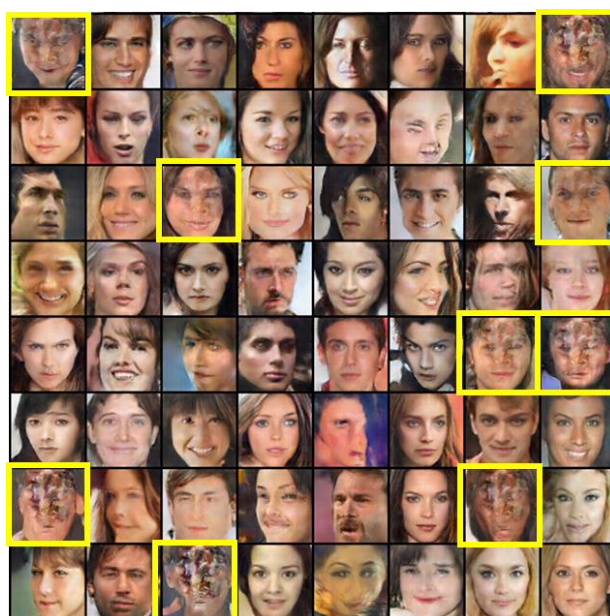


fig2\_4

**Problem 3. ACGAN (4%)**

- 1) The architecture of our VAE network is shown in fig3\_1, where the red box is the generator and the blue box is the discriminator. For generator, the optimizer uses Adam with learning rate  $G\_lr=2e-4$ ,  $\text{betas}=(0.5,0.999)$ . For discriminator, the optimizer uses Adam with learning rate  $D\_lr=2e-4$ ,  $\text{betas}=(0.5,0.999)$ .  $\text{batch\_size}=128$  and  $z\_dim=100$ . We use BCE loss as in Problem 2 for discriminator and CrossEntropyLoss for classification.

```

netG_ACGAN(
  (tconv1): Sequential(
    (0): ConvTranspose2d(101, 384, kernel_size=(4, 4), stride=(1, 1), bias=False)
    (1): BatchNorm2d(384, eps=1e-05, momentum=0.1, affine=True)
    (2): ReLU(inplace)
  )
  (tconv2): Sequential(
    (0): ConvTranspose2d(384, 192, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (1): BatchNorm2d(192, eps=1e-05, momentum=0.1, affine=True)
    (2): ReLU(inplace)
  )
  (tconv3): Sequential(
    (0): ConvTranspose2d(192, 96, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (1): BatchNorm2d(96, eps=1e-05, momentum=0.1, affine=True)
    (2): ReLU(inplace)
  )
  (tconv4): Sequential(
    (0): ConvTranspose2d(96, 48, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (1): BatchNorm2d(48, eps=1e-05, momentum=0.1, affine=True)
    (2): ReLU(inplace)
  )
  (tconv5): Sequential(
    (0): ConvTranspose2d(48, 3, kernel_size=(4, 4), stride=(2, 2), padding=(1, 1), bias=False)
    (1): Tanh()
  )
)

```

```

netD_ACGAN(
  (conv1): Sequential(
    (0): Conv2d(3, 16, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (1): LeakyReLU(0.2, inplace)
    (2): Dropout(p=0.5)
  )
  (conv2): Sequential(
    (0): Conv2d(16, 32, kernel_size=(3, 3), stride=(1, 1), padding=(1, 1), bias=False)
    (1): BatchNorm2d(32, eps=1e-05, momentum=0.1, affine=True)
    (2): LeakyReLU(0.2, inplace)
    (3): Dropout(p=0.5)
  )
  (conv3): Sequential(
    (0): Conv2d(32, 64, kernel_size=(3, 3), stride=(2, 2), padding=(1, 1), bias=False)
    (1): BatchNorm2d(64, eps=1e-05, momentum=0.1, affine=True)
    (2): LeakyReLU(0.2, inplace)
    (3): Dropout(p=0.5)
  )
  (conv4): Sequential(
    (0): Conv2d(64, 128, kernel_size=(3, 3), stride=(1, 1), bias=False)
    (1): BatchNorm2d(128, eps=1e-05, momentum=0.1, affine=True)
    (2): LeakyReLU(0.2, inplace)
    (3): Dropout(p=0.5)
  )
  (conv5): Sequential(
    (0): Conv2d(128, 256, kernel_size=(3, 3), stride=(2, 2), bias=False)
    (1): BatchNorm2d(256, eps=1e-05, momentum=0.1, affine=True)
    (2): LeakyReLU(0.2, inplace)
    (3): Dropout(p=0.5)
  )
  (conv6): Sequential(
    (0): Conv2d(256, 512, kernel_size=(3, 3), stride=(1, 1), bias=False)
    (1): BatchNorm2d(512, eps=1e-05, momentum=0.1, affine=True)
    (2): LeakyReLU(0.2, inplace)
    (3): Dropout(p=0.5)
  )
  (fc_dis): Linear(in_features=8192, out_features=1, bias=True)
  (fc_aux): Linear(in_features=8192, out_features=2, bias=True)
  (softmax): Softmax()
  (sigmoid): Sigmoid()
)

```

fig3\_1

- 2) Learning curve is shown in fig3\_2. (left) shows the discriminator loss against steps. We can observe that the generator loss and the discriminator reach to balance after about 10k steps. (right) shows the classify accuracy against steps. We can observe that the classifier performs better on fake images than real images.

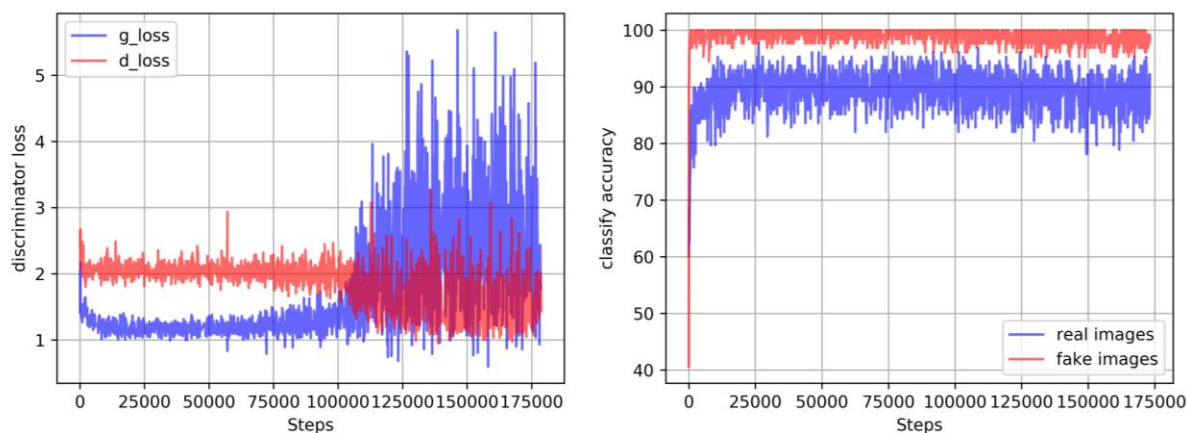


fig3\_2

- 3) 10 pairs of random generated images are shown in fig3\_3. The first row has the attribute “smile” and the second row does not have this attribute.



fig3\_3