

# Conditional Variational Autoencoder vs Conditional Generative Adversarial Network

## Introduction

### CVAE

Conditional Variational Autoencoder (CVAE) is an extension of Variational Autoencoder (VAE).

On VAE, the objective is:

$$\log P(X) - D_{KL}[Q(z|X) \| P(z|X)] = E[\log P(X|z)] - D_{KL}[Q(z|X) \| P(z)]$$

The original VAE model has two parts: the encoder  $Q(z|X)$  and the decoder  $P(X|z)$ . The reason why VAE can't generate specific data is because the encoder models the latent variable  $z$  directly based on  $X$ , it doesn't care about the different type of  $X$ .

Hence, VAE could be improved by conditioning the encoder and decoder to  $C$ , so the encoder is now conditioned to two variables  $X$  and  $C$ :  $Q(z|X, C)$ . The same with the decoder, it's now conditioned to two variables  $z$  and  $C$ :  $P(X|z, C)$ .

Finally, the variational lower bound objective is in this following form:

$$\log P(X | c) - D_{KL}[Q(z | X, c) || P(z | X, c)] = E[\log P(X | z, c)] - D_{KL}[Q(z | X, c) || P(z | c)]$$

Now, the real latent variable is distributed under  $P(z | c)$  and for each possible value of  $C$ , we would have a  $P(z)$ . We could also use this form of thinking for the decoder.

## CGAN

Similarly, Conditional Generative Adversarial Network (CGAN) is an extension of Generative Adversarial Network (GAN)

GAN consists of two ‘adversarial’ models: a generative model  $G$  that captures the data distribution, and a discriminative model  $D$  that estimates the probability that a sample came from the training data rather than  $G$ . Both  $G$  and  $D$  could be a non-linear mapping function, such as a multi-layer perceptron.

To learn a generator distribution  $p_g$  over data  $x$ , the generator builds a mapping function from a prior noise distribution  $p_z(z)$  to data space as  $G(z; \theta_g)$ . And the discriminator,  $D(x; \theta_d)$ , outputs a single scalar representing the probability that  $x$  came from training

data rather than  $p_g$ .  $G$  and  $D$  are both trained simultaneously: we adjust parameters for  $G$  to minimize  $\log(1 - D(G(z)))$  and adjust parameters for  $D$  to minimize  $\log D(X)$ , as if they are following the two-player min-max game with value function  $V(G, D)$ :

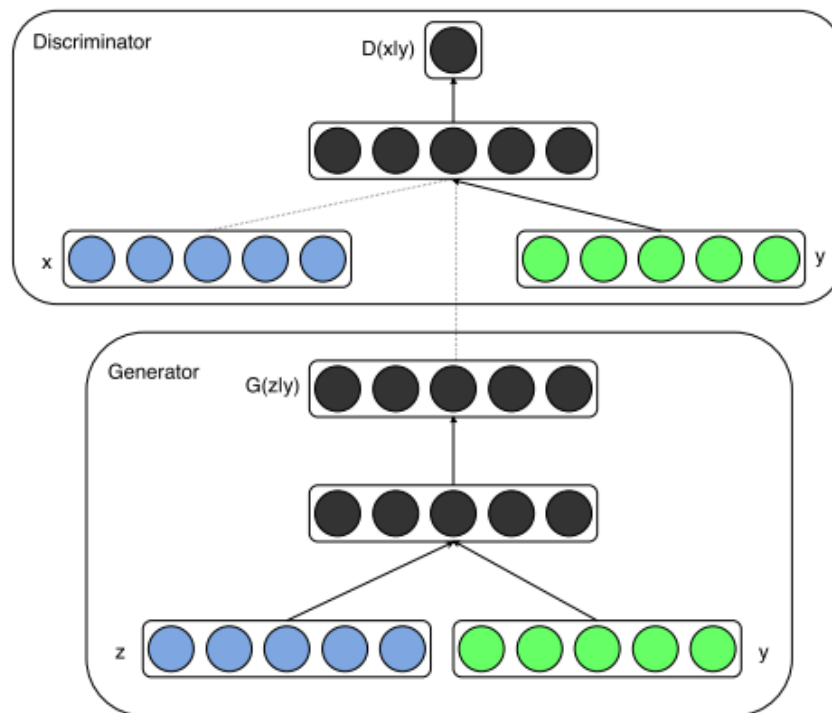
$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)} [\log D(x)] + E_{z \sim p_z(z)} [\log(1 - D(G(z)))].$$

Generative adversarial nets can be extended to a conditional model if both the generator and discriminator are conditioned on some extra information  $y$ .  $y$  could be any kind of auxiliary information, such as class labels or data from other modalities.

We can perform the conditioning by feeding  $y$  into the both the discriminator and generator as additional input layer. In the generator the prior input noise  $p_z(z)$ , and  $y$  are combined in joint hidden representation, and the adversarial training framework allows for considerable flexibility in how this hidden representation is composed. In the discriminator  $x$  and  $y$  are presented as inputs and to a discriminative function (embodied again by a MLP in this case).

The objective function of a two-player minimax game would be as

$$\min_G \max_D V(D, G) = E_{x \sim p_{data}(x)} [\log D(x|y)] + E_{z \sim p_z(z)} [\log(1 - D(G(z|y)))].$$



Conditional adversarial net

## Implementation

### CVAE

Use MNIST for example. We could use the label as our conditional variable  $C$ . In this case,  $C$  is categorically distributed.

```
1. mnist = input_data.read_data_sets('MNIST_data', one_hot=True)
2. X_train, y_train = mnist.train.images, mnist.train.labels
3. X_test, y_test = mnist.test.images, mnist.test.labels
4.
5. m = 50
6. n_x = X_train.shape[1]
7. n_y = y_train.shape[1]
8. n_z = 2
9. n_epoch = 20
```

```

10. # Q(z|X,y) -- encoder
11. X = Input(batch_shape=(m, n_x))
12. cond = Input(batch_shape=(m, n_y))

```

## Concatenation

```

1. inputs = merge([X, cond], mode='concat', concat_axis=1)
2.
3. h_q = Dense(512, activation='relu')(inputs)
4. mu = Dense(n_z, activation='linear')(h_q)
5. log_sigma = Dense(n_z, activation='linear')(h_q)
6.
7. def sample_z(args):
8.     mu, log_sigma = args
9.     eps = K.random_normal(shape=(m, n_z), mean=0., std=1.)
10.    return mu + K.exp(log_sigma / 2) * eps
11.
12. # Sample z ~ Q(z|X,y)
13. z = Lambda(sample_z)([mu, log_sigma])
14. z_cond = merge([z, cond], mode='concat', concat_axis=1) # <--- NEW!
15.
16. # P(X|z,y) -- decoder
17. decoder_hidden = Dense(512, activation='relu')
18. decoder_out = Dense(784, activation='sigmoid')
19.
20. h_p = decoder_hidden(z_cond)
21. outputs = decoder_out(h_p)

```

## Loss function

```

1. def vae_loss(y_true, y_pred):
2.     """ Calculate loss = reconstruction loss + KL loss for each data in mini
    batch """
3.     # E[log P(X|z,y)]
4.     recon = K.sum(K.binary_crossentropy(y_pred, y_true), axis=1)
5.     # D_KL(Q(z|X,y) || P(z|X)); calculate in closed form as both dist. are G
    aussian
6.     kl = 0.5 * K.sum(K.exp(log_sigma) + K.square(mu) - 1. - log_sigma, axis=
    1)
7.
8.     return recon + kl

```

# CGAN

Use MNIST for example as well.

```
1. mnist = input_data.read_data_sets('MNIST_data', one_hot=True)
2.
3. if not os.path.exists("logdir"):
4.     os.makedirs("logdir")
5.
6. LOGDIR = "logdir"
7. real_img_size = mnist.train.images[0].shape[0]
8. noise_size = 100
9. noise = 'normal0-1'
10. alpha = 0.1
11. learning_rate = 0.001
12. smooth = 0.05
13.
14. batch_size = 100
15. k = 10
16. epochs = 120
```

## Generator

```
1. def get_generator(digit, noise_img, reuse = False):
2.     with tf.variable_scope("generator", reuse = reuse):
3.         concatenated_img_digit = tf.concat([digit, noise_img], 1)
4.
5.         # output = tf.layers.dense(concatenated_img_digit, 256)
6.
7.         output = fully_connected('gf1', concatenated_img_digit, 128)
8.         output = leakyRelu(output)
9.         output = tf.layers.dropout(output, rate = 0.5)
10.
11.         # output = tf.layers.dense(output, 128)
12.
13.         output = fully_connected('gf2', output, 128)
14.         output = leakyRelu(output)
15.         output = tf.layers.dropout(output, rate = 0.5)
16.
17.         logits = tf.layers.dense(output, 784)
18.         logits = fully_connected('gf3', output, 784)
19.         outputs = tf.tanh(logits)
20.         return logits, outputs
```

## Discriminator

```
1. def get_discriminator(digit, img, reuse = False):
2.     with tf.variable_scope("discriminator", reuse=reuse):
3.         concatenated_img_digit = tf.concat([digit, img], 1)
4.
5.         # output = tf.layers.dense(concatenated_img_digit, 256)
6.         output = fully_connected('df1', concatenated_img_digit, 128)
7.         output = leakyRelu(output)
8.         output = tf.layers.dropout(output, rate = 0.5)
9.
10.        # output = tf.layers.dense(concatenated_img_digit, 128)
11.        output = fully_connected('df2', output, 128)
12.        output = leakyRelu(output)
13.        output = tf.layers.dropout(output, rate = 0.5)
14.
15.        logits = tf.layers.dense(output, 1)
16.        logits = fully_connected('df3', output, 1)
17.        outputs = tf.sigmoid(logits)
18.
19.        return logits, outputs
```

## Conditional MNIST

Test our CVAE model and CGAN model to generate MNIST data.

After reconstructing the images, the results are as followed.

Name	Epoch 1	Epoch 10	Epoch 25
CVAE			
			
			
			
			
			
			
			
			
			
CGAN			
			
			
			
			
			
			
			
			
			

We can find that, in general, for CVAE, the generated image is relatively stable, but the generated image diversity is bad, but for CGAN, the generated image diversity is better, but the generation process is random .

## Conclusion

In this post, we compared the CVAE to the CGAN. In CVAE, we could generate data with specific attribute, an operation that can't be done with the vanilla VAE. We showed this by applying CVAE on MNIST data and conditioned the model to the images' labels. The resulting model allows us to sample data under specific label. While in CGAN, compared with the original GAN network, the CGAN network has not changed, only the input data of the generator G and the discriminator D are changed, which allows CGAN to be embedded in other GAN networks as a general strategy.

All in all, we can know that the CVAE optimization process forcibly fits the data to a finite-dimensional Gaussian mixture or other distributions, resulting in: information loss during the mapping process, especially the loss of secondary information, and information that does not conform to the preset distribution. The encoding and recovery effect of the system is poor. In contrast, there is no such strong preset in the training of CGAN. It is optimized through the



discriminator network, so that the distribution of the data generated by the generator directly fits the distribution of the training data, and there is nothing special about the specific distribution.