311553060

周睦鈞

1. Introduction

Dec

這次作業要實作 Condition Variational Auto Encoder,我們想要透過前兩個 frame 的資訊來生成後面 10 個 frame 的結果,整 體的架構圖如下圖,會由一個 Encoder 和 Decoder 組成,其中在 latent space 中會加入 Condition 的資訊 (Action and Position) 在 Training 的過程中我們會透過調整 Teacher forcing ratio 以及使用不同的 KL weight type 來實驗看看各自的效果對於訓練結

果的差異

```
g_t
                                                              Prior loss
                       LSTM_{\theta}
                                                                 \mathbf{D}_{\mathrm{KL}}
                                                                                                       Dec
                                                N(\mu_{\phi(t)},\sigma_{\phi(t)})
                                                                         N(0, 1)
                        h_{t-1}
          Action
            and
                                                                          Fixed
          Position
                                                                                                        g_t
                         Enc
                                                                          prior
                                                                                                     LSTM_{\theta}
                        x_{t-1}
                                                     LSTM_{\phi}
                                                                                                                    \overline{z_t} \sim N(0,1)
                                                                                         Action
                      Reconstruction
                                                                                          and
                             loss
                                                                                        Position
                                                       Enc
                                                                                                       Enc
                                                       x_t
                             L2
                                                                                                      x_{t-1}
                                                                                                            (b)
                                           (a)
      Figure 1: The illustration of overall framework. (a) Training procedure (b) Generation procedure
2. Derivation of CVAE
```

 $P(x|c) = \int P(z|c) P(x|z,c) dz$

$\Rightarrow \int_{0.5}^{0.5} f(X|C) = \int_{0.5}^{0.5} f(X|C) \int_{0.5}^{0.5} f(X|C) dx = \int_{0.5}^{0.5} \frac{f(x|C)}{f(x|C)} dx$

```
= \begin{cases} 9(2|x,c) & f(2|x,c) \\ 9(2|x,c) & f(2|x,c) \\ \hline 9(2|x,c) & f(2|x,c) \\ \hline \end{cases}
                 = \int \frac{2(2|x|c)}{2(2|x|c)} \int \frac{P(2\cdot x|c)}{2(2|x|c)} dz + \int \frac{2(3|x|c)}{2(3|x|c)} dz
              = \int \frac{2(2|x_1c)}{2(2|x_1c)} \int \frac{P(2|x_1c)}{2(2|x_1c)} d2 + k \left[ \left( \frac{2(2|x_1c)}{2(2|x_1c)} \right) \right] P(2|x_1c)
               = /25 P(X|4) = Lb + KL (2(3|x.4) || P(2)x.4))
\frac{1}{2} = \int_{0}^{2} \frac{1}{4} \left( \frac{1}{2} | x \cdot c \right) \int_{0}^{2} \frac{1}{4} \left( \frac{1}{2} | x \cdot c \right) \frac{1}{4} \left( \frac{1}{2} | x \cdot c \right)
     = ) 8(\frac{1}{2}|x\cdot c) log \frac{p(\frac{1}{2}|c)}{2(\frac{1}{2}|x\cdot c)} dz + ) 8(\frac{1}{2}|x\cdot c) log p(x|2\cdot c) dz
     = - K [ (8(3|x,c) || P(3|4)) + Eg(3|x,c) [ log P(x|3,c)]
3. Implementation details
3-1. Describe how you implement your model (encoder, decoder, reparameterization
trick, dataloader, etc.)
Encoder 和 Decoder 的部分我並沒有做太多修改,使用原本 sample code 提供的 vgg encoder 和 decoder
 class vgg_encoder(nn.Module):
     def __init__(self, dim):
```

self.c2 = nn.Sequential(vgg_layer(64, 128), vgg_layer(128, 128),

self.upc2 = nn.Sequential(

self.upc3 = nn.Sequential(

16 x 16

return output

def get_csv(self, index):

def __getitem__(self, index):

seq = self.get_seq(index)

cond = self.get_csv(index)

mse += mse_criterion(x_pred, x[i])

super(lstm, self).__init__()

for _ in range(self.n_layers):

def forward(self, cond, input):

return self.output(h_in)

cond_embedded = self.cond_embed(cond)

if epoch >= args.tfr_start_decay_epoch:

args.tfr -= args.tfr_decay_step

args.tfr = max(args.tfr, 0)

Update teacher forcing ratio

for i in range(self.n_layers):

h_in = self.hidden[i][0]

embedded = self.embed(torch.cat([cond_embedded, input], dim = 1))

self.hidden[i] = self.lstm[i](h_in, self.hidden[i])

的過程中會慢慢減少使用 GT 的方式並且改使用 model 自己生成的輸出來當作下一層的輸入

--- Teacher ratio

total loss

--- KL weight

250

self.device = device

loss = mse + kld * beta

loss.backward()

class lstm(nn.Module):

hidden = []

return hidden

h_in = embedded

TODO

Teacher forcing ratio:

20

loss/psnr 15

10

20

25

20

loss/psnr 12

10

0

0

50

100

4-3. Discussion

都是以 cyclical 作為標準

kld += kl_criterion(mu, logvar, args)

beta = kl_anneal.get_beta(args, epoch, start_epoch)

訊和 h 以及 z concatenate 再一起之後按到 lstm 的方式得到要送到 decoder 前的 latent code

self.set_seed(index)

return seq, cond

前的 latent code

start = 0

gradient 的問題

vgg_layer(512*2, 512),

vgg_layer(256*2, 256),

vgg_layer(256, 256),

vgg_layer(512, 512),

vgg_layer(512, 256)

self.c1 = nn.Sequential(

self.dim = dim

64 x 64

32 x 32

super(vgg_encoder, self).__init__()

vgg_layer(3, 64),

vgg_layer(64, 64),

16 x 16 self.c3 = nn.Sequential(vgg_layer(128, 256), vgg_layer(256, 256), vgg_layer(256, 256),

```
self.c4 = nn.Sequential(
        vgg_layer(256, 512),
```

```
vgg_layer(512, 512),
                 vgg_layer(512, 512),
        # 4 x 4
        self.c5 = nn.Sequential(
                nn.Conv2d(512, dim, 4, 1, 0),
                nn.BatchNorm2d(dim),
                nn.Tanh()
        self.mp = nn.MaxPool2d(kernel_size=2, stride=2, padding=0)
    def forward(self, input):
        h1 = self.c1(input) # 64 -> 32
        h2 = self.c2(self.mp(h1)) # 32 -> 16
        h3 = self.c3(self.mp(h2)) # 16 -> 8
        h4 = self.c4(self.mp(h3)) # 8 \rightarrow 4
        h5 = self.c5(self.mp(h4)) # 4 \rightarrow 1
        return h5.view(-1, self.dim), [h1, h2, h3, h4]
class vgg_decoder(nn.Module):
    def __init__(self, dim):
        super(vgg_decoder, self).__init__()
        self.dim = dim
        self.upc1 = nn.Sequential(
                nn.ConvTranspose2d(dim, 512, 4, 1, 0),
                nn.BatchNorm2d(512),
                nn.LeakyReLU(0.2, inplace=True)
```

```
vgg_layer(256, 128)
    # 32 x 32
   self.upc4 = nn.Sequential(
            vgg_layer(128*2, 128),
            vgg_layer(128, 64)
   # 64 x 64
   self.upc5 = nn.Sequential(
            vgg_layer(64*2, 64),
            nn.ConvTranspose2d(64, 3, 3, 1, 1),
            nn.Sigmoid()
   self.up = nn.UpsamplingNearest2d(scale_factor=2)
def forward(self, input):
   vec, skip = input
   d1 = self.upc1(vec.view(-1, self.dim, 1, 1)) # 1 -> 4
   up1 = self.up(d1) # 4 -> 8
   d2 = self.upc2(torch.cat([up1, skip[3]], 1)) # 8 x 8
   up2 = self.up(d2) # 8 -> 16
   d3 = self.upc3(torch.cat([up2, skip[2]], 1)) # 16 x 16
   up3 = self.up(d3) # 8 -> 32
   d4 = self.upc4(torch.cat([up3, skip[1]], 1)) # 32 x 32
   up4 = self.up(d4) # 32 -> 64
```

output = self.upc5(torch.cat([up4, skip[0]], 1)) # 64 x 64

position = self.positions[index, :self.seq_length] action = self.actions[index, :self.seq_length] condition = np.hstack((position, action))

return torch.from_numpy(condition).to(torch.float32)

use_teacher_forcing = True if random.random() < args.tfr else False</pre>

 $h_{seq} = [modules['encoder'](x[i]) for i in range(args.n_past+args.n_future)]$

```
def reparameterize(self, mu, logvar):
              std = torch.exp(0.5*logvar)
              eps = torch.randn_like(std)
              return mu + eps*std
Dataloader 的部分,由於我們最後是想要用前兩個 frame 來預測後面 10 個 frame ,因此我改寫成只會讀取前 12 張 frame 送
到 Network 做訓練,另外 condition 的部分,會將 action 和 position concatenate 在一起
   def __len__(self):
      return self.length
   def load_img(self, video, frame):
      img = cv2.imread(self.data_path + video + f'/{frame}.png')
      return self.transform((cv2.cvtColor(img, cv2.COLOR_BGR2RGB))/255.0)
   def get_seq(self, index):
      video = self.videos[index]
      seq = torch.stack([self.load_img(video, i).to(torch.float32) for i in range(self.seq_length)], dim=0
      return seq
```

Reparameterization Trick: 將 Network 預測出來的 μ 加上 var * $\mathcal{N}(0,1)$ (normal distribution) 來解決 sample 沒辦法做

```
## TODO
for i in range(1, args.n_past + args.n_future):
    if args.last_frame_skip or i < args.n_past:</pre>
        h, skip = h_seq[i-1]
    else:
        if use_teacher_forcing:
            h, skip = h_seq[i-1]
        else:
            h = modules['encoder'](x_pred)[0]
   h_target = h_seq[i][0]
    z_t, mu, logvar = modules['posterior'](h_target)
    h_pred = modules['frame_predictor'](cond[i-1], torch.cat([h, z_t], dim = 1))
    x_pred = modules['decoder']([h_pred, skip])
```

當要使用 Teacher forcing 的時候,我會將上一個 frame 由 model predict 出來的結果來通過 encoder 當作下一層要送到 Istm

```
self.input_size = input_size
   self.output_size = output_size
   self.hidden_size = hidden_size
   self.batch_size = batch_size
   self.n_layers = n_layers
   self.cond_embed = nn.Linear(condition_size, 2)
   self.embed = nn.Linear(2 + input_size, hidden_size)
   self.lstm = nn.ModuleList([nn.LSTMCell(hidden_size, hidden_size) for i in range(self.n_layers)])
    self.output = nn.Sequential(
           nn.Linear(hidden_size, output_size),
           nn.BatchNorm1d(output_size),
            nn.Tanh())
   self.hidden = self.init_hidden()
def init_hidden(self):
```

hidden.append((Variable(torch.zeros(self.batch_size, self.hidden_size).to(self.device)),

Variable(torch.zeros(self.batch_size, self.hidden_size).to(self.device))))

Condition 的部分,我會將 condition 先經過 nn.Linear 將 dim 從 7 → 2 (Action + Position),之後將 embedded 的 condition 資

def __init__(self, condition_size, input_size, output_size, hidden_size, n_layers, batch_size, device):

```
Drawbacks: 由於 model 會過度依賴 GT 的結果,因此當訓練和預測的 dataset 分布差異過大時,model 會因為受到 GT 的影
響過渡矯正,造成在預測結果上可能會表現不好
4. Results and discussion
4-1. Show your results of video prediction
(a) Make videos or gif images for test result
Show at demo
(b) Output the prediction at each time step
4-2. Plot the KL loss and PSNR curves during training
                                          (b) cyclical
(a) monotonic
             Training loss/ratio curve
                                                       Training loss/ratio curve
                                                  25 -
                                            25
            0.8
```

20

usd/s

10

50

100

150

epochs

Training loss/ratio curve

Training loss/ratio curve

mse

_ _psnr___

250

0.0

300

200

0.6

0.4

0.2

1.0

0.6

0.4

-- Teacher ratio

KL weight

total loss

250

args.tfr_decay_step = 1.0 / (args.tfr_truncate_epoch - args.tfr_start_decay_epoch)

3-2. Describe the teacher forcing (including main idea, benefits and drawbacks.)

Teacher forcing 的想法是我們在訓練序列生成的時候,會想要將 GT 來取代 Decoder 的輸出並且送到下一層的輸入,在訓練

Benefits: 利用 teacher forcing 的方式可以讓 model 比較好收斂,因為 model 可以在訓練初期學到正確的輸入送到下一層,

反之,若不使用 teacher forcing 容易造成 model 預測的和 GT 差距過大造成 loss 容易發散並且不容易收斂

Training loss/ratio curve 25

慢開始上升,代表不再受到 Training GT 的影響效果會開始變

Training loss/ratio curve

100

epochs

20 0.6 --- Teacher ratio usd/ssol KL weight 0.4 10 10 0.2 kld - kld 5 total loss total loss Teacher ratio mse KL weight 100 100 50 150 200 150 200 250 250 300 epochs epochs 在下面兩張圖,我分別針對 teacher forcing ratio 從 100 → 200 降到 0 以及 150 → 300 降到 0 ,在過程中我發現如果都使用

teacher forcing 來訓練模型,可能會因為 training 和 validation 的資料集分佈有差距,如果 model 被強制學成跟 training 一樣

的分布畫,在 validation 的效果可能不會太好,因此可以看到當開始可以使用模型自己的輸出當作下一層輸入時,PSNR 會慢

25 -

Jusd/ss

10

1.0

0.6

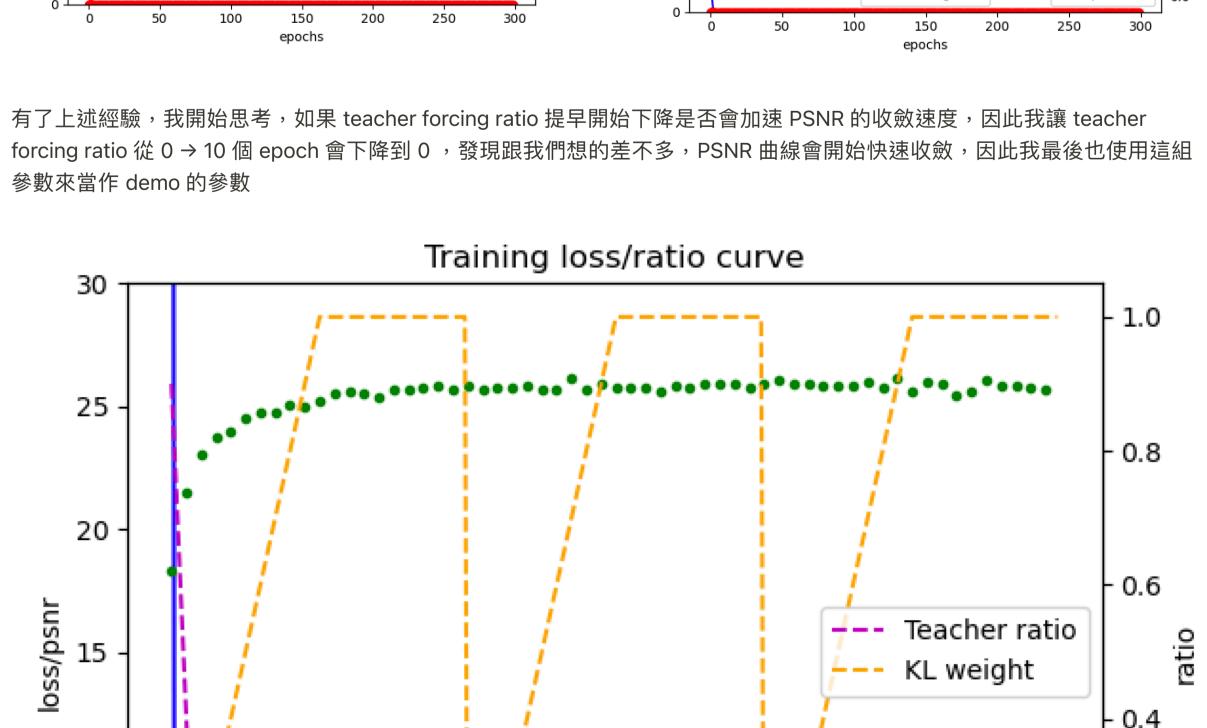
0.4

我一開始使用的 teacher forcing ratio 是到 150個 epoch 會線性遞減到 0 , 在這個過程中我發現 cyclical 的收斂速度比

monotonic 還要快,推測使用 cyclical 的方式對我們這次的 task 幫助較大 ,因此後面我嘗試改變不同的 teacher forcing ratio

25

0.2 total loss total loss Teacher ratio Teacher ratio mse 📏 KL weight KL weight psnr 250 100 150 200 300 50 100 150 200 250 300 epochs epochs



0.4 10 - 0.2 kld 5 total loss

150

epochs

200