

Conditional VAE for Video Prediction  
Lab Report # 4

By  
312581020  
許瀚丰

Deep Learning  
Spring 2024  
Date Submitted: May 7, 2024

# Contents

<b>1</b>	<b>Derivate conditional VAE formula</b>	<b>2</b>
<b>2</b>	<b>Introduction</b>	<b>3</b>
2.1	Problem Statement . . . . .	3
<b>3</b>	<b>Implementation Details</b>	<b>3</b>
3.1	How do you write your training protocol . . . . .	3
3.2	How do you implement reparameterization tricks . . . . .	5
3.3	How do you set your teacher forcing strategy . . . . .	5
3.4	How do you set your kl annealing ratio . . . . .	5
<b>4</b>	<b>Analysis &amp; Discussion</b>	<b>7</b>
4.1	Plot Teacher forcing ratio . . . . .	7
4.2	Plot the loss curve while training with different settings . . . . .	8
4.3	Plot the PSNR-per frame diagram in validation dataset . . . . .	9
4.4	Other training strategy analysis . . . . .	9
4.4.1	Model Modification . . . . .	9

# Listings

1	Training . . . . .	4
2	reparameterization tricks . . . . .	5
3	teacher forcing . . . . .	5
4	KL annealing . . . . .	6
5	Generator . . . . .	9

# 1 Derivate conditional VAE formula

對於一個conditional的機率分布 $p(x|c; \theta)$ ，我們目標是最大化其的log-likelihood，原式如下

$$\max_{\theta} \int \log p(x|c; \theta) dx \quad (1)$$

假設有一個任意的機率分布 $q(z|c)$ ，我們可以將原式修改為以下

$$\begin{aligned} \log p(x|c; \theta) &= \int q(z|c) \log p(x|c; \theta) dz \\ &= \int q(z|c) \log \left( \frac{p(x, z|c; \theta)}{p(z|x, c; \theta)} \right) dz \\ &= \int q(z|c) \log \left( \frac{p(x, z|c; \theta)}{q(z|c)} \frac{q(z|c)}{p(z|x, c; \theta)} \right) dz \\ &= \int q(z|c) \log \left( \frac{p(x, z|c; \theta)}{q(z|c)} \right) dz + \int q(z|c) \log \left( \frac{q(z|c)}{p(z|x, c; \theta)} \right) dz \\ &= \int q(z|c) \log \left( \frac{p(x, z|c; \theta)}{q(z|c)} \right) dz + KL(q(z|c) || p(z|x, c; \theta)) \end{aligned} \quad (2)$$

由於我們知道對於任意分布 $p, q$ ， $KL(p||q) \geq 0$ ，因此我們可以得知 $\int q(z|c) \log \left( \frac{p(x, z|c; \theta)}{q(z|c)} \right) dz$ 必為 $\log p(x|c; \theta)$ 的下界，也就是

$$\log p(x|c; \theta) \geq \int q(z|c) \log \left( \frac{p(x, z|c; \theta)}{q(z|c)} \right) dz \quad (3)$$

我們將右式定義為 $L(X, c, q, \theta)$ 並其展開，可以得到以下

$$\begin{aligned} L(X, c, q, \theta) &= \int q(z|c) \log \left( \frac{p(x, z|c; \theta)}{q(z|c)} \right) dz \\ &= \int q(z|c) \log \left( \frac{p(x|z, c; \theta) p(z|c; \theta)}{q(z|c)} \right) dz \\ &= \int q(z|c) \log p(x|z, c; \theta) dz + \int q(z|c) \log \left( \frac{p(z|c; \theta)}{q(z|c)} \right) dz \\ &= \int q(z|c) \log p(x|z, c; \theta) dz - KL(q(z|c) || p(z|c; \theta)) \end{aligned} \quad (4)$$

若此時的 $q(z|c)$ 是由另一個網路 $\phi$ 與輸入資料 $x$ 產生，也就是 $q(z|x, c; \phi)$ ，我們可將原本的結果改寫為以下

$$\begin{aligned}
& \int q(z|c) \log p(x|z, c; \theta) dz - KL(q(z|c) || p(z|c; \theta)) \\
&= \int q(z|c, x; \phi) \log p(x|z, c; \theta) dz - KL(q(z|c, x; \phi) || p(z|c; \theta)) \\
&= E_{z \sim q(z|x, c; \phi)} \log p(x|z, c; \theta) - KL(q(z|x, c; \phi) || p(z|c; \theta))
\end{aligned} \tag{5}$$

## 2 Introduction

### 2.1 Problem Statement

在本次實驗中，我們需要實作一個cVAE的模型來完成Video Prediction的任務，在訓練時共有兩個輸入，分別為目前人物在上一幀 $X_{t-1}$ 中畫面的影像與其當前這一幀的Pose影像 $P_t$ ，而我們需要做的是利用這兩者來去預測出在此時的畫面 $\hat{X}_t$ 。而整體訓練的目標是希望能讓預測結果 $\hat{X}_t$ 與實際的 $X_t$ 越像越好，並透過的設定不同Teacher Forcing與KL Annealing的方式來幫助訓練，最後會由PSNR作為此次實驗評分的指標。

## 3 Implementation Details

### 3.1 How do you write your training protocol

在訓練中，每次輸入進來的shape為(Batch\_size, Time\_step, Channel, Height, Width)，因此第一個迴圈對於每個Batch，而對於每個Batch的訓練，我將Frame與Pose的組合 $(X_0, P_1)$ 經由模型所生成 $\hat{X}_1$ 與真實的 $X_1$ 計算MSELoss，之後則使用預測出的Frame  $\hat{X}_1$ 取代真實的 $X_1$ 作為新的輸入 $(\hat{X}_1, P_2)$ 並計算Loss，一直持續到 $(\hat{X}_{14}, P_{15})$ 為止。而KLD的部分與前面類似，差別僅在於我們並不需要生成 $\hat{X}_0$ ，因此使用的是 $(X_1, P_1)$ 到 $(X_{15}, P_{15})$ 的組合來計算KLD，最終只需要將兩者相加並更新即可。而Teacher Forcing的部分僅需要簡單的將 $\hat{X}_i$ 替換為真實答案 $X_i$ ，而KL annealing的部分就是直接在一開始先取得beta值，最後再加總Loss Function時將其乘上KLD即可，如程式碼1所示。

```

1 def training_one_step(self, img_batch, label_batch, adapt_TeacherForcing):
2     # img_batch: (Batch_size, Time_step, Channel, Height, Width) = (2, 16, 3, 32, 64)
3     batch_size = img_batch.shape[0] # 2
4     total_loss = 0
5     beta = self.kl_annealing.get_beta()
6
7     for i in range(batch_size):
8         img = img_batch[i] # img: (Time_step, Channel, Height, Width)
9         label = label_batch[i] # label: (Time_step, Channel, Height, Width)
10        mse = 0
11        kld = 0
12        x_hat = img[0].unsqueeze(0)
13
14        for j in range(1, img.size(0)):
15            if adapt_TeacherForcing:
16                prev_img = img[j - 1].unsqueeze(0)
17            else:
18                prev_img = x_hat
19
20            frame_enc = self.frame_transformation(img[j].unsqueeze(0))
21            label_enc = self.label_transformation(label[j].unsqueeze(0))
22            z, mu, logvar = self.Gaussian_Predictor(frame_enc, label_enc)
23            kld += kl_criterion(mu, logvar, batch_size)
24
25            prev_frame_enc = self.frame_transformation(prev_img).detach()
26            decoded = self.Decoder_Fusion(prev_frame_enc, label_enc, z)
27            x_hat = self.Generator(decoded)
28            mse += self.mse_criterion(x_hat, img[j].unsqueeze(0))
29        loss = mse + beta * kld
30
31        self.optim.zero_grad()
32        loss.backward()
33        self.optimizer_step()
34
35        total_loss += loss
36
37    return total_loss / batch_size

```

Listing 1: Training

### 3.2 How do you implement reparameterization tricks

我們假設模型輸出結果為 $\mu$ 與 $\log \sigma^2$ ，為何使用 $\log \sigma^2$ 而非直接使用 $\sigma$ 是因為後者的輸出要求為非負實數，此項constraint會使模型不好做訓練。若我們要得到 $\sigma$ 只需要將 $\log \sigma^2$ 除以2在做exp並即可。

整體reparameterization tricks的部分，我們需要在給定的 $\mu$ 與 $\sigma$ 的Normal Distribution中做抽樣，因此我先使用randn\_like從 $N(0, 1)$ 抽樣，並將其乘以 $\sigma$ 再加上 $\mu$ 就是 $N(\mu, \sigma)$ 了，如程式碼2所示。

```
1 def reparameterize(self, mu, logvar):
2     std = torch.exp(0.5 * logvar)
3     eps = torch.randn_like(std)
4     return eps * std + mu
```

Listing 2: reparameterization tricks

### 3.3 How do you set your teacher forcing strategy

要設定Teacher Forcing Ratio(tfr)共有三個參數，分別為tfr, tfr\_sde, tfr\_d\_step，分別代表目前的tfr，多少epoch要降低tfr，每次降低多少tfr。而在我的實作中，就是依照上述的設定去撰寫，每次更新皆會先判斷當前的epoch是否需要更新，而更新時就直接將tfr減去tfr\_d\_step，需要特別注意不要讓其降到0以下，如程式碼3所示。

```
1 def teacher_forcing_ratio_update(self):
2     if self.current_epoch >= self.tfr_sde and self.current_epoch % self.tfr_sde == 0:
3         self.tfr -= self.tfr_d_step
4         self.tfr = max(0, self.tfr)
```

Listing 3: teacher forcing

### 3.4 How do you set your kl annealing ratio

在KL annealing的部分，我直接參考了原始論文的做法，在實驗中我共使用了三種不同的方法，分別為Cyclical, Monotonic, Without(直接設定為1)，且由於在實驗中我發現，由於神經網路較為複雜，因此一開始若完全不使用KL Divergence來更新網路，

可能會導致後續計算Loss產生Nan，因此在一開始我會先會讓計算出的beta不要直接為0來防止Nan的發生。

```
1 class kl_annealing():
2     def __init__(self, args, current_epoch=0):
3         self.annealing_type = args.kl_anneal_type
4         assert self.annealing_type in ["Cyclical", "Monotonic", "Without"]
5         self.iter = current_epoch + 1
6
7         if self.annealing_type == "Cyclical":
8             self.L = self.frange_cycle_linear(num_epoch=args.num_epoch, start=0.0,
9 stop=1.0, n_cycle=args.kl_anneal_cycle, ratio=args.kl_anneal_ratio)
10        elif self.annealing_type == "Monotonic":
11            self.L = self.frange_cycle_linear(num_epoch=args.num_epoch, start=0.0,
12 stop=1.0, n_cycle=1, ratio=args.kl_anneal_ratio)
13        else:
14            self.L = np.ones(args.num_epoch + 1)
15
16    def update(self):
17        self.iter += 1
18
19    def get_beta(self):
20        return self.L[self.iter]
21
22    def frange_cycle_linear(self, num_epoch, start=0.0, stop=1.0, n_cycle=1, ratio=1):
23        # adapted from https://github.com/haofuml/cyclical_annealing
24        L = np.ones(num_epoch + 1)
25        period = num_epoch / n_cycle
26        step = (stop - start) / (period * ratio)
27        for c in range(n_cycle):
28            v, i = start, 0
29            while v <= stop and (int(i+c*period) < num_epoch):
30                L[int(i + c * period)] = v
31                v += step
32                i += 1
33        return L
```

Listing 4: KL annealing

## 4 Analysis & Discussion

### 4.1 Plot Teacher forcing ratio

在本次實驗中，在沒有使用KL aneling的情況(kl-type為Without)下我共測試以下四種不同的Teacher forcing ratio(tfr)設定：

- tfr完全為0：(kl-type\_Without\_tfr\_0.0\_teacher-decay\_0.1)
- tfr完全為0.1：(kl-type\_Without\_tfr\_0.1\_teacher-decay\_0.0)
- tfr一開始為0.5，之後每10個epoch下降0.1：(kl-type\_Without\_tfr\_0.5\_teacher-decay\_0.1)
- tfr一開始為1，之後每5個epoch下降0.5：(kl-type\_Without\_tfr\_1.0\_teacher-decay\_0.5)

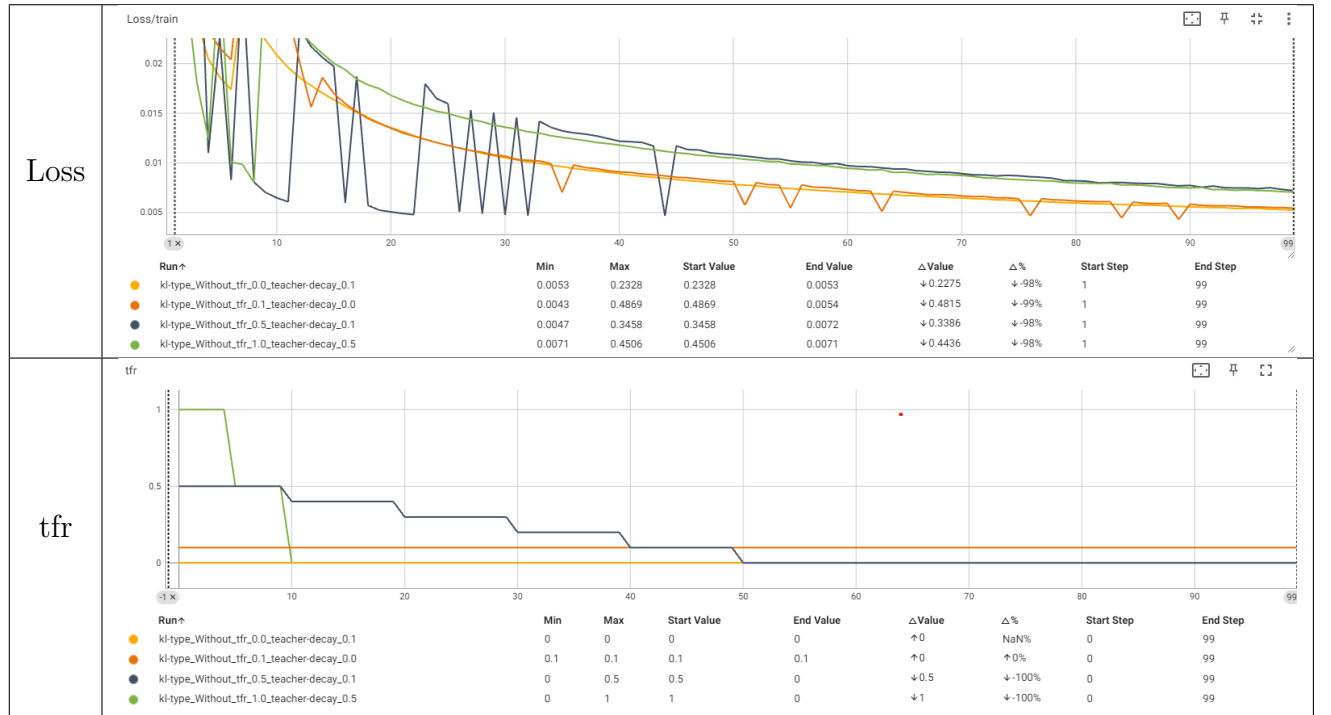


Table 1: Teacher forcing ratio Comparison

實驗結果如表1，從實驗中可以發現，對於本次任務來說，使用Teacher Forcing的效果並不好，甚至容易讓模型提早收斂。我認為其原因可能是因為在Training時，同一個Epoch只會有全部都使用Teacher Forcing或全部都不使用Teacher Forcing，因



此可能會讓模型在過程中陷入local minimum，而在因此在切換時(從使用tfr轉換成不使用tfr, vice versa)就容易產生震盪，進而導致結果受到影響。且由於本次任務的condition(Pose)與要預測的結果(Frame)較為類似，因此儘管不使用Teacher Forcing也能產生出不錯的結果。

## 4.2 Plot the loss curve while training with different settings

在本次實驗中，在固定Teacher forcing ratio為0的情況下，我共測試以下三種不同的設定：

- KL annealing (Monotonic): kl-type\_Monotonic\_tfr\_0.0\_teacher-decay\_0.0
- KL annealing (Cyclical): kl-type\_Cyclical\_tfr\_0.0\_teacher-decay\_0.0
- Without KL Annealing: kl-type\_Without\_tfr\_0.0\_teacher-decay\_0.0

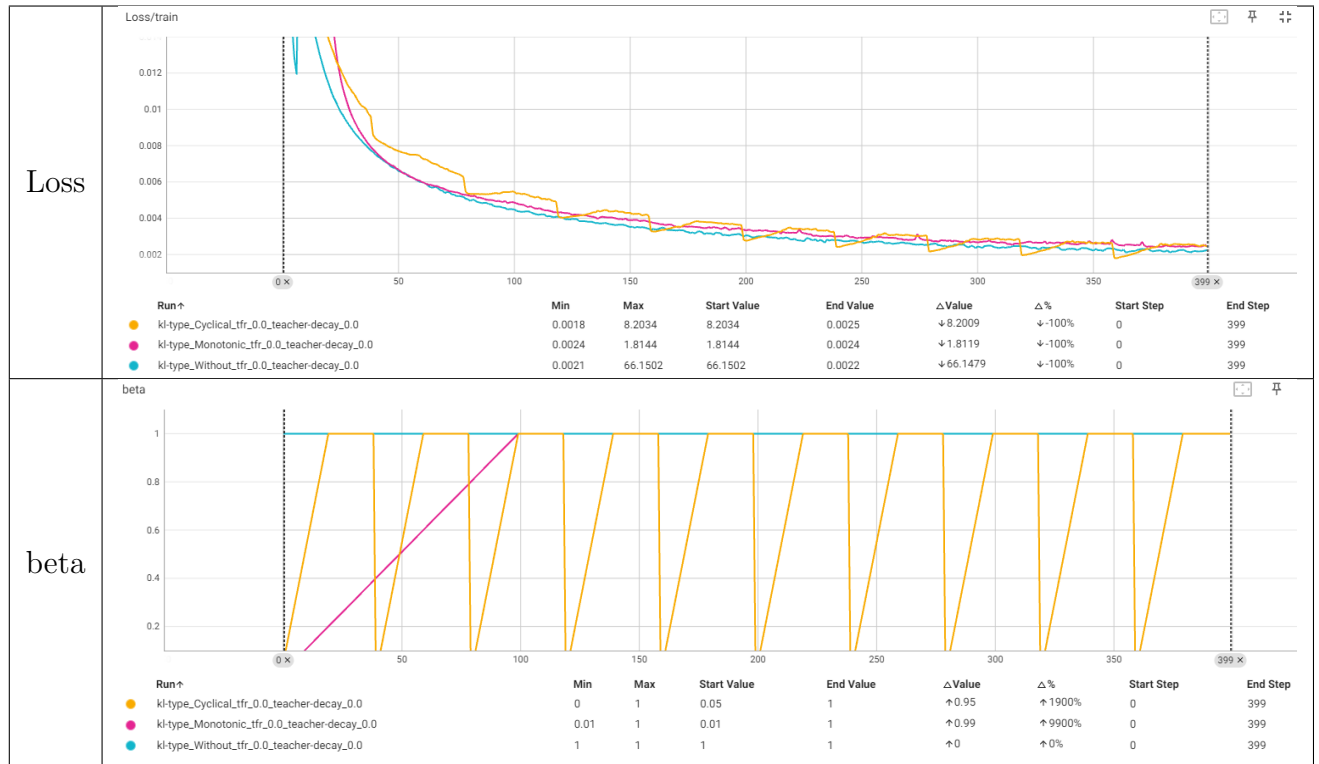


Table 2: KL annealing Comparison

實驗結果如表2，從實驗中可以發現，對於此任務來說，使用不同的KL annealing策略基本上差異不大，但仍可以觀察到使用Cyclical的方式Loss確實會有如設定上的震盪。結束訓練時三種方式的Loss都降至0.002左右，可見三種方式都能有效的讓模型收斂。

### 4.3 Plot the PSNR-per frame diagram in validation dataset

實驗結果如3，三者皆是固定Teacher forcing ratio為0且訓練400 epoch並比較更改不同KL Annealing所產生出的結果。可以觀察三者的平均PSNR皆有37以上，且每個Frame的PSNR差異並不大。

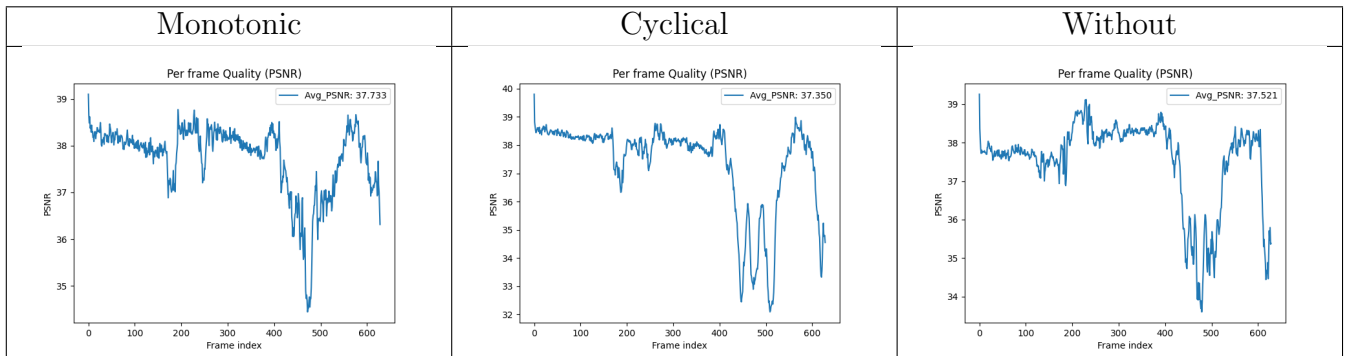


Table 3: PSNR

## 4.4 Other training strategy analysis

### 4.4.1 Model Modification

我們對圖片做前處理時，會將其normalize到0到1之間，然而在原始程式中，Generator在產生出結果時是透過最後的nn.Conv2d直接產生出結果，因此並不會限制pixel皆在0到1之間，因此我在Generator輸出前加上一個Sigmoid Function，在確保圖片可以限制在0到1之間的同時，也可以讓模型在訓練前期更加穩定，如程式碼5。

```

1 class Generator(nn.Sequential):
2     def __init__(self, input_nc, output_nc):
3         super(Generator, self).__init__(
4             DepthConvBlock(input_nc, input_nc),
5             ResidualBlock(input_nc, input_nc//2),

```

```

6         DepthConvBlock(input_nc//2, input_nc//2),
7         ResidualBlock(input_nc//2, input_nc//4),
8         DepthConvBlock(input_nc//4, input_nc//4),
9         ResidualBlock(input_nc//4, input_nc//8),
10        DepthConvBlock(input_nc//8, input_nc//8),
11        nn.Conv2d(input_nc//8, 3, 1),
12
13        nn.Sigmoid()
14    )
15
16    def forward(self, input):
17        return super().forward(input)

```

Listing 5: Generator