

HW10 Image captioning with RNNs and LSTMs

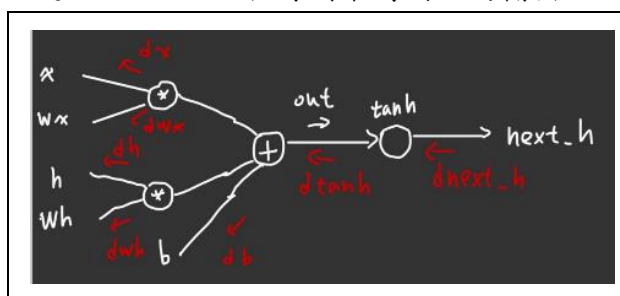
M113040105 劉東霖

壹. 使用到的 function:

一. rnn_step_forward:

首先，先來講解 RNN 單部 forward 的部分。

如下圖所示，將輸入 x 乘上權重 W_x 加上上一個隱藏層狀態 $prev_h$ 乘上權重 W_h 再加上 bias，再將結果經過 \tanh 激活，就得到本時間段的輸出。

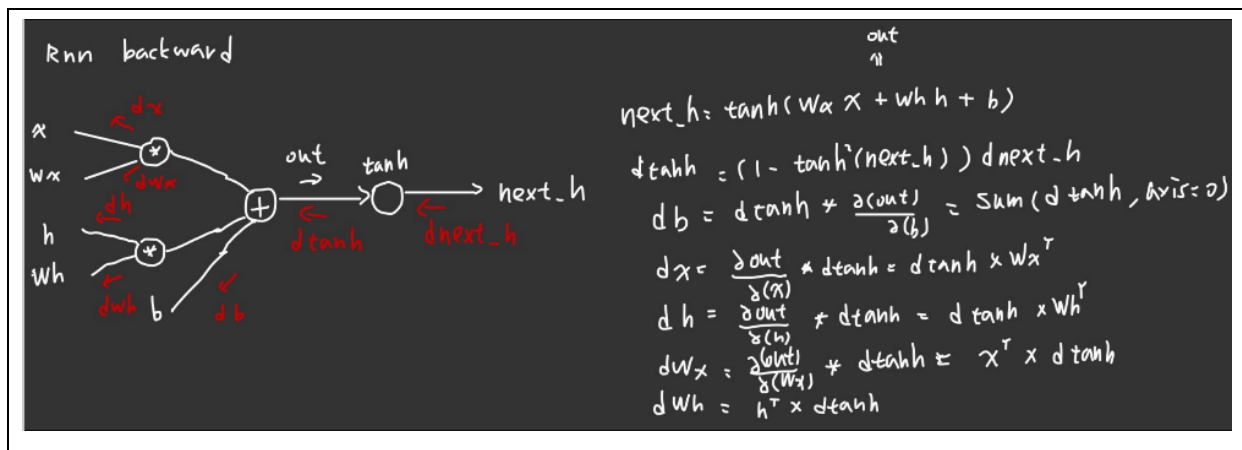


程式碼如下，最後把參數用 cache 存起來方便 backward。

```
next_h=torch.tanh((x @ Wx) + (prev_h @ Wh)+b)
cache=(x, prev_h, Wx, Wh, next_h)
```

二. rnn_step_backward:

下圖為 RNN 單部 backward 的公式推導

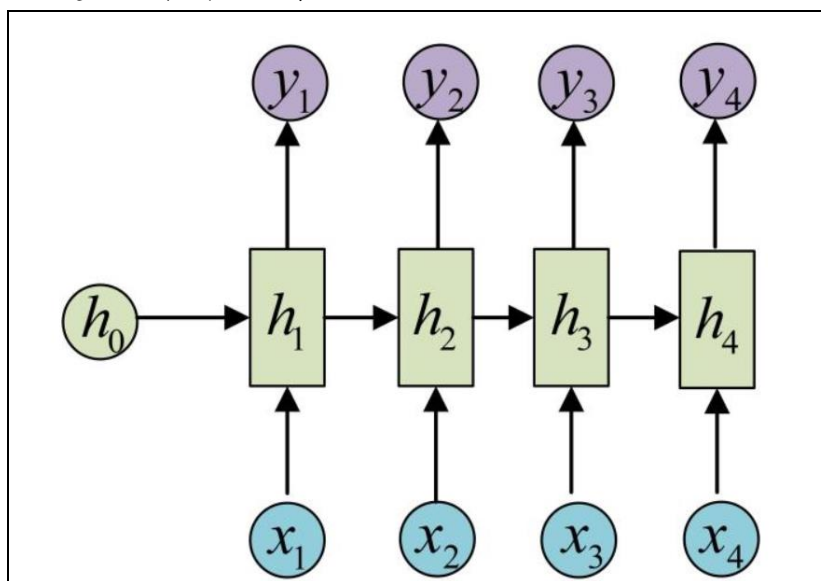


程式碼如下，首先把參數從 cache 裡面拿出來，再根據公式完成 backward。

```
(x, prev_h, Wx, Wh, next_h)=cache
dtanh=dnext_h*(1-torch.pow(next_h, 2))
dx=dtanh.mm(Wx.T)
dprev_h=dtanh.mm(Wh.T)
dWx=x.T.mm(dtanh)
dWh=prev_h.T.mm(dtanh)
db=dtanh.sum(axis=0)
```

三. rnn_forward:

前面已經完成單部 forward 了，接下來要把前面所寫的單部 forward 組合起來實現 RNN，這樣就可以處理一系列的數據。



首先定義一些變數， N 為 batch size， H 為 hidden size 的維度， T 為此段 sequence 的長度， D 為每一段 feature 的維度。

```
H=h0.shape[1]  
N,T,D=x.shape
```

創建一個 tensor h ，為了記錄每一次的狀態，並用一個 list cache 儲存中間的參數。再來把初始狀態設定成 h_0 。

```
h=torch.zeros((N,T,H),dtype=x.dtype,device=x.device)  
cache=[]  
prev_h=h0
```

再來進行 RNN forward 計算，用一個 for 對每一個時間進行迭代。 x_t 為輸入 x 在第 i 個時間的值，並把 x_t 、 $prev_h$ 、 W_x 、 W_h 、 b 代入之前寫的 `rnn_step_forward` 得到現在的狀態 th 和現在的參數 $tcache$ 。做完 forward 後，把 $prev_h$ 更新成 th ，並把每一次的參數和狀態記錄起來。程式碼如下：

```
for i in range(T):  
    xt=x[:,i,:]  
    th,tcache=rnn_step_forward(xt,prev_h,Wx,Wh,b)  
    prev_h=th  
    h[:,i,:]=th  
    cache.append(tcache)
```

四. rnn_backward:

前面已經完成單部 backward 了，接下來要把前面所寫的單部 backward 組合起來實現 RNN。

首先定義一些變數。 D 為輸入 data 的維度， N 和 T 和 H 分別為 batch size，hidden size，sequence 的長度，並把一些需要更新的參數初始化。

```

D=cache[0][0].shape[1]
N, T, H=dh.shape
dx=torch.zeros((N, T, D), dtype=dh.dtype, device=dh.device)
dh0=torch.zeros((N, H), dtype=dh.dtype, device=dh.device)
dWx=torch.zeros((D, H), dtype=dh.dtype, device=dh.device)
dWh=torch.zeros((H, H), dtype=dh.dtype, device=dh.device)
db=torch.zeros((H, ), dtype=dh.dtype, device=dh.device)
tdprev_h=torch.zeros((N, H), dtype=dh.dtype, device=dh.device)

```

再來進行 RNN backward 計算，用一個 for 對每一個時間進行迭代。先把當前狀態加前一個狀態和當前的參數放到 rnn_step_backward 計算 backward，因每個時間的輸入 x 和前一個時間的 hidden state 都對當前時間的輸出產生影響，因此在計算權重的梯度時需要將所有時間步的貢獻加總起來，所以 dWx 、 dWh 、 db 要累加 backward 算出來的梯度，並將當前時間的梯度 tdx 用 dx 存起來。

最後把最後一個時間的 hidden state 的梯度 $tdprev_h$ 儲存到 $dh0$ 就完成 backward。程式碼如下：

```

for i in range(T-1, -1, -1):
    tdx, tdprev_h, tdWx, tdWh, tdb=rnn_step_backward(dh[:, i, :]+tdprev_h, cache[i])
    dx[:, i, :]+=tdx
    dWx+=tdWx
    dWh+=tdWh
    db+=tdb
dh0=tdprev_h

```

完整程式碼如下：

```

D=cache[0][0].shape[1]
N, T, H=dh.shape
dx=torch.zeros((N, T, D), dtype=dh.dtype, device=dh.device)
dh0=torch.zeros((N, H), dtype=dh.dtype, device=dh.device)
dWx=torch.zeros((D, H), dtype=dh.dtype, device=dh.device)
dWh=torch.zeros((H, H), dtype=dh.dtype, device=dh.device)
db=torch.zeros((H, ), dtype=dh.dtype, device=dh.device)
tdprev_h=torch.zeros((N, H), dtype=dh.dtype, device=dh.device)
for i in range(T-1, -1, -1):
    tdx, tdprev_h, tdWx, tdWh, tdb=rnn_step_backward(dh[:, i, :]+tdprev_h, cache[i])
    dx[:, i, :]+=tdx
    dWx+=tdWx
    dWh+=tdWh
    db+=tdb
dh0=tdprev_h

```

五. word_embedding:

題目要求我們實作 class WordEmbedding 裡面 forward 的部分。這裡只需要一程式碼，用 `self.W_embed[x]` 就可以把一個詞彙轉成對應的向量。程式碼如下：

```

def forward(self, x):
    out = None
    =====
    # TODO: Implement the forward pass for word embeddings.
    =====
    # Replace "pass" statement with your code
    out=self.W_embed[x]
    =====
    #                                     END OF YOUR CODE
    =====
    return out

```

六. temporal_softmax_loss:

我們使用 `nn.functional.cross_entropy` 計算模型預測結果與實際結果之間的差異。輸入 `x` 的維度為 (N, T, V) ，必須先把 `x` 的第 1 和 2 的維度交換，才能在每個時間上進行 `cross_entropy`。

將 `reduction` 設為 `sum` 來算所有時間的損失和忽略 `ignore_index` 所指定的索引，最後再除以 mini batch size 求得平均 loss。

完整程式碼如下：

```
loss = nn.functional.cross_entropy(torch.transpose(x, 1, 2), y,
                                     ignore_index=ignore_index, reduction='sum') / x.shape[0]
```

七. LSTM step_forward:

這裡主要是在實現 LSTM 的核心計算。

1. 計算當前時間的輸出 `a`，`a` 的定義如下圖：

$$a \in \mathbb{R}^{4H} \text{ as } a = W_x x_t + W_h h_{t-1} + b.$$

2. 把 `a` 分成 4 等分，第一份為 `ai`，第 2 份為 `af`，第 3 份為 `ao`，第 4 份為 `ag`。再來把這 4 等份根據下圖的方式帶入不同的 activation function。

$$i = \sigma(a_i) \quad f = \sigma(a_f) \quad o = \sigma(a_o) \quad g = \tanh(a_g)$$

3. 最後在根據下圖的公式來計算下一個 cell state 和 hidden state。

$$c_t = f \odot c_{t-1} + i \odot g \quad h_t = o \odot \tanh(c_t)$$

where \odot is the elementwise product of vectors.

完整程式碼如下：

```
a=prev_h @ self.Wh+x @ self.Wx+self.b
h=a.shape[1]//4
ai,af,ao,ag=a[:, :h], a[:, h:2*h], a[:, 2*h:3*h], a[:, 3*h:]
i,f,o,g=torch.sigmoid(ai), torch.sigmoid(af), torch.sigmoid(ao), torch.tanh(ag)
next_c=f*prev_c+i*g
next_h=o*tanh(torch.tanh(next_c))
```

八. LSTM forward:

1. 首先，建立一個空的 list `h` 來儲存每個隱藏層的狀態，並將最後一個設為初始狀態 `h0`。

2. 將 x 進行轉置，讓 T 變成第一維度，以便我們可以沿著時間計算每個隱藏層的狀態，並將初始 cell state 設為 c_0

3. 沿著每一個時間，代入剛剛寫的 `step_forward` 得出每一個隱藏層狀態和 cell state，並把每一個 hidden state 儲存起來。

4. 用 `torch.stack` 將 h 變成一個張量，並將張量的長度重新排列，使得第一維度為 `batch_size`。

完整程式碼如下：

```
T=x.shape[1]
N,H=h0.shape
h=[None]*T
h[-1]=h0
c=c0
x=x.permute(1,0,2)
for i in range(T):
    h[i],c=self.step_forward(x[i],h[i-1],c)
hn=torch.stack(h).permute(1,0,2)
```

九. CaptioningRNN `__init__`:

首先看目前的 model 是 RNN 還是 LSTM，並定義 model。Model 的輸入維度為 `wordvec_dim`，輸出維度為 `hidden_dim`。程式碼如下：

```
if cell_type == "rnn":
    self.rnn=RNN(wordvec_dim,hidden_dim)
elif cell_type == "lstm":
    self.rnn=LSTM(wordvec_dim,hidden_dim)
```

再來定義圖像編碼的部分。題目使用了 tiny ResNet-X 400MF model 來提取影像的特徵再用一個 feature projection 於將提取出的圖像特徵映射到 RNN/LSTM 的維度。程式碼如下：

```
self.image_encoder=ImageEncoder(pretrained_image_encoder_pretrained)

if cell_type == "rnn" or cell_type == "lstm":
    self.feature_projection=nn.Linear(input_dim*16, hidden_dim)
```

最後定義了兩個子模組的模型。WordEmbedding 用來將輸入的單詞索引轉換成對應的單詞向量，其輸入維度為 `vocab_size`，輸出維度為 `wordvec_dim`。Linear 層用來將 RNN 的輸出向量映射到單詞分類的得分向量，其輸入維度為 `hidden_dim`，輸出維度為 `vocab_size`。程式碼如下：

```
self.embed=WordEmbedding(vocab_size,wordvec_dim)
self.fc=nn.Linear(hidden_dim,vocab_size)
```

完整程式碼如下：

```

if cell_type == "rnn":
    self.rnn=RNN(wordvec_dim,hidden_dim)
elif cell_type == "lstm":
    self.rnn=LSTM(wordvec_dim,hidden_dim)
self.image_encoder=ImageEncoder(pretrained_image_encoder_pretrained)
self.embed=WordEmbedding(vocab_size,wordvec_dim)
self.fc=nn.Linear(hidden_dim,vocab_size)
if cell_type == "rnn" or cell_type == "lstm":
    self.feature_projection=nn.Linear(input_dim*16, hidden_dim)

```

十. CaptioningRNN forward:

程式碼分為以下步驟:

1. 首先，使用圖像編碼器對圖像進行編碼，得到圖像特徵向量。
2. 使用 WordEmbedding 將輸入的語句轉換為序列 embedding 向量 x，以便與圖像特徵向量相連接。
3. 如果模型類型為 RNN 或 LSTM，則使用 feature projection 對圖像特徵向量進行維度轉換，然後將其作為 RNN/LSTM 的初始 hidden state。
4. 將序列 embedding 向量 x 和初始 hidden state 輸入 RNN/LSTM，獲得 RNN/LSTM 的輸出。
5. 將 RNN/LSTM 的輸出通過 self.fc，得到每個單詞在詞彙表中的概率得分。
6. 最後，使用 softmax 損失函數計算預測語句的損失值。

完整程式碼如下:

```

features=self.image_encoder(images)

x=self.embed(captions_in)
if self.cell_type == 'rnn' or self.cell_type == 'lstm':
    h0=self.feature_projection(features.reshape(features.shape[0],-1))
    h=self.rnn(x,h0)
scores=self.fc(h)
loss=temporal_softmax_loss(scores,captions_out,self.ignore_index)

```

十一. CaptioningRNN sample:

程式碼分為以下步驟:

1. 輸入圖片到 image_encoder，產生影像特徵 features。
2. 如果 cell_type 是 rnn，將影像特徵 features reshape 成(N, -1)並通過 feature_projection 產生初始狀態 h。
3. 如果 cell_type 是 lstm，將影像特徵 features reshape 成(N, -1)並通過 feature_projection 產生初始 hidden state 狀態 h，初始 cell state 值為 0。
4. 將每個單詞初始化為起始單詞 self._start，並用 embed 方法轉換成詞向量 w。
5. 如果 cell_type 是 rnn，通過 step_forward 方法產生下一個時間的輸出 h。

6. 如果 cell_type 是 lstm，通過 step_forward 方法產生下一個時間的輸出 hidden state h 和下一個 cell state c。
7. 將 h 通過 self.fc 產生單詞的分數。
8. 用 torch.argmax 方法找到分數最高的單詞並將其作為下一個時間的輸入 words。
9. 將生成的單詞存儲在 captions 中，並
10. 重複步驟 4-9，直到生成的字幕的最大長度。
11. 最後在 captions 的最前面添加起始單詞 self._start。

程式碼如下：

```
features=self.image_encoder(images)
if self.cell_type == 'rnn':
    h=self.feature_projection(features.reshape(features.shape[0],-1))
elif self.cell_type == 'lstm':
    h=self.feature_projection(features.reshape(features.shape[0],-1))
    c=torch.zeros_like(h)
words=self._start*images.new(N,1).long()
for i in range(max_length):
    w=self.embed(words).reshape(N,-1)
    if self.cell_type == 'rnn':
        h=self.rnn.step_forward(w,h)
    elif self.cell_type == 'lstm':
        h,c=self.rnn.step_forward(w,h,c)
    scores=self.fc(h)
    words=torch.argmax(scores,dim=1)
    captions[:,i]=words
starts=torch.full((captions.shape[0], 1), self._start, device=captions.device).long()
captions=torch.cat([starts, captions],dim=1)
```

貳. 程式執行結果：

1. 計算出 WordEmbedding 與預期輸出結果相比，發現差距小。

```
expected_out = torch.tensor(
  [
    [
      [0.0, 0.07142857, 0.14285714],
      [0.64285714, 0.71428571, 0.78571429],
      [0.21428571, 0.28571429, 0.35714286],
      [0.42857143, 0.5, 0.57142857],
    ],
    [
      [0.42857143, 0.5, 0.57142857],
      [0.21428571, 0.28571429, 0.35714286],
      [0.0, 0.07142857, 0.14285714],
      [0.64285714, 0.71428571, 0.78571429],
    ],
  ],
  **to_double
)

print("out error: ", rel_error(expected_out, out))
out error: 2.727272753724473e-09
```

2. 在不同的形狀和條件下多次呼叫 check_loss() 函數來比較 softmax loss，發現輸出都符合預期。


```

check_loss(1000, 1, 10, 1.0) # Should be about 2.00-2.11
check_loss(1000, 10, 10, 1.0) # Should be about 20.6-21.0
check_loss(5000, 10, 10, 0.1) # Should be about 2.00-2.11

2.0354809761047363
20.68878173828125
2.093963861465454

```

3. 代入 CaptioningRNN 裡面計算出 loss，發現沒有差異。

```

loss = model(images, captions).item()
expected_loss = 150.6090393066

print("loss: ", loss)
print("expected loss: ", expected_loss)
print("difference: ", rel_error(torch.tensor(loss), torch.tensor(expected_loss)))

For input images in NCHW format, shape (2, 3, 224, 224)
Shape of output c5 features: torch.Size([2, 400, 7, 7])
loss: 150.60903930664062
expected loss: 150.6090393066
difference: 0.0

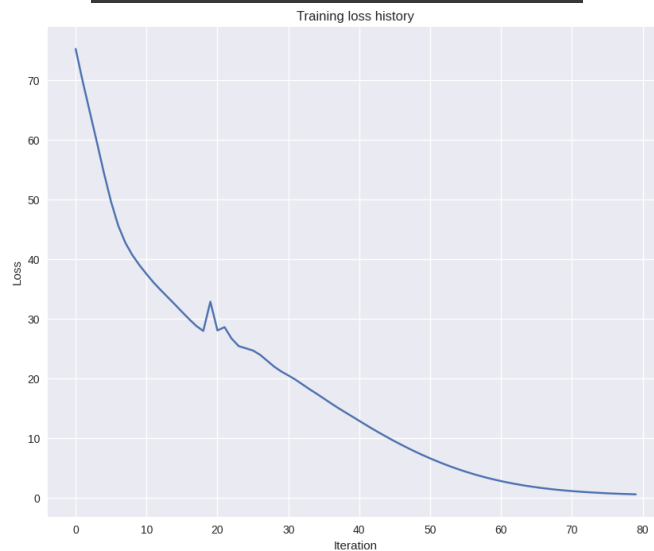
```

4. 使用 train_captioner 來訓練 CaptioningRNN，cell type 為 rnn，且只有使用 50 個 example 來訓練。在 learning rate 為 $1e-3$ 的情況下，發現 loss 小於 1。

```

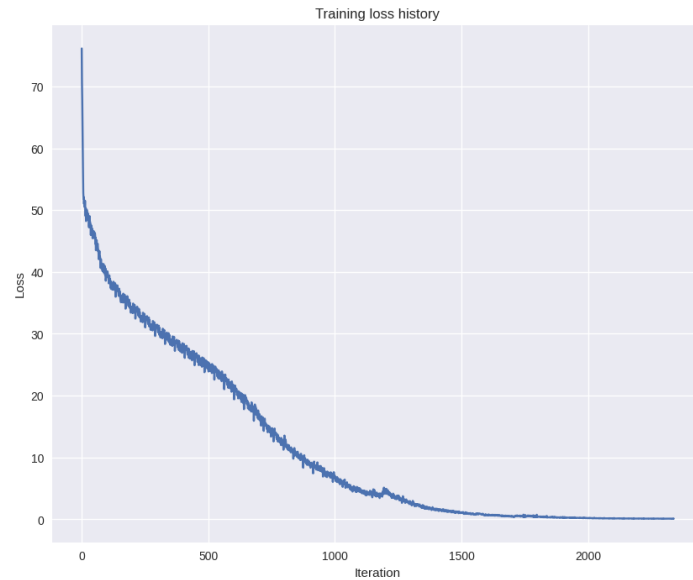
(Epoch 68 / 80) loss: 1.3781 time per epoch: 0.1s
(Epoch 69 / 80) loss: 1.2660 time per epoch: 0.1s
(Epoch 70 / 80) loss: 1.1653 time per epoch: 0.1s
(Epoch 71 / 80) loss: 1.0753 time per epoch: 0.1s
(Epoch 72 / 80) loss: 0.9946 time per epoch: 0.1s
(Epoch 73 / 80) loss: 0.9223 time per epoch: 0.1s
(Epoch 74 / 80) loss: 0.8575 time per epoch: 0.1s
(Epoch 75 / 80) loss: 0.7993 time per epoch: 0.1s
(Epoch 76 / 80) loss: 0.7470 time per epoch: 0.1s
(Epoch 77 / 80) loss: 0.6999 time per epoch: 0.1s
(Epoch 78 / 80) loss: 0.6574 time per epoch: 0.1s
(Epoch 79 / 80) loss: 0.6190 time per epoch: 0.1s

```



5. 使用 train_captioner 來訓練 CaptioningRNN，cell type 為 rnn。在 learning rate 為 $1e-3$ 和 epoch=60 的情況下，發現 loss 小於 1。但有一個缺點，因為 rnn 不能平行運算，導致每一次都會算很久。


```
(Epoch 49 / 60) loss: 0.2676 time per epoch: 12.7s
(Epoch 50 / 60) loss: 0.2211 time per epoch: 12.7s
(Epoch 51 / 60) loss: 0.1910 time per epoch: 12.7s
(Epoch 52 / 60) loss: 0.1681 time per epoch: 12.8s
(Epoch 53 / 60) loss: 0.1517 time per epoch: 12.7s
(Epoch 54 / 60) loss: 0.1403 time per epoch: 12.7s
(Epoch 55 / 60) loss: 0.1321 time per epoch: 12.8s
(Epoch 56 / 60) loss: 0.1243 time per epoch: 12.7s
(Epoch 57 / 60) loss: 0.1177 time per epoch: 12.7s
(Epoch 58 / 60) loss: 0.1119 time per epoch: 12.7s
(Epoch 59 / 60) loss: 0.1067 time per epoch: 12.8s
```

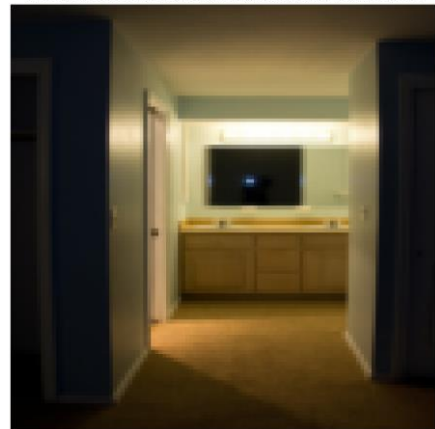


6. captions 出來的結果

[train] RNN Generated: <START> a group of people stand in a <UNK> room area with a frisee <END>
GT: <START> this is a street with a clock in the middle and people on a bench <END>



[train] RNN Generated: <START> a small bathroom with a sink and <UNK> the mirror <END>
GT: <START> an <UNK> to a kitchen with a microwave at the center <END>



[train] RNN Generated: <START> a city street <UNK> machine is parked by fire hydrant <END>
GT: <START> a city street <UNK> machine is parked by fire hydrant <END>



[train] RNN Generated: <START> has couple a green and baby vase sitting on a horse <END>
GT: <START> a group of people riding motorcycles on a street <END>



[val] RNN Generated: <START> a pizza is a slice of pizza and a man is eating it <END>
GT: <START> a pizza is a slice of pizza and a man is eating it <END>



[val] RNN Generated: <START> a person is surfing on a wave <END>
GT: <START> a person is surfing on a wave <END>



[val] RNN Generated: <START> at the pizza is a lady putting pizza sauce on some <END>
GT: <START> there is a lady putting pizza sauce on some <END>



[val] RNN Generated: <START> a man with a tie and glasses is standing by a door <END>
GT: <START> a man wearing a tie and glasses is standing by a door <END>



7. 再來進入 lstm 的部分。計算 next_h 和 next_c 和下一個隱藏狀態和 loss 與預期輸出結果相比，發現差距小。

```
expected_next_h = torch.tensor(
    [
        [0.24635157, 0.28610883, 0.32240467, 0.35525807, 0.38474904],
        [0.49223563, 0.55611431, 0.61507696, 0.66844003, 0.71591811],
        [0.56735664, 0.66310127, 0.74419266, 0.80889665, 0.8582991],
    ],
    **to_double
)
expected_next_c = torch.tensor(
    [
        [0.32986176, 0.39145139, 0.451556, 0.51014116, 0.56717407],
        [0.66382255, 0.76674007, 0.87195994, 0.97902709, 1.08751345],
        [0.74192008, 0.90592151, 1.07717006, 1.25120233, 1.42395676],
    ],
    **to_double
)

print("next_h error: ", rel_error(expected_next_h, next_h))
print("next_c error: ", rel_error(expected_next_c, next_c))

next_h error: 2.606541143878533e-09
next_c error: 1.7376745523804369e-09

hm = model(x, h0)

expected_hm = torch.tensor(
    [
        [
            [0.01764008, 0.01823233, 0.01882671, 0.01942332],
            [0.11287491, 0.12146228, 0.13018446, 0.13902939],
            [0.31358768, 0.33338627, 0.35304453, 0.37250975],
        ],
        [
            [0.45767879, 0.4761092, 0.4936887, 0.51041945],
            [0.6704845, 0.69350089, 0.71486014, 0.7346449],
            [0.81733511, 0.83677871, 0.85403753, 0.86935314],
        ],
    ],
    **to_double
)

print("hm error: ", rel_error(expected_hm, hm))

hm error: 2.668523515654886e-09
```

```

loss = model(images.to(DEVICE), captions.to(DEVICE))
expected_loss = torch.tensor(146.3161468505)

print("loss: ", loss.item())
print("expected loss: ", expected_loss.item())
print("difference: ", rel_error(loss, expected_loss))

For input images in NCHW format, shape (2, 3, 224, 224)
Shape of output c5 features: torch.Size([2, 400, 7, 7])
loss: 146.316162109375
expected loss: 146.31614685058594
difference: 5.214321112077035e-08

```

8. 使用 train_captioner 來訓練 CaptioningRNN，cell type 為 lstm，且只有使用 50 個 example 來訓練。在 learning rate 為 $1e-2$ 的情況下，發現 loss 小於 4。

```

(Epoch 69 / 80) loss: 3.5036 time per epoch: 1.1s
(Epoch 70 / 80) loss: 3.4996 time per epoch: 1.1s
(Epoch 71 / 80) loss: 3.4948 time per epoch: 1.0s
(Epoch 72 / 80) loss: 3.4920 time per epoch: 0.8s
(Epoch 73 / 80) loss: 3.4929 time per epoch: 0.8s
(Epoch 74 / 80) loss: 3.4946 time per epoch: 0.8s
(Epoch 75 / 80) loss: 3.4930 time per epoch: 0.8s
(Epoch 76 / 80) loss: 3.4890 time per epoch: 0.9s
(Epoch 77 / 80) loss: 3.4869 time per epoch: 0.9s
(Epoch 78 / 80) loss: 3.4874 time per epoch: 0.9s
(Epoch 79 / 80) loss: 3.4877 time per epoch: 0.9s

```

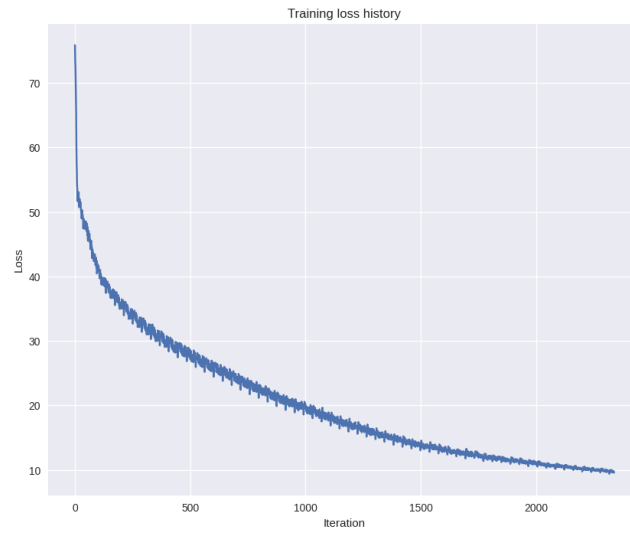


9. 使用 train_captioner 來訓練 CaptioningRNN，cell type 為 lstm。在 learning rate 為 $1e-3$ 和 epoch=60 的情況下，發現 loss 可以一直下降。

```

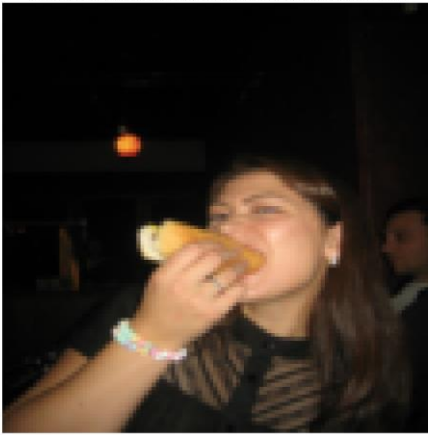
(Epoch 49 / 60) loss: 11.0403 time per epoch: 13.3s
(Epoch 50 / 60) loss: 10.8568 time per epoch: 13.3s
(Epoch 51 / 60) loss: 10.7346 time per epoch: 13.3s
(Epoch 52 / 60) loss: 10.5620 time per epoch: 13.3s
(Epoch 53 / 60) loss: 10.4227 time per epoch: 13.3s
(Epoch 54 / 60) loss: 10.2826 time per epoch: 13.3s
(Epoch 55 / 60) loss: 10.1743 time per epoch: 13.3s
(Epoch 56 / 60) loss: 9.9719 time per epoch: 13.3s
(Epoch 57 / 60) loss: 9.8601 time per epoch: 13.3s
(Epoch 58 / 60) loss: 9.7169 time per epoch: 13.3s
(Epoch 59 / 60) loss: 9.5804 time per epoch: 13.3s

```

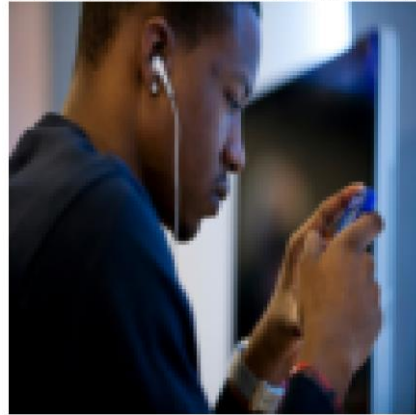


10. captions 出來的結果

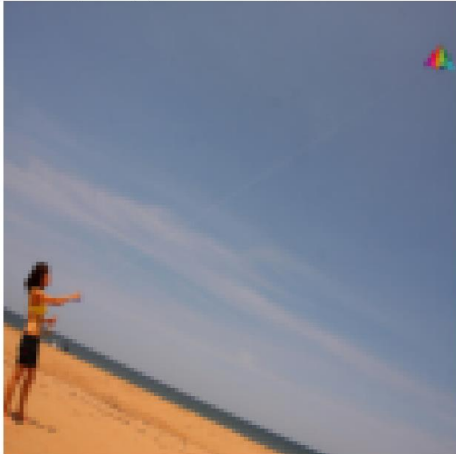
[train] LSTM Generated: <START> a man in a suit holding something in his hand <END>
GT: <START> a person eating a hot dog with a dark background <END>



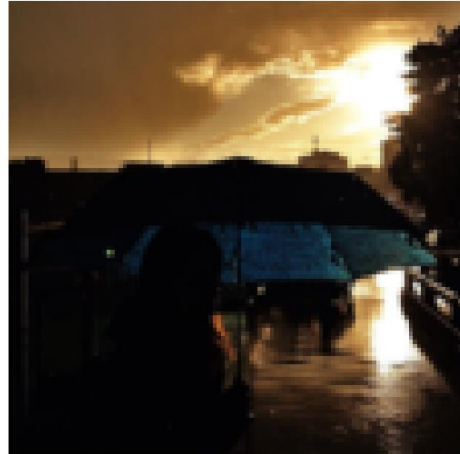
[train] LSTM Generated: <START> a man in a striped shirt cutting a cake <END>
GT: <START> a guy looking at a small <UNK> and wearing <UNK> phones <END>



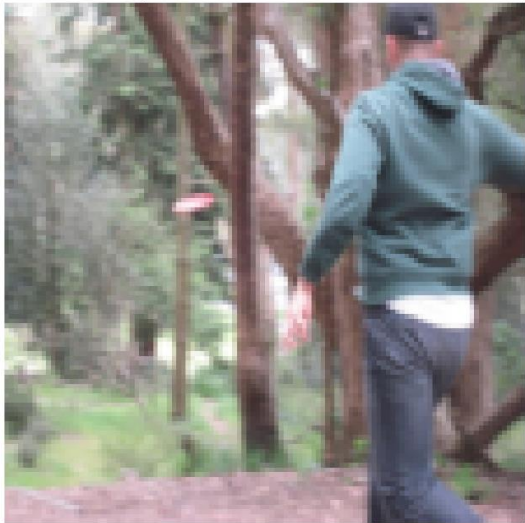
[train] LSTM Generated: <START> a man is riding a wave in the ocean <END>
GT: <START> a girl wearing a yellow shirt and black pants on the beach with a kite <END>



[train] LSTM Generated: <START> a man in a suit holding something in his hand <END>
GT: <START> a woman standing under an umbrella in a <UNK> city <END>



[val] LSTM Generated: <START> a man is walking on a beach with a surfboard <END>
GT: <START> a person <UNK> in the <UNK> next to some trees <END>



[val] LSTM Generated: <START> a man is riding a wave in the ocean <END>
GT: <START> the giraffes are eating from the many tall trees <END>



[val] LSTM Generated: <START> a man in a suit holding something in his hand <END>
GT: <START> a living room and dining room with <UNK> <UNK> <END>



[val] LSTM Generated: <START> a couple of blue street signs sitting on the side of a road <END>
GT: <START> a man is <UNK> in the water on a sunny day <END>

