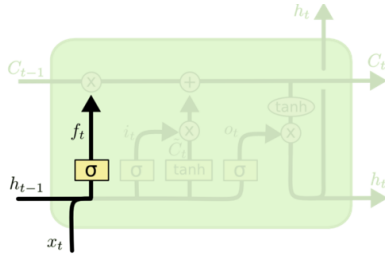
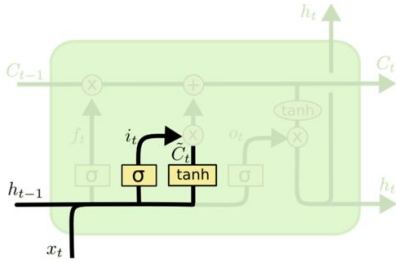


## Assignment 3

### Part 1 Differentiate LSTM

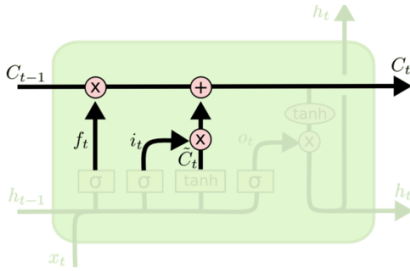


$$f_t = \sigma(W_f \cdot [h_{t-1}, x_t] + b_f)$$

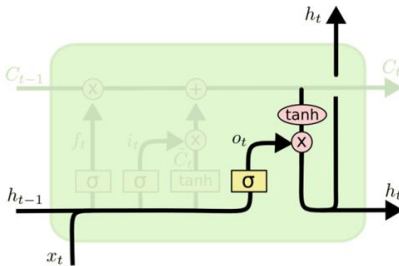


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$



$$o_t = \sigma(W_o [h_{t-1}, x_t] + b_o)$$

$$h_t = o_t * \tanh(C_t)$$

#### Lemma:

1.  $\sigma(z) = y = \frac{1}{1+e^{-z}}$
2.  $\sigma'(z) = y(1-y)$
3.  $\tanh(z) = y = \frac{e^z - e^{-z}}{e^z + e^{-z}}$
4.  $\tanh'(z) = 1 - y^2$
5.  $\text{diag}[a]X = a * X$ , where  $X$  is a matrix, and  $a$  is the vector composed of all the diag elements of  $\text{diag}[a]$
6.  $a^T \text{diag}[b] = a * b$ , where  $a$  is a column vector and  $b$  is the vector composed of all the diag elements of  $\text{diag}[b]$

Define the loss as  $E$ ,  $\delta_t \stackrel{\text{def}}{=} \frac{\partial E}{\partial h_t}$

$$\mathbf{net}_{o,t} = W_o[h_{t-1}, x_t] + b_o = W_{oh}h_{t-1} + W_{ox}x_t + b_o$$

$$\delta_{o,t} \stackrel{\text{def}}{=} \frac{\partial E}{\partial net_{o,t}}$$

$$\frac{\partial h_t}{\partial W_o} = \frac{\partial h_t}{\partial o_t} \frac{\partial o_t}{\partial net_{o,t}} \frac{\partial net_{o,t}}{\partial W_o} = \text{diag}[\tanh(c_t)] \text{diag}[o_t * (1 - o_t)] [h_{t-1}, x_t]$$

$$\begin{aligned}
\frac{\partial h_t}{\partial b_f} &= \frac{\partial h_t}{\partial f_t} \frac{\partial f_t}{\partial net_{f,t}} \frac{\partial net_{f,t}}{\partial b_f} = \text{diag}[o_t * (1 - \tanh(c_t)^2)] \text{diag}[c_{t-1}] \text{diag}[f_t * (1 - f_t)] \\
\frac{\partial h_t}{\partial b_i} &= \frac{\partial h_t}{\partial i_t} \frac{\partial i_t}{\partial net_{i,t}} \frac{\partial net_{i,t}}{\partial b_i} = \text{diag}[o_t * (1 - \tanh(c_t)^2)] \text{diag}[\bar{c}_t] \text{diag}[i_t * (1 - i_t)] \\
\frac{\partial h_t}{\partial b_c} &= \frac{\partial h_t}{\partial \bar{c}_t} \frac{\partial \bar{c}_t}{\partial net_{\bar{c},t}} \frac{\partial net_{\bar{c},t}}{\partial b_c} = \text{diag}[o_t * (1 - \tanh(c_t)^2)] \text{diag}[i_t] \text{diag}[1 - \bar{c}_t^2] \\
\frac{\partial h_t}{\partial b_o} &= \frac{\partial h_t}{\partial o_t} \frac{\partial o_t}{\partial net_{o,t}} \frac{\partial net_{o,t}}{\partial b_o} = \text{diag}[\tanh(c_t)] \text{diag}[o_t * (1 - o_t)]
\end{aligned}$$

To compute BPTT, first we need to compute  $\delta_{t-1}^T = \frac{\partial E}{\partial h_{t-1}} = \frac{\partial E}{\partial h_t} \frac{\partial h_t}{\partial h_{t-1}} = \delta_t^T \frac{\partial h_t}{\partial h_{t-1}}$

$$\begin{aligned}
\delta_t^T \frac{\partial h_t}{\partial h_{t-1}} &= \delta_t^T \frac{\partial h_t}{\partial o_t} \frac{\partial o_t}{\partial net_{o,t}} \frac{\partial net_{o,t}}{\partial h_{t-1}} + \delta_t^T \frac{\partial h_t}{\partial c_t} \frac{\partial c_t}{\partial f_t} \frac{\partial f_t}{\partial net_{f,t}} \frac{\partial net_{f,t}}{\partial h_{t-1}} + \delta_t^T \frac{\partial h_t}{\partial c_t} \frac{\partial c_t}{\partial i_t} \frac{\partial i_t}{\partial net_{i,t}} \frac{\partial net_{i,t}}{\partial h_{t-1}} + \delta_t^T \frac{\partial h_t}{\partial c_t} \frac{\partial c_t}{\partial \bar{c}_t} \frac{\partial \bar{c}_t}{\partial net_{\bar{c},t}} \frac{\partial net_{\bar{c},t}}{\partial h_{t-1}} \\
&= \delta_{o,t}^T \frac{\partial net_{o,t}}{\partial h_{t-1}} + \delta_{f,t}^T \frac{\partial net_{f,t}}{\partial h_{t-1}} + \delta_{i,t}^T \frac{\partial net_{i,t}}{\partial h_{t-1}} + \delta_{\bar{c},t}^T \frac{\partial net_{\bar{c},t}}{\partial h_{t-1}} = \delta_{o,t}^T W_{oh} + \delta_{f,t}^T W_{fh} + \delta_{i,t}^T W_{ih} + \delta_{\bar{c},t}^T W_{ch}
\end{aligned}$$

So finally we get:

$$\begin{aligned}
\delta_{t-1}^T &= \delta_t^T * [\tanh(c_t) * o_t * (1 - o_t) W_{oh} + o_t * (1 - \tanh(c_t)^2 * c_{t-1} * f_t * (1 - f_t) W_{fh} + o_t \\
&\quad * (1 - \tanh(c_t)^2 * \bar{c}_t * i_t * (1 - i_t) W_{ih} + o_t * (1 - \tanh(c_t)^2 * i_t * (1 - \bar{c}_t^2)]
\end{aligned}$$

At time t,

$$\frac{\partial E}{\partial W_{oh,t}} = \delta_{o,t} h_{t-1}^T, \frac{\partial E}{\partial W_{fh,t}} = \delta_{f,t} h_{t-1}^T, \frac{\partial E}{\partial W_{ih,t}} = \delta_{i,t} h_{t-1}^T, \frac{\partial E}{\partial W_{ch,t}} = \delta_{\bar{c},t} h_{t-1}^T$$

Add all the gradients in each time step:

$$\begin{aligned}
\frac{\partial E}{\partial W_{oh}} &= \sum_{j=1}^t \delta_{o,j} h_{j-1}^T, \frac{\partial E}{\partial W_{fh}} = \sum_{j=1}^t \delta_{f,j} h_{j-1}^T, \frac{\partial E}{\partial W_{ih}} = \sum_{j=1}^t \delta_{i,j} h_{j-1}^T, \frac{\partial E}{\partial W_{ch}} = \sum_{j=1}^t \delta_{\bar{c},j} h_{j-1}^T \\
\frac{\partial E}{\partial W_{ox}} &= \delta_{o,t} x_t^T, \frac{\partial E}{\partial W_{fx}} = \delta_{f,t} x_t^T, \frac{\partial E}{\partial W_{ix}} = \delta_{i,t} x_t^T, \frac{\partial E}{\partial W_{cx}} = \delta_{\bar{c},t} x_t^T \\
\frac{\partial E}{\partial b_o} &= \sum_{j=1}^t \delta_{o,j}, \frac{\partial E}{\partial b_i} = \sum_{j=1}^t \delta_{i,j}, \frac{\partial E}{\partial b_f} = \sum_{j=1}^t \delta_{f,j}, \frac{\partial E}{\partial b_c} = \sum_{j=1}^t \delta_{\bar{c},j}
\end{aligned}$$

## Part 2 Autograd Training of LSTM

**Dataset:** 《全唐诗》， fetched from <https://github.com/chinese-poetry/chinese-poetry>, 10k poems with length ranged from 10 to 63. The dataset is split into training dataset and development dataset (80%:20%) (Please include the json file folder of this github repository in the project as source of data.)

**Dictionary:** all words in training set, EOS, OOV, START, PAD.

### Hyperparameters and training setting:

- Vocabulary size,  $|V|$ : 5031
- Batch size, bs: 128
- Sentence length, sl: 64
- Hidden size, hs: 256
- Input size, is: 256

### Data preparation:

1. Fetch the data, data cleaning.

Download the dataset from github, decode the json file and extract the poems with proper length(10-63). The data cleaning procedures include deleting all the content in (), [], {}, 《》 and all the numbers or strange symbols.

2. Construct vocabulary dictionary, and save it as file for future reference.

```
1. def wordDic(self, data): # 增加padding
2.     words = sorted(set([character for sent in data for character in sent]))
3.     words.extend(['<EOS>', '<OOV>', '<START>', '<PAD>'])
4.     word_to_idx = {word: idx for idx, word in enumerate(words)}
5.     idx_to_word = {idx: word for idx, word in enumerate(words)}
6.     # save dict
7.     with open('wordDic', 'wb') as f:
8.         pickle.dump(word_to_idx, f)
9.     return word_to_idx, idx_to_word
```

3. Map sentences to sequences of integers.

Append EOS to the end of the poem. Add paddings to extend the poem to standard length. Mask unknown words in development data with OOV.

### LSTM model:

```
1. class LSTM(nn.Module):
2.     def __init__(self, vocab_size, embedding_dim, hidden_dim):
3.         super(LSTM, self).__init__()
4.         # embedding layer
5.         self.embeddings = nn.Embedding(vocab_size, embedding_dim)
6.         self.hidden_dim = hidden_dim
7.         # weights for inputs
8.         self.weight_ih = Parameter(torch.Tensor(embedding_dim, hidden_dim*4))
9.         # weights for hidden states
10.        self.weight_hh = Parameter(torch.Tensor(hidden_dim, hidden_dim*4))
11.        self.bias = Parameter(torch.Tensor(hidden_dim*4))
12.        self.init_weights()
13.        # linear transformation
14.        self.linear = nn.Linear(self.hidden_dim, vocab_size)
15.        # logsoftmax, for loss calculation
16.        self.softmax = nn.LogSoftmax()
17.
18.    def init_weights(self):
19.        for p in self.parameters():
20.            if p.data.ndimension() >= 2:
21.                nn.init.xavier_uniform_(p.data)
22.            else:
23.                nn.init.zeros_(p.data)
24.
```

```

25.     def forward(self, input, hidden=None):
26.         # input(batch, seq_length, input_size)
27.         length = input.size()[1]
28.         batch_size = input.size()[0]
29.         embeds = self.embeddings(input).view((batch_size, length, -1))
30.         bs, seq_sz, _ = embeds.size()
31.         hidden_seq = []
32.         if hidden is None:
33.             h_t, c_t = (torch.zeros(self.hidden_dim).to(embeds.device), torch.zeros(self.hidden_dim).
to(embeds.device))
34.         else:
35.             h_t, c_t = hidden
36.
37.         for t in range(seq_sz):
38.             x_t = embeds[:, t, :]
39.             gates = x_t @ self.weight_ih + h_t @ self.weight_hh + self.bias
40.             i_t, f_t, g_t, o_t = gates.chunk(4, 1)
41.             i_t = torch.sigmoid(i_t)
42.             f_t = torch.sigmoid(f_t)
43.             g_t = torch.tanh(g_t)
44.             o_t = torch.sigmoid(o_t)
45.             c_t = f_t * c_t + i_t * g_t
46.             h_t = o_t * torch.tanh(c_t)
47.             hidden_seq.append(h_t.unsqueeze(0))
48.
49.         hidden_seq = torch.cat(hidden_seq, dim=0)
50.         # resize to shape(batch, seq_len, input_size)
51.         hidden_seq = hidden_seq.transpose(0, 1)
52.
53.         output = F.relu(self.linear(hidden_seq))
54.         if output.size()[1] == 1: # seq_length=1, 会导致softmax出问题 [1,1,V] 要先改成[1,V]
55.             output = output[:, 0, :]
56.             output = self.softmax(output)
57.             output = output.view(1, 1, -1)
58.             return output, (h_t, c_t)
59.         output = self.softmax(output)
60.         return output, (h_t, c_t)

```

## Train & Test:

```

1.     D = data_processing.Dataset_batched_padding()
2.     src = './model_padding_10k/' # 存放模型
3.
4.     model = LSTM(len(D.word_to_ix), 256, 256)
5.     optimizer = optim.Adam(model.parameters(), lr=0.01, weight_decay=0.0001)
6.     # together with the softmax layer in the model, compute CrossEntropy loss
7.     criterion = nn.NLLLoss()
8.
9.
10.    epochNum = 50
11.    TRAINSIZE = len(D.train)
12.    TESTSIZE = len(D.develop)
13.    batch = 128
14.
15.    def compute_perplexity(output, target): # 计算每个句子的困惑度
16.        N = len(target)
17.        return 2**(-1/N*sum([output[j][target[j]] for j in range(len(target))]))
18.
19.

```

```

20. def test():
21.     model.eval()
22.     loss = 0
23.     perplexity = 0
24.
25.     for batchIndex in range(int(TESTSIZE / batch)):
26.         if batchIndex*batch >= TESTSIZE:
27.             continue
28.         t, o = D.makeForOneBatch(D.develop[batchIndex*batch:min((batchIndex+1)*batch, TESTSIZE)])
29.         output, hidden = model(t)
30.         for i in range(o.size()[0]):
31.             loss += criterion(output[i, :, :], o[i, :])
32.             perplexity += compute_perplexity(output[i, :, :], o[i, :])
33.
34.     loss = loss / TESTSIZE
35.     perplexity = perplexity / TESTSIZE
36.     print("=====", loss.item())
37.     print("+++++", perplexity.item())
38.     return loss.item()
39.
40.
41. print("start training")
42. for epoch in range(epochNum):
43.     model.train()
44.     for batchIndex in range(int(TRAINSIZE / batch)):
45.         if batchIndex*batch >= TRAINSIZE:
46.             continue
47.         model.zero_grad()
48.         t, o = D.makeForOneBatch(D.train[batchIndex*batch:min((batchIndex+1)*batch, TRAINSIZE)])
49.         output, hidden = model(t)
50.         loss = 0
51.         for i in range(o.size()[0]):
52.             loss += criterion(output[i, :, :], o[i, :])
53.         loss = loss / o.size()[0]
54.         loss.backward()
55.         print(epoch, loss.item())
56.         optimizer.step()
57.
58.     test_loss = test()
59.     torch.save(model, src+'poetry-gen-epoch%d-loss%f.pt' % (epoch+1, test_loss))

```

**Poetry generation:**

```

1. def generate(startWord='<START>', max_length=95):
2.     model = torch.load('./model_padding_10k/poetry-gen-epoch38-loss4.066066.pt')
3.     with open('wordDic', 'rb') as f:
4.         word_to_ix = p.load(f)
5.         ix_to_word = invert_dict(word_to_ix)
6.         input = make_index(startWord, word_to_ix)
7.         poem = ""
8.         if startWord != '<START>':
9.             poem = startWord
10.        hidden = None
11.        for i in range(max_length):
12.            output, hidden = model(input, hidden)
13.            topvs, topis = output[0].data.topk(2)
14.            topi = topis[0][0].item()
15.            topi1 = topis[0][1].item()
16.            w = ix_to_word[topi]
17.            w1 = ix_to_word[topi1]
18.            if w == '<EOS>': # generate new sentence by random sampling
19.                w = ix_to_word[numpy.random.randint(low=2, high=len(ix_to_word)-4)]
20.                hidden = None
21.            if w != '<EOS>':
22.                poem += w
23.                input = make_index(w, word_to_ix)
24.        return poem

```

### Initialization: why we shouldn't initialize all the parameters to zero?

When we are training a model, in fact we are trying to update and find the best weights so that the model can find the hidden relationship between input and output. **If all the weights are initialized to zero, then every neuron in the same layer will have the same output value. When we do back propagation, all the gradients are the same, too. Therefore, the neurons cannot learn different features, and the model is likely to fail.**

### Proper way to initialize parameters:

A proper initialization should avoid two things: excessive saturation of activation functions(the gradients will not propagate well), and overly linear units. **Xavier initialization** ensures that the variances of the gradients on the weights is the same for all layers. Assume that the size of layer i is






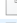
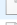







$n_i$ , the initialization can be in the form:  $W \sim U[-\frac{\sqrt{6}}{\sqrt{n_j+n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j+n_{j+1}}}]$

### Final perplexity and some of the models:

```

48 4.403109550476074
48 4.105821132659912
48 4.1353888511657715
48 4.054264545440674
48 4.121681213378906
48 4.101590633392334
48 4.087003707885742
48 4.108267784118652
48 4.042253494262695
48 4.250035285949707
48 4.16510534286499
48 4.160760402679443
48 4.147387504577637
48 4.411069393157959
48 4.185972213745117
48 4.450275421142578
48 4.17116117477417
48 4.194272041320801
48 4.112797737121582
=====loss: 4.074352264404297
+++++perplexity: 18.452558517456055

```

名称	修改日期	类型
 poetry-gen-epoch48-loss4.074352.pt	2019/5/21 20:57	PT 文件
 poetry-gen-epoch47-loss4.056911.pt	2019/5/20 11:22	PT 文件
 poetry-gen-epoch46-loss4.068842.pt	2019/5/20 10:47	PT 文件
 poetry-gen-epoch45-loss4.072767.pt	2019/5/20 10:13	PT 文件
 poetry-gen-epoch44-loss4.066272.pt	2019/5/20 9:38	PT 文件
 poetry-gen-epoch43-loss4.063836.pt	2019/5/20 9:03	PT 文件
 poetry-gen-epoch42-loss4.066301.pt	2019/5/20 0:24	PT 文件
 poetry-gen-epoch41-loss4.062868.pt	2019/5/20 0:02	PT 文件
 poetry-gen-epoch40-loss4.060951.pt	2019/5/19 23:40	PT 文件
 poetry-gen-epoch39-loss4.059323.pt	2019/5/19 23:18	PT 文件
 poetry-gen-epoch38-loss4.066066.pt	2019/5/19 22:56	PT 文件
 poetry-gen-epoch37-loss4.062064.pt	2019/5/19 22:33	PT 文件
 poetry-gen-epoch36-loss4.061742.pt	2019/5/19 22:11	PT 文件
 poetry-gen-epoch35-loss4.067415.pt	2019/5/19 21:48	PT 文件

The poems generated by the model: (at different epochs)

```
C:\Myprogram\python3.6.7\python.exe "C:/Myprogram/Pycharm Projects/Autograd_lstm/sample_batched.py"
C:\Myprogram\Pycharm Projects\Autograd_lstm\model_batched.py:107: UserWarning: Implicit dimension choice for log_softmax has been deprecated. Change the call to include dim=X as an argument.
    output = self.softmax(output)
日暮雲霞雪，天明月下山。凍衣袈落日，不得此時時。湛露寒雲下，風風動雪聲。焚天涯遠望，天子是誰家。蒲葉落花落，江南風柳風。補風吹雪雪，江上見天風。捫君不見此，不得不知君。殘風吹落葉，江水不飛天。
紅燭滿堂花，風風動白雲。童子子家家，天涯不得來。又逢君子不，不得此時時。舊國風吹落，天風風動來。貴君不見此，不得不知君。胤君不見此，不得不知君。攝君王子，不得其天。陵路遠山水，青山不可飛。掉首
山川水上天，天子不能聞。鴻君不見此，不得不知君。救君不見此，不得不知君。鴻風吹落葉，落葉不成風。堤上天高望，天涯不得歸。真子不知此，我心不知君。捫君此去，不得其然。攻君不見此，不得不知君。誠知
夜月明明月，天明不見人。驕君不見此，不得不知君。豎風吹落葉，江水不飛天。形城南望天，天子不能聞。印衣袈上，天子不成。攤子風吹，不復大夫。寒花落花落，不得風風吹。懿君王子，不得天子。伯馬行人不，
湖上天涯遠，天涯此日來。隋風吹雪雪，江上見天風。遑君不見此，不得不能聞。嚶有詩人事，不知此時時。慨然此去，誰復不知。貴君不見此，不得不知君。伴有詩人事，誰家不得時。吸日暮雲去，天涯不可飛。輻心
海上天涯遠，天涯此日來。蕙子風吹，天子其然。撫君不見此，不得不知君。蒙詩書自得，我士不能言。證道風流落，天明月下山。畝中有所思，我有此時時。盱君不見此，不得不知君。振衣袈上，天子不成。衰聲落日
月明天子，天子是天。迪君不見此，不得不知君。批君不見此，不得不知君。馳道風流水，天涯不見人。竭風吹落葉，江水不飛天。面下天高望，天明月下來。敵風吹落葉，江水不飛天。獫子風吹，天子其然。孰知此時
春風吹柳色，江上見春風。袖中有雪下，不得一聲聲。坦衣袈上，天子不成。組不知君不得，誰家不得知君。倪君不見此，不得不知君。源中有所見，我有此時時。圖吉自有所，我士不知君。曠子風吹，天子其然。晦日
花落花開樹，風風柳上風。元君不見此，不得不能聞。宙上天高望，天涯不得歸。狄下天高望，天明月下來。虞馬行人不見君，不知天子不能聞。沌風吹落葉，江水不飛天。黠落花開處，風風柳上風。弘者不知此，天子
秋風吹落葉，江水不飛天。聲聲聲落，天子不成。建章天子，天子是天。透簾聞翠露，天子不能聞。沃陵風景已，天子月明明。鳴馬行人不見君，不知天子不能聞。枚君不見此，不得不知君。葡君不見此，不得不知君。

Process finished with exit code 0
```

```
C:\Myprogram\python3.6.7\python.exe "C:/Myprogram/Pycharm Projects/Autograd_lstm/sample_batched.py"
C:\Myprogram\Pycharm Projects\Autograd_lstm\model_batched.py:107: UserWarning: Implicit dimension choice for log_softmax has been deprecated. Change the call to include dim=X as an argument.
    output = self.softmax(output)
日暮雲霞雪，風風動雪聲。背風吹落葉，江水不悠悠。捫子風吹露，風風動此聲。湊風吹落葉，江水不悠悠。阮家家子，無事不知。崑風吹落葉，江水不悠悠。萌馬行人不，此時不可憐。曠子
紅燭滿堂中，風風動玉衣。柚花落葉落花落，不得風聲入玉樓。故中有所思，我我不知此。冊君不見，我心不知。惻君王子，不得其然。故人聞有所，我事不知君。搢得知君子，不知此去年。愚子不知，我其其然。貴君
山川風吹雪，江上月明明。遑啼花落落，不得風風吹。暗風吹落葉，江水不悠悠。彰君不見，我心不知。縫衣袈上，天子不成。惑子風吹，不見天子。昇天子子，我我其然。荷風吹落葉，江水不悠悠。錢子風吹露，風風
夜月明明月，風風動玉衣。舛馬行人不，此時不可憐。小人心自有，我此不知君。曠君不見，我心不知。植花開處處，風柳半天風。插風吹落葉，江水不悠悠。儼雪中天下，天風動雪飛。漢陽城下有餘風，天子無聲不得
湖上天涯遠，天涯不可憐。輅下天涯遠，天南天子歸。臆君不見，我心不知。伐風吹落葉，江水不悠悠。孫子不知，我其其然。鵠風吹落葉，江水不悠悠。衛風吹落葉，江水不悠悠。澣衣袈上雪，天子不知人。壁衣袈落
海上風吹雪，江南風柳風。激君不知，我心不知。掉首望鄉處，不知此去時。撫君不見，我心不知。諧君不見，我心不知。助有餘年少，不知此去來。庵下有餘風，天子不得來。亘子風吹露，風風動此聲。毗酒不知，我
月明明月照，江上天風風。倏下天涯遠，天涯不可憐。袋風吹落葉，江水不悠悠。匪君不見，我心不知。誕風吹雪雪，江上雪天秋。癰聲聲遠，天子不知。西南望鄉路，不得一聲聲。騰風吹落葉，江水不悠悠。幟城南望
春風吹柳色，江水綠花開。淺不知君不得，不知此去不知。穉花落葉落花落，不得風聲入玉樓。剗風吹落葉，江水不悠悠。頤子風吹露，風風動此聲。厭君王子，不得其然。譏風吹落葉，江水不悠悠。痿玉琴聲，不見天
花落花開盡，風風柳色飛。龍風吹落葉，江水不悠悠。俊人不見，我心不知。勉君不知，我心不知。靚風吹落葉，江水不悠悠。挂衣袈上，風吹雪天。着子君無事，我心不得知。珠簾下翠露，天子不知人。珣君不見，我
秋風吹落葉，江水不悠悠。楠風吹落葉，江水不悠悠。晴君不見，我心不知。連年年少少，不得不能歸。涓風吹落葉，江水不悠悠。謨下天涯遠，天涯不可憐。復聞聲聲，不見天子。閆門外人不，此日不能歸。鵲聲不可

Process finished with exit code 0
```

日暮雲霞雪，天明月下山。俄有餘花落，不知此去來。

紅燭下落落，不知風露清。柳色青青草，風流不可飛。隕水無窮事，何年不得歸。漢陽城下有，我有此中來。

山川風吹雪，江上月明明。莫風吹落葉，落葉不蕭條。

夜雨深山月暮，無人不待風。煖風吹落露枝，不知不得離。

夜雨聲聲落露，不知不得離人。無事不知此，不知無所思。

湖上有餘，此心無所。荻花落盡日，風景不可飛。慘別離心不，此時無所情。

海上天涯遠，天涯此去來。船下有餘風，天明不可見。

月明朝日暮，天子不知來。荊州城下雪，不見青青天。蜚風吹落日，不得一聲聲。莫問君王子，誰家不見人。

月明明月下，不見此時來。疎柳花枝落日斜，不知此日望南來。

numpy implementation:

```
1. class Sigmoid(object):
2.     def forward(self, weighted_input):
3.         return 1.0 / (1.0 + np.exp(-weighted_input))
4.
5.     def backward(self, output):
6.         return output * (1 - output)
7.
8.
9. class Tanh(object):
10.     def forward(self, weighted_input):
11.         return 2.0 / (1.0 + np.exp(-2 * weighted_input)) - 1.0
12.
13.     def backward(self, output):
14.         return 1 - output * output
15.
16.
17. class Identity(object):
18.     def forward(self, weighted_input):
19.         return weighted_input
20.
21.     def backward(self, output):
22.         return 1
```



```

1. class LstmLayer(object):
2.     def __init__(self, input_size, hidden_size,
3.                   learning_rate):
4.         self.input_size = input_size
5.         self.hidden_size = hidden_size
6.         self.learning_rate = learning_rate
7.         # 门的激活函数
8.         self.gate_activator = Sigmoid()
9.         # 输出的激活函数
10.        self.output_activator = Tanh()
11.        # 当前时刻初始化为t0
12.        self.times = 0
13.        # 各个时刻的单元状态向量c
14.        self.c_list = self.init_state_list()
15.        # 各个时刻的输出向量h
16.        self.h_list = self.init_state_list()
17.        # 各个时刻的遗忘门f
18.        self.f_list = self.init_state_list()
19.        # 各个时刻的输入门i
20.        self.i_list = self.init_state_list()
21.        # 各个时刻的输出门o
22.        self.o_list = self.init_state_list()
23.        # 各个时刻的即时状态c~
24.        self.ct_list = self.init_state_list()
25.        # 遗忘门权重矩阵Wfh, Wfx, 偏置项bf
26.        self.Wfh, self.Wfx, self.bf = (
27.            self.init_weight_mat())
28.        # 输入门权重矩阵Wih, Wix, 偏置项bi
29.        self.Wih, self.Wix, self.bi = (
30.            self.init_weight_mat())
31.        # 输出门权重矩阵Woh, Wox, 偏置项bo
32.        self.Woh, self.Wox, self.bo = (
33.            self.init_weight_mat())
34.        # 单元状态权重矩阵Wch, Wcx, 偏置项bc
35.        self.Wch, self.Wcx, self.bc = (
36.            self.init_weight_mat())
37.
38.    def init_state_list(self):
39.        # 初始化保存状态的向量
40.        state_vec_list = []
41.        state_vec_list.append(np.zeros(
42.            (self.hidden_size, 1)))
43.        return state_vec_list
44.
45.    def init_weight_mat(self):
46.        # 初始化权重矩阵
47.        Wh = np.random.uniform(-1e-4, 1e-4,
48.                                (self.hidden_size, self.hidden_size))
49.        Wx = np.random.uniform(-1e-4, 1e-4,
50.                                (self.hidden_size, self.input_size))
51.        b = np.zeros((self.hidden_size, 1))
52.        return Wh, Wx, b
53.

```

```

54.     def forward(self, x):
55.         # 前向计算
56.         self.times += 1
57.         # 遗忘门
58.         fg = self.calc_gate(x, self.Wfx, self.Wfh,
59.                             self.bf, self.gate_activator)
60.         self.f_list.append(fg)
61.         # 输入门
62.         ig = self.calc_gate(x, self.Wix, self.Wih,
63.                             self.bi, self.gate_activator)
64.         self.i_list.append(ig)
65.         # 输出门
66.         og = self.calc_gate(x, self.Wox, self.Woh,
67.                             self.bo, self.gate_activator)
68.         self.o_list.append(og)
69.         # 即时状态
70.         ct = self.calc_gate(x, self.Wcx, self.Wch,
71.                             self.bc, self.output_activator)
72.         self.ct_list.append(ct)
73.         # 单元状态
74.         c = fg * self.c_list[self.times - 1] + ig * ct
75.         self.c_list.append(c)
76.         # 输出
77.         h = og * self.output_activator.forward(c)
78.         self.h_list.append(h)
79.
80.     def calc_gate(self, x, Wx, Wh, b, activator):
81.         # 计算门
82.         h = self.h_list[self.times - 1] # 上次的LSTM输出
83.         net = np.dot(Wh, h) + np.dot(Wx, x) + b
84.         gate = activator.forward(net)
85.         return gate
86.
87.     def backward(self, x, delta_h, activator):
88.         # 反向传播, 实现LSTM训练算法
89.         self.calc_delta(delta_h, activator)
90.         self.calc_gradient(x)
91.

```

```

92. def update(self):
93.     # 按照梯度下降, 更新权重
94.     self.Wfh -= self.learning_rate * self.Whf_grad
95.     self.Wfx -= self.learning_rate * self.Whx_grad
96.     self.bf -= self.learning_rate * self.bf_grad
97.     self.Wih -= self.learning_rate * self.Whi_grad
98.     self.Wix -= self.learning_rate * self.Whi_grad
99.     self.bi -= self.learning_rate * self.bi_grad
100.    self.Woh -= self.learning_rate * self.Wof_grad
101.    self.Wox -= self.learning_rate * self.Wox_grad
102.    self.bo -= self.learning_rate * self.bo_grad
103.    self.Wch -= self.learning_rate * self.Wcf_grad
104.    self.Wcx -= self.learning_rate * self.Wcx_grad
105.    self.bc -= self.learning_rate * self.bc_grad
106.
107.    def calc_delta(self, delta_h, activator):
108.        # 初始化各个时刻的误差项
109.        self.delta_h_list = self.init_delta() # 输出误差项
110.        self.delta_o_list = self.init_delta() # 输出误差项
111.        self.delta_i_list = self.init_delta() # 输入误差项
112.        self.delta_f_list = self.init_delta() # 遗忘误差项
113.        self.delta_ct_list = self.init_delta() # 即时输出误差项
114.
115.        # 保存从上一层传递下来的当前时刻的误差项
116.        self.delta_h_list[-1] = delta_h
117.
118.        # 迭代计算每个时刻的误差项
119.        for k in range(self.times, 0, -1):
120.            self.calc_delta_k(k)
121.
122.    def init_delta(self):
123.        # 初始化误差项
124.        delta_list = []
125.        for i in range(self.times + 1):
126.            delta_list.append(np.zeros(
127.                (self.hidden_size, 1)))
128.        return delta_list
129.

```

```

130.     def calc_delta_k(self, k):
131.         # 根据k时刻的delta_h, 计算k时刻的delta_f、delta_i、delta_o、delta_ct, 以及k-1时刻的delta_h
132.         # 获得k时刻前向计算的值
133.         ig = self.i_list[k]
134.         og = self.o_list[k]
135.         fg = self.f_list[k]
136.         ct = self.ct_list[k]
137.         c = self.c_list[k]
138.         c_prev = self.c_list[k - 1]
139.         tanh_c = self.output_activator.forward(c)
140.         delta_k = self.delta_h_list[k]
141.
142.         delta_o = (delta_k * tanh_c *
143.                    self.gate_activator.backward(og))
144.         delta_f = (delta_k * og *
145.                    (1 - tanh_c * tanh_c) * c_prev *
146.                    self.gate_activator.backward(fg))
147.         delta_i = (delta_k * og *
148.                    (1 - tanh_c * tanh_c) * ct *
149.                    self.gate_activator.backward(ig))
150.         delta_ct = (delta_k * og *
151.                     (1 - tanh_c * tanh_c) * ig *
152.                     self.output_activator.backward(ct))
153.         delta_h_prev = (
154.             np.dot(delta_o.transpose(), self.Woh) +
155.             np.dot(delta_i.transpose(), self.Wih) +
156.             np.dot(delta_f.transpose(), self.Wfh) +
157.             np.dot(delta_ct.transpose(), self.Wch)
158.         ).transpose()
159.
160.         # 保存全部delta值
161.         self.delta_h_list[k - 1] = delta_h_prev
162.         self.delta_f_list[k] = delta_f
163.         self.delta_i_list[k] = delta_i
164.         self.delta_o_list[k] = delta_o
165.         self.delta_ct_list[k] = delta_ct
166.

```

```

167. def calc_gradient(self, x):
168.     # 初始化遗忘门权重梯度矩阵和偏置项
169.     self.Wfh_grad, self.Wfx_grad, self.bf_grad = (
170.         self.init_weight_gradient_mat())
171.     # 初始化输入门权重梯度矩阵和偏置项
172.     self.Wih_grad, self.Wix_grad, self.bi_grad = (
173.         self.init_weight_gradient_mat())
174.     # 初始化输出门权重梯度矩阵和偏置项
175.     self.Woh_grad, self.Wox_grad, self.bo_grad = (
176.         self.init_weight_gradient_mat())
177.     # 初始化单元状态权重梯度矩阵和偏置项
178.     self.Wch_grad, self.Wcx_grad, self.bc_grad = (
179.         self.init_weight_gradient_mat())
180.
181.     # 计算对上一次输出h的权重梯度
182.     for t in range(self.times, 0, -1):
183.         # 计算各个时刻的梯度
184.         (Wfh_grad, bf_grad,
185.          Wih_grad, bi_grad,
186.          Woh_grad, bo_grad,
187.          Wch_grad, bc_grad) = (
188.             self.calc_gradient_t(t))
189.         # 实际梯度是各时刻梯度之和
190.         self.Wfh_grad += Wfh_grad
191.         self.bf_grad += bf_grad
192.         self.Wih_grad += Wih_grad
193.         self.bi_grad += bi_grad
194.         self.Woh_grad += Woh_grad
195.         self.bo_grad += bo_grad
196.         self.Wch_grad += Wch_grad
197.         self.bc_grad += bc_grad
198.
199.     # 计算对本次输入x的权重梯度
200.     xt = x.transpose()
201.     self.Wfx_grad = np.dot(self.delta_f_list[-1], xt)
202.     self.Wix_grad = np.dot(self.delta_i_list[-1], xt)
203.     self.Wox_grad = np.dot(self.delta_o_list[-1], xt)
204.     self.Wcx_grad = np.dot(self.delta_ct_list[-1], xt)
205.

```

```

206.     def init_weight_gradient_mat(self):
207.         # 初始化权重矩阵
208.         Wh_grad = np.zeros((self.hidden_size,
209.                             self.hidden_size))
210.         Wx_grad = np.zeros((self.hidden_size,
211.                             self.input_size))
212.         b_grad = np.zeros((self.hidden_size, 1))
213.         return Wh_grad, Wx_grad, b_grad
214.
215.     def calc_gradient_t(self, t):
216.         # 计算每个时刻t权重的梯度
217.         h_prev = self.h_list[t - 1].transpose()
218.         Wfh_grad = np.dot(self.delta_f_list[t], h_prev)
219.         bf_grad = self.delta_f_list[t]
220.         Wih_grad = np.dot(self.delta_i_list[t], h_prev)
221.         bi_grad = self.delta_f_list[t]
222.         Woh_grad = np.dot(self.delta_o_list[t], h_prev)
223.         bo_grad = self.delta_f_list[t]
224.         Wch_grad = np.dot(self.delta_ct_list[t], h_prev)
225.         bc_grad = self.delta_ct_list[t]
226.         return Wfh_grad, bf_grad, Wih_grad, bi_grad, \
227.             Woh_grad, bo_grad, Wch_grad, bc_grad
228.
229.     def reset_state(self):
230.         # 当前时刻初始化为t0
231.         self.times = 0
232.         # 各个时刻的单元状态向量c
233.         self.c_list = self.init_state_list()
234.         # 各个时刻的输出向量h
235.         self.h_list = self.init_state_list()
236.         # 各个时刻的遗忘门f
237.         self.f_list = self.init_state_list()
238.         # 各个时刻的输入门i
239.         self.i_list = self.init_state_list()
240.         # 各个时刻的输出门o
241.         self.o_list = self.init_state_list()
242.         # 各个时刻的即时状态c~
243.         self.ct_list = self.init_state_list()

```

### Optimization: Comparing two optimizers

- Adagrad: Adagrad is an Adaptive Gradient Method that implies different adaptive learning rates for each feature. Hence it is more intuitive for especially sparse problems and it is likely to find more discriminative features. Although we provide an initial learning rate, Adagrad tunes it regarding the history of the gradients for each feature dimension. The formulation for Adagrad is as below:

$$W_{t+1} = W_t - \frac{\eta_0}{\sqrt{\sum_{t'=1}^t (g_{t',i})^2 + \epsilon}} \odot g_{t,i}$$

The upper formula states that, for each feature dimension, learning rate  $\eta_0$  is divided by all the squared root gradient history. **The main strengths of Adagrad is that we don't need to adjust the learning rate by ourselves. However, as the iteration times become larger, the learning rate will become smaller and finally approaches zero.**

- Adam: Adam is derived from adaptive moment estimation, and is very beneficial when it's used on non-convex optimization problems. Adam computes individual adaptive learning rates for different parameters from estimates of first and second moments of the gradients. The formulation of Adam is as below:

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\widehat{m}_t = \frac{m_t}{1 - \beta_1^t}, \widehat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$W_{t+1} = W_t - \frac{\eta}{\sqrt{\widehat{v}_t + \epsilon}} \widehat{m}_t$$

$m_t, v_t$  is first and second moments of the gradients.

Specifically, the algorithm calculates an exponential moving average of the gradient and the squared gradient, and the parameters beta1 and beta2 control the decay rates of these moving averages.

The initial value of the moving averages and beta1 and beta2 values close to 1.0 (recommended) result in a bias of moment estimates towards zero. This bias is overcome by first calculating the biased estimates before then calculating bias-corrected estimates.

Optimizer	Strength	Weakness
Adagrad	In the context of sparse data, Adagrad can better use the information from sparse gradients, so it can converge quicker than standard SGD.	As iteration times become larger, the denominator becomes larger and larger, finally the learning rate will become too small and the gradients cannot be updated efficiently.
Adam	Adam is easy to implement and efficient in calculation. The hyperparameters have great interpretability. It can also change the learning rate automatically, so it's especially suitable for unstable target function.	

Therefore, Adam might be the best overall choice.

**These two optimizers will influence our gradient calculation in that, the gradients update are not merely determined by initial learning rate, but also by other variables like moment or history weights. So when we are calculating the gradients, we have to store additional parameters.**