

# 深度学习作业——神经网络语言模型

Clark

## 1 实验目的

本实验旨在基于 PyTorch 深度学习框架，构建并训练一个基于 LSTM（长短期记忆网络）的神经网络语言模型。通过本实验，主要达到以下目的：

1. 深入理解循环神经网络（RNN）及其变体 LSTM 在处理序列数据中的工作机制；
2. 掌握语言模型的基本原理，即如何根据历史上下文预测下一个词的概率分布；
3. 学习并实践文本数据的预处理流程，包括分词、构建词典、数据批处理（Batchify）等；
4. 掌握在训练过程中应用正则化技术（如 Locked Dropout、权重绑定）和优化策略（如梯度裁剪、学习率退火）以提升模型性能。

## 2 概述

语言模型是自然语言处理（NLP）中的核心任务之一，其目标是计算一个单词序列的概率。本实验使用经典的 Penn Treebank（PTB）数据集，构建一个多层 LSTM 网络来训练语言模型。实验涵盖了从数据下载、预处理、模型构建、训练循环到最终评估的全过程。通过计算困惑度（Perplexity, PPL）来衡量模型的性能，并分析不同训练策略对模型收敛速度和泛化能力的影响。

## 3 实验要求

1. 基于 Python 语言和 PyTorch 深度学习框架，完成 PTB 数据集的下载与读取；
2. 构建基于 LSTM 的语言模型，包含 Embedding 层、LSTM 层和全连接解码层；
3. 实现模型的训练与评估流程，使用困惑度（PPL）作为评价指标；
4. 应用 Dropout、梯度裁剪（Gradient Clipping）等技术防止过拟合和梯度爆炸；
5. 绘制训练损失和验证集困惑度的变化曲线，分析实验结果。

## 4 实验原理

### 4.1 语言模型

语言模型的目标是计算一个词序列  $w_1, w_2, \dots, w_T$  的联合概率  $P(w_1, w_2, \dots, w_T)$ 。根据链式法则，该概率可以分解为条件概率的乘积：

$$P(w_1, w_2, \dots, w_T) = \prod_{t=1}^T P(w_t | w_1, w_2, \dots, w_{t-1}). \quad (1)$$

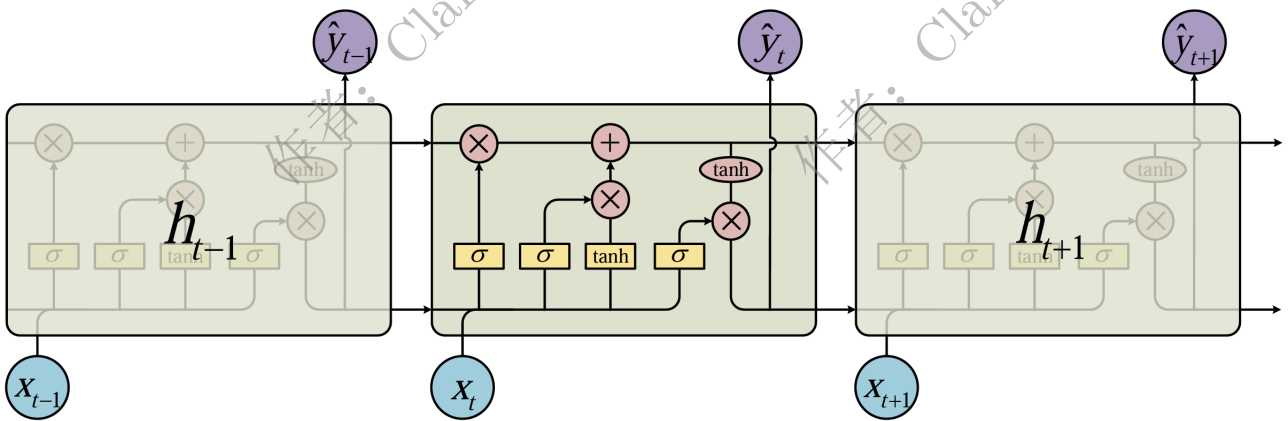
在神经网络语言模型中，我们通常使用 RNN 或 LSTM 来编码历史信息  $w_1, \dots, w_{t-1}$  为一个隐藏状态  $h_{t-1}$ ，然后预测下一个词  $w_t$  的概率分布。

### 4.2 长短期记忆网络 (LSTM)

LSTM 通过引入门控机制解决了传统 RNN 存在的梯度消失和梯度爆炸问题。一个标准的 LSTM 单元包含遗忘门  $f_t$ 、输入门  $i_t$  和输出门  $o_t$ 。其状态更新公式如下：

$$\begin{aligned} i_t &= \sigma(W_{ii}x_t + b_{ii} + W_{hi}h_{t-1} + b_{hi}), \\ f_t &= \sigma(W_{if}x_t + b_{if} + W_{hf}h_{t-1} + b_{hf}), \\ g_t &= \tanh(W_{ig}x_t + b_{ig} + W_{hg}h_{t-1} + b_{hg}), \\ o_t &= \sigma(W_{io}x_t + b_{io} + W_{ho}h_{t-1} + b_{ho}), \\ c_t &= f_t \odot c_{t-1} + i_t \odot g_t, \\ h_t &= o_t \odot \tanh(c_t), \end{aligned} \quad (2)$$

其中， $x_t$  是当前时刻的输入， $h_{t-1}$  是上一时刻的隐藏状态， $c_t$  是细胞状态， $\sigma$  是 Sigmoid 激活函数， $\odot$  表示逐元素乘法。最后我们给出 LSTM 的网络结构图。



### 4.3 困惑度 (Perplexity)

困惑度是评价语言模型好坏的常用指标，定义为交叉熵损失的指数形式：

$$\text{PPL} = e^{\text{CrossEntropyLoss}} = e^{-\frac{1}{N} \sum_{i=1}^N \ln P(w_i | w_{<i})}, \quad (3)$$

PPL 越低，表示模型对下一个词的预测越准确。

### 4.4 词嵌入与概率生成

在模型输入端，词嵌入层 (Embedding Layer) 将离散的单词索引映射为稠密的实值向量。数学上，这等价于一个矩阵乘法：

$$\mathbf{e}_t = \mathbf{W}_{\text{emb}} \cdot \mathbf{x}_t, \quad (4)$$

其中  $\mathbf{x}_t \in \{0, 1\}^{|V|}$  是单词的 One-hot 向量， $\mathbf{W}_{\text{emb}} \in \mathbb{R}^{d \times |V|}$  是嵌入矩阵。在模型输出端，解码层将 LSTM 的隐藏状态  $\mathbf{h}_t$  映射回词表空间，并通过 Softmax 函数生成概率分布：

$$P(w_{t+1} = k | w_{1:t}) = \frac{e^{\mathbf{w}_k^T \mathbf{h}_t + b_k}}{\sum_{j=1}^{|V|} e^{\mathbf{w}_j^T \mathbf{h}_t + b_j}}, \quad (5)$$

其中  $\mathbf{w}_k$  是解码矩阵  $\mathbf{W}_{\text{dec}}$  的第  $k$  行。本实验中采用了权重绑定 (Weight Tying) 策略，即约束  $\mathbf{W}_{\text{dec}} = \mathbf{W}_{\text{emb}}^T$ ，这不仅减少了参数量，还赋予了嵌入向量更明确的语义解释。

### 4.5 随时间反向传播 (BPTT)

循环神经网络的训练依赖于随时间反向传播算法 (Backpropagation Through Time, BPTT)。BPTT 本质上是将 RNN 在时间维度上展开，将其视为一个层数等于时间步长  $T$  的深层前馈网络。总损失函数  $L$  是每个时间步损失  $L_t$  的总和：

$$L = \sum_{t=1}^T L_t(\theta), \quad (6)$$

其中  $\theta$  是网络参数。参数的梯度是所有时间步梯度的累加：

$$\frac{\partial L}{\partial \theta} = \sum_{t=1}^T \frac{\partial L_t}{\partial \theta}. \quad (7)$$

在普通 RNN 中，梯度计算涉及连乘项  $\prod_{k=t}^T \frac{\partial h_k}{\partial h_{k-1}}$ ，这容易导致梯度消失或爆炸。LSTM 通过细胞状态  $c_t$  的加性更新规则 (即  $c_t = f_t \odot c_{t-1} + \dots$ ) 引入了一条“恒等映射”路径，使得梯度能够更长时间地维持其幅度，从而有效缓解了梯度消失问题。

## 5 实验数据集及工具

### 5.1 数据预处理

实验使用 PTB 数据集。首先定义 “Dictionary” 类来构建词表，将单词映射为索引；“Corpus” 类负责读取训练、验证和测试数据并进行 Tokenization。

```

1 class Dictionary(object):
2     def __init__(self):
3 self.word2idx = {}      self.idx2word = []
4
5     def add_word(self, word):
6         if word not in self.word2idx:
7             self.idx2word.append(word)
8             self.word2idx[word] = len(self.idx2word) - 1
9         return self.word2idx[word]
10
11     def __len__(self):
12         return len(self.idx2word)
13
14 class Corpus(object):
15     def __init__(self, path):
16         self.dictionary = Dictionary()
17         self.train = self.tokenize(os.path.join(path, 'ptb.train.txt'))
18         self.valid = self.tokenize(os.path.join(path, 'ptb.valid.txt'))
19         self.test = self.tokenize(os.path.join(path, 'ptb.test.txt'))
20
21     def tokenize(self, path):
22         """Tokenizes a text file."""
23         assert os.path.exists(path)
24         # Add words to the dictionary
25         with open(path, 'r', encoding='utf-8') as f:
26             tokens = 0
27             for line in f:
28                 words = line.split() + ['<eos>']
29                 tokens += len(words)
30                 for word in words:
31                     self.dictionary.add_word(word)
32
33         # Tokenize file content
34         with open(path, 'r', encoding='utf-8') as f:
35             ids = torch.LongTensor(tokens)
36             token = 0
37             for line in f:
38                 words = line.split() + ['<eos>']
39                 for word in words:
40                     ids[token] = self.dictionary.word2idx[word]
41                     token += 1
42         return ids

```

为了利用 GPU 进行并行计算，我们需要将数据整理成“batch\_size”列的矩阵。“batchify”函数将长序

列数据切割并堆叠:

```
1 def batchify(data, bsz):
2     # Work out how cleanly we can divide the dataset into bsz parts.
3     nbatch = data.size(0) // bsz
4     # Trim off any extra elements that wouldn't cleanly fit (remainders).
5     data = data.narrow(0, 0, nbatch * bsz)
6     # Evenly divide the data across the bsz batches.
7     data = data.view(bsz, -1).t().contiguous()
8     return data.to(device)
```

## 5.2 模型构建

本实验构建了一个基于多层 LSTM (Long Short-Term Memory) 的神经网络语言模型。模型主要由三部分组成: 编码器 (Encoder)、LSTM 核心层和解码器 (Decoder)。其中, LSTM 层负责捕捉序列数据中的长距离依赖关系。为了防止过拟合, 模型特别采用了“LockedDropout”机制, 这是一种专门针对 RNN 设计的 Dropout 变体。与在每个时间步独立采样掩码的标准 Dropout 不同, LockedDropout 在整个时间步序列中对相同的神经元应用相同的 Dropout 掩码, 从而在破坏神经元间共适应性的同时, 最大程度地保留了时间维度上的信息流。

```
1 class LockedDropout(nn.Module):
2     def __init__(self):
3         super(LockedDropout, self).__init__()
4
5     def forward(self, x, dropout=0.5):
6         if not self.training or not dropout:
7             return x
8         m = x.data.new(1, x.size(1), x.size(2)).bernoulli_(1 - dropout)
9         mask = m.div_(1 - dropout)
10        mask = mask.expand_as(x)
11        return x * mask
12
13 class LSTModel(nn.Module):
14     def __init__(self, vocab_size, embed_size, hidden_size, num_layers, dropout=0.5):
15         super(LSTModel, self).__init__()
16         self.locked_drop = LockedDropout()
17         self.drop = nn.Dropout(dropout)
18         self.encoder = nn.Embedding(vocab_size, embed_size)
19         self.lstm = nn.LSTM(embed_size, hidden_size, num_layers, dropout=dropout)
20         self.decoder = nn.Linear(hidden_size, vocab_size)
21
22         # Tie weights if dimensions match
23         if embed_size == hidden_size:
24             self.decoder.weight = self.encoder.weight
25
```

```

26         self.init_weights()
27 self.hidden_size = hidden_size          self.num_layers = num_layers
28
29     def init_weights(self):
30         initrange = 0.1
31         self.encoder.weight.data.uniform_(-initrange, initrange)
32         self.decoder.bias.data.zero_()
33         self.decoder.weight.data.uniform_(-initrange, initrange)
34
35         # Initialize LSTM forget gate bias to 1
36         for names in self.lstm._all_weights:
37             for name in filter(lambda n: "bias" in n, names):
38                 bias = getattr(self.lstm, name)
39                 n = bias.size(0)
40                 start, end = n // 4, n // 2
41                 bias.data[start:end].fill_(1.)
42
43     def forward(self, input, hidden):
44         emb = self.encoder(input)
45         emb = self.locked_drop(emb, dropout=self.drop.p)
46
47         output, hidden = self.lstm(emb, hidden)
48
49         output = self.locked_drop(output, dropout=self.drop.p)
50         decoded = self.decoder(output.view(output.size(0)*output.size(1), output.size(2))
51                                 )
52         return decoded.view(output.size(0), output.size(1), decoded.size(1)), hidden
53
54     def init_hidden(self, bsz):
55         weight = next(self.parameters())
56         return (weight.new_zeros(self.num_layers, bsz, self.hidden_size),
57                 weight.new_zeros(self.num_layers, bsz, self.hidden_size))

```

### 5.3 模型训练与评估

训练过程中使用了梯度裁剪（Gradient Clipping）来防止梯度爆炸，并实现了学习率退火策略。

```

1 def train(model, train_data, criterion, optimizer, epoch):
2     model.train()
3     total_loss = 0.
4     start_time = time.time()
5     hidden = model.init_hidden(BATCH_SIZE)
6
7     batch_losses = []
8

```

```

9     for batch, i in enumerate(range(0, train_data.size(0) - 1, SEQ_LEN)):
10         data, targets = get_batch(train_data, i)           # Starting each batch, we detach the
                                                             hidden state from how it was previously produced
11         hidden = repackage_hidden(hidden)
12
13         model.zero_grad()
14         output, hidden = model(data, hidden)
15         loss = criterion(output.view(-1, vocab_size), targets)
16         loss.backward()
17
18         # `clip_grad_norm` helps prevent the exploding gradient problem in RNNs / LSTMs.
19         torch.nn.utils.clip_grad_norm_(model.parameters(), GRAD_CLIP)
20         optimizer.step()
21
22         total_loss += loss.item()
23
24     if batch % 200 == 0 and batch > 0:
25         cur_loss = total_loss / 200
26         elapsed = time.time() - start_time
27         print('|_epoch_|{:3d}|_|_|{:5d}|/{:5d}|_batches_|_|lr_|{:02.5f}|_|_|ms/batch_|{:5.2f}|_|_|'
28               '|_loss_|{:5.2f}|_|_|app|_|{:8.2f}|'.format(
29                 epoch, batch, len(train_data) // SEQ_LEN, LEARNING_RATE,
30                 elapsed * 1000 / 200, cur_loss, math.exp(cur_loss)))
31         batch_losses.append(cur_loss)
32         total_loss = 0
33         start_time = time.time()
34
35 return batch_losses

```

## 6 实验结果与分析

### 6.1 训练过程分析

最后实验的结果如下图所示，我们经过 55 个 Epoch 的训练，模型在训练集上的 ppl 得分为 79.22，损失函数为 4.37；在测试集上的 ppl 得分为 75.88，损失函数为 4.33。

```

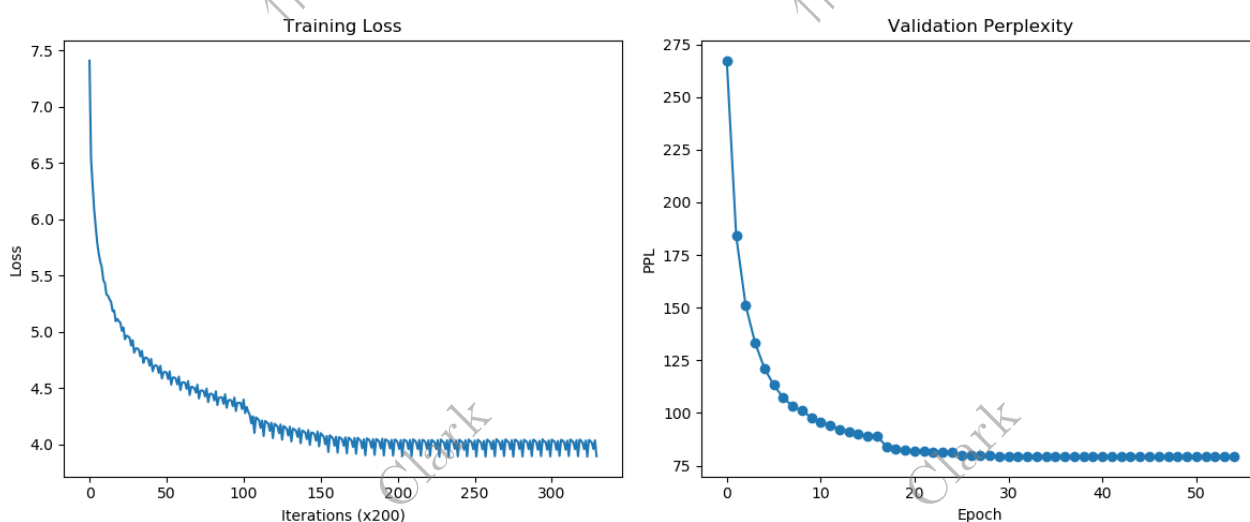
-----
| epoch 55 | 200/ 1327 batches | lr 20.00000 | ms/batch 17.19 | loss 4.04 | ppl 56.66
| epoch 55 | 400/ 1327 batches | lr 20.00000 | ms/batch 17.12 | loss 4.03 | ppl 56.16
| epoch 55 | 600/ 1327 batches | lr 20.00000 | ms/batch 17.16 | loss 4.01 | ppl 55.24
| epoch 55 | 800/ 1327 batches | lr 20.00000 | ms/batch 17.19 | loss 3.96 | ppl 52.29
| epoch 55 | 1000/ 1327 batches | lr 20.00000 | ms/batch 17.19 | loss 4.03 | ppl 56.50
| epoch 55 | 1200/ 1327 batches | lr 20.00000 | ms/batch 17.21 | loss 3.89 | ppl 49.01
-----

| end of epoch 55 | time: 23.27s | valid loss 4.37 | valid ppl 79.22
-----

====
| End of training | test loss 4.33 | test ppl 75.88
=====
Plot saved to result.png
chenpeng@amax:~/test$

```

实验设置了 55 个 Epoch，Batch Size 为 20，序列长度为 35。从生成的训练曲线可以看出，随着训练的进行，Training Loss 呈现稳步下降趋势，表明模型正在有效地拟合训练数据。同时，Validation Perplexity（验证集困惑度）也随之降低，说明模型的泛化能力在不断增强。



## 6.2 关键改进措施

- **Locked Dropout:** 与标准的 Dropout 不同，Locked Dropout 在每个时间步采样相同的 Dropout 掩码。这对于 RNN 尤为重要，因为它避免了破坏循环连接中的长期依赖关系，从而允许模型更有效地学习长序列信息。
- **权重绑定 (Weight Tying):** 在“LSTMMModel”中，当“embed\_size”等于“hidden\_size”时，我们将 Embedding 层的权重与输出解码层 (Linear) 的权重共享。这不仅显著减少了模型的参数量，降低了过拟合风险，还利用了 Embedding 层学到的语义信息来辅助解码。



- 梯度裁剪 (Gradient Clipping): 在 “train” 函数中使用了 “`torch.nn.utils.clip_grad_norm_`”, 设定阈值为 0.25。这一措施有效防止了 RNN 训练中常见的梯度爆炸问题, 确保了训练过程的数值稳定性。
- 学习率退火 (Learning Rate Annealing): 在主循环中, 如果验证集 Loss 没有下降, 则将学习率除以 4。这种动态调整策略使得模型在训练初期能快速收敛, 而在后期能进行更精细的参数搜索, 从而达到更优的局部极小值。

## 7 实验总结

本次实验成功构建并训练了一个基于 LSTM 的神经网络语言模型。通过对 PTB 数据集的分析和处理, 验证了 LSTM 在序列建模任务中的有效性。实验结果表明, 结合 Locked Dropout、权重绑定和梯度裁剪等优化技术, 模型能够达到较低的困惑度 (Perplexity), 具备良好的预测能力。

实验过程中, 我深入理解了语言模型的评价指标 PPL 的物理含义, 以及如何通过 Batchify 技术高效利用 GPU 资源。此外, 通过对比不同正则化手段的作用, 进一步掌握了深度学习模型调优的方法论。这些经验为后续进行更复杂的自然语言处理任务 (如机器翻译、文本生成) 奠定了坚实的基础。