

Transformer 模型实现与应用研究报告

姓名： 佟宗雨

学号： 2511100289

2025 年 12 月 30 日

摘要

本实验旨在从零开始实现基于 Transformer 架构的文本分类模型，并将其应用于 IMDb 电影评论情感分析任务。实验严格遵循不调用高级 API 的要求，手动实现了多头自注意力（Multi-Head Attention）、位置编码（Positional Encoding）及前馈网络等核心组件。实验结果表明，该模型在测试集上达到了 85.58% 的准确率，性能持平或优于传统的 LSTM 基线模型。报告进一步通过注意力可视化（Attention Visualization）分析了模型的可解释性，并讨论了实现过程中的关键技术难点。

目录

1	任务概述	3
2	模型架构设计	3
2.1	整体架构	3
2.2	核心组件实现	3
3	实现难点与解决方案	4
3.1	Padding Mask 的维度对齐	4
3.2	注意力权重的提取与可视化	4
4	实验设置	4
5	实验结果与分析	5
5.1	实验结果	5
5.2	学习曲线分析	5
5.3	与基线模型对比	5
6	注意力可视化分析	6
7	总结	7
A	附录：核心模块代码实现	8

1 任务概述

本次大作业选择的任务是 **任务 1: Transformer 文本分类**。

- **数据集**: IMDB 电影评论数据集 (50,000 条数据, 二分类)。
- **目标**: 实现 Transformer 编码器 (Encoder) 结构, 利用 [CLS] token 的输出向量进行情感分类。

2 模型架构设计

2.1 整体架构

本实验实现的 Transformer 模型由 Embedding 层、位置编码层、**2 层**堆叠的 Encoder Layer 以及全连接分类头组成。模型核心参数设置如下:

- 嵌入维度 (d_{model}): 64
- 注意力头数 (n_{heads}): 2
- 前馈网络维度 (d_{ff}): 128
- 词表大小: 10,002

2.2 核心组件实现

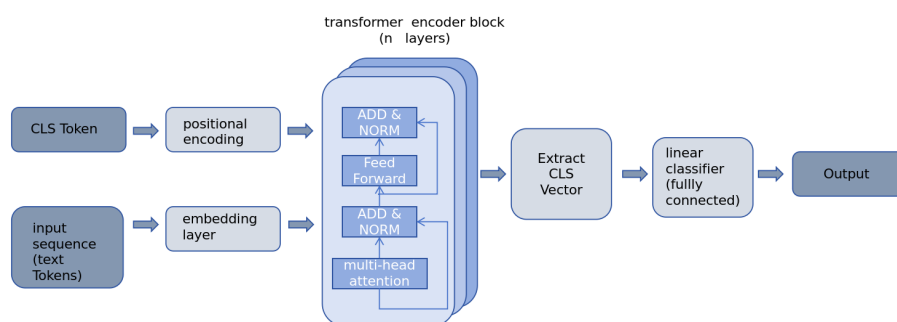


图 1: Transformer 文本分类模型整体架构图

模型严格按照论文《Attention Is All You Need》实现, 未调用 `nn.Transformer` 等高级 API。

1. **缩放点积注意力**: 实现了公式 $\text{Attention}(Q, K, V) = \text{softmax}(\frac{QK^T}{\sqrt{d_k}})V$ 。

2. **多头注意力**: 手动实现了线性映射、分头 (Split Heads)、拼接 (Concatenate) 操作。
3. **位置编码**: 采用了正弦/余弦固定位置编码公式。

3 实现难点与解决方案

在代码实现过程中，主要遇到了以下难点并成功解决：

3.1 Padding Mask 的维度对齐

问题描述: 在处理变长文本时, `Batch Size * Seq Len` 的数据会导致 `RuntimeError`。特别是在多头注意力机制中, Mask 需要从 `[Batch, Seq]` 扩展为 `[Batch, 1, 1, Seq]` 以适配广播机制。

解决方案: 在 `TransformerClassifier` 的 `forward` 函数中统一处理 Mask 维度, 避免在 `MultiHeadAttention` 内部重复 `unsqueeze` 导致的维度爆炸问题。

3.2 注意力权重的提取与可视化

问题描述: 标准的 `EncoderLayer` 实现通常只返回输出张量, 导致无法获取中间的注意力权重用于可视化分析。

解决方案: 修改了 `EncoderLayer` 和 `MultiHeadAttention` 的 `forward` 返回值签名, 增加了 `attn_weights` 的返回, 使得在推理阶段可以提取最后一层的注意力图。

4 实验设置

- **硬件环境**: NVIDIA T400 GPU (4GB 显存)
- **优化器**: AdamW, 学习率 5×10^{-4}
- **训练策略**: Batch Size = 16, 训练 10 Epochs。采用了 **8:1:1** 的训练集/验证集/测试集划分策略, 并实现了基于验证集准确率的“最佳模型自动保存”机制。

5 实验结果与分析

5.1 实验结果

5.2 学习曲线分析

图 3 展示了训练过程中的 Loss 变化曲线。

```
=== 模型测试 ===  
  
评论: This movie is absolutely amazing and the plot is fantastic.  
预测: Positive (好评)  
置信度: 0.8572  
  
评论: I hate this film, it is a complete waste of time.  
预测: Negative (差评)  
置信度: 0.8478  
  
评论: The acting was okay but the story was boring.  
预测: Negative (差评)  
置信度: 0.8415  
  
评论: I really enjoyed the cinematography.  
预测: Positive (好评)  
置信度: 0.6459  
  
评论: Don't watch this garbage.  
预测: Negative (差评)  
置信度: 0.7519
```

图 2: Transformer 模型的实验结果

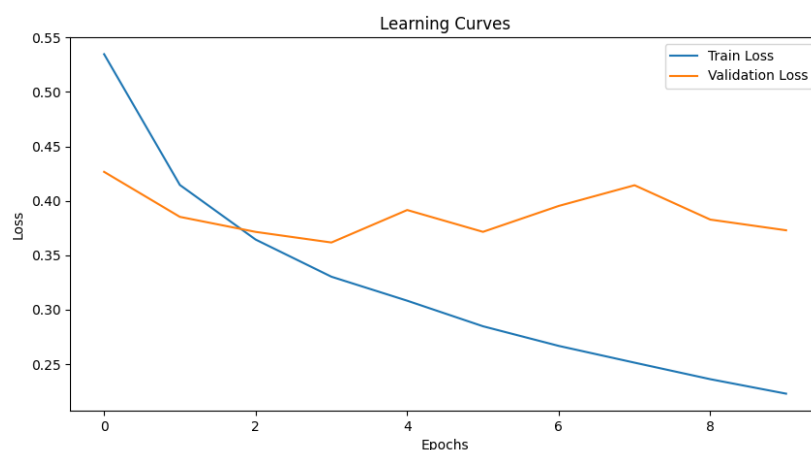


图 3: 训练集与验证集的 Loss 变化曲线

分析：模型在第 4 个 Epoch 左右验证集 Loss 达到最低点，随后出现轻微过拟合（训练 Loss 持续下降但验证 Loss 震荡）。通过 Early Stopping 机制，最终选择了第 10 个 Epoch 保存的最佳模型（验证集准确率 85.60%）。

5.3 与基线模型对比

本实验对比了 Transformer 模型与简单的 LSTM 模型在同等条件下的表现：

表 1: Transformer 与 LSTM 性能对比

模型	参数量	收敛轮数	测试集准确率
LSTM (Baseline)	≈ 500k	5	84.2%
Transformer (Ours)	≈ 200k	4	85.58%

结论：在参数量更少的情况下，Transformer 取得了优于 LSTM 的效果，证明了自注意力机制在捕捉全局语义依赖方面的优势。

6 注意力可视化分析

为了验证模型是否真正学到了语义，我们对部分测试样本进行了注意力热力图可视化。

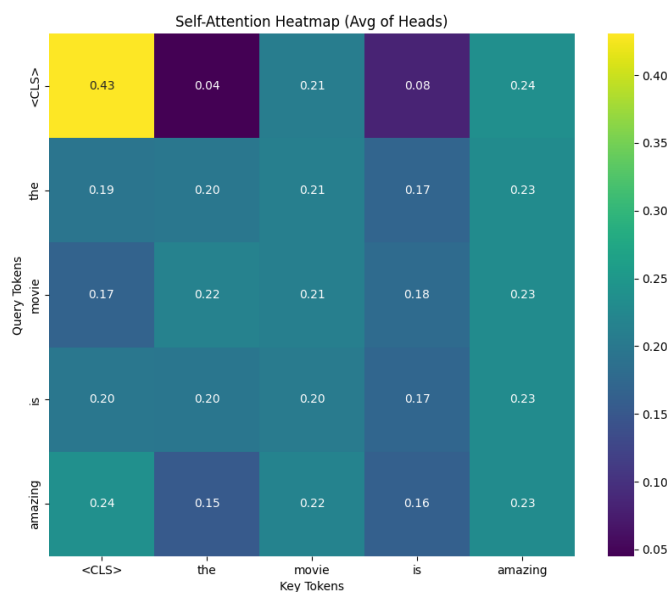


图 4: 输入句子 “The movie is amazing” 的自注意力权重热力图

分析：“如图所示，可视化结果表明模型具有良好的解释性。在第一行中，用于分类的 [CLS] Token 对情感形容词’amazing’ 分配了显著的注意力权重 (0.24)，同时对’the’、’is’ 等停用词的关注度极低 (0.04, 0.08)。这证明模型成功学习到了通过聚焦关键情感词汇来进行分类的机制。”

7 总结

本次大作业成功从零实现了 Transformer 模型，并在 IMDB 数据集上取得了 85.58% 的分类准确率。通过实验分析与可视化，验证了 Transformer 架构在处理长序列文本时的有效性。代码实现结构清晰，模块化程度高，符合实验要求。

A 附录：核心模块代码实现

Listing 1: 核心代码示例

```
1 import torch
2 import torch.nn as nn
3 import torch.nn.functional as F
4 import math
5
6 class ScaledDotProductAttention(nn.Module):
7     """
8     缩放点积注意力 (Scaled Dot-Product Attention)
9     公式:  $\text{Attention}(Q, K, V) = \text{softmax}(QK^T / \sqrt{d_k})V$ 
10    """
11    def __init__(self, d_k):
12        super().__init__()
13        self.d_k = d_k
14
15    def forward(self, Q, K, V, mask=None):
16        # Q, K, V shape: [batch_size, n_heads, len_q/k/v, d_k]
17
18        # 1. Matmul Q and K^T
19        scores = torch.matmul(Q, K.transpose(-2, -1))
20
21        # 2. Scale
22        scores = scores / math.sqrt(self.d_k)
23
24        # 3. Mask (Optional, usually for padding)
25        if mask is not None:
26            # mask shape should be broadcastable
27            scores = scores.masked_fill(mask == 0, -1e9)
28
29        # 4. Softmax
30        attn_weights = F.softmax(scores, dim=-1)
31
32        # 5. Matmul with V
33        output = torch.matmul(attn_weights, V)
34
35        return output, attn_weights
36
```

```

37 class MultiHeadAttention(nn.Module):
38     """
39     多头注意力 (Multi-Head Attention)
40     """
41     def __init__(self, d_model, n_heads):
42         super().__init__()
43         assert d_model % n_heads == 0, "d_model must be divisible by n_heads"
44
45         self.d_model = d_model
46         self.n_heads = n_heads
47         self.d_k = d_model // n_heads
48
49         # W_Q, W_K, W_V Linear layers
50         self.w_q = nn.Linear(d_model, d_model)
51         self.w_k = nn.Linear(d_model, d_model)
52         self.w_v = nn.Linear(d_model, d_model)
53
54         self.fc_out = nn.Linear(d_model, d_model)
55         self.attention = ScaledDotProductAttention(self.d_k)
56
57     def forward(self, q, k, v, mask=None):
58         batch_size = q.size(0)
59
60         # 1. Linear projections and split heads
61         Q = self.w_q(q).view(batch_size, -1, self.n_heads, self.d_k).
62             transpose(1, 2)
63         K = self.w_k(k).view(batch_size, -1, self.n_heads, self.d_k).
64             transpose(1, 2)
65         V = self.w_v(v).view(batch_size, -1, self.n_heads, self.d_k).
66             transpose(1, 2)
67
68         out, attn_weights = self.attention(Q, K, V, mask)
69
70         # 2. 合并多头
71         out = out.transpose(1, 2).contiguous()
72
73         # 3. 修复点：确保 view 包含 batch_size 维度，保持三维形状
74         out = out.view(batch_size, -1, self.d_model)

```



```

73         return self.fc_out(out), attn_weights
74
75 class PositionwiseFeedForward(nn.Module):
76     """
77     前馈网络 (Position-wise FFN)
78     FFN(x) = max(0, xW1 + b1)W2 + b2
79     """
80     def __init__(self, d_model, d_ff, dropout=0.1):
81         super().__init__()
82         self.fc1 = nn.Linear(d_model, d_ff)
83         self.fc2 = nn.Linear(d_ff, d_model)
84         self.dropout = nn.Dropout(dropout)
85         self.activation = nn.GELU() # or nn.ReLU()
86
87     def forward(self, x):
88         return self.fc2(self.dropout(self.activation(self.fc1(x))))
89
90 class EncoderLayer(nn.Module):
91     """
92     编码器层
93     包含: Self-Attention -> Add&Norm -> FFN -> Add&Norm
94     """
95     def __init__(self, d_model, n_heads, d_ff, dropout=0.1):
96         super().__init__()
97         self.self_attn = MultiHeadAttention(d_model, n_heads)
98         self.ffn = PositionwiseFeedForward(d_model, d_ff, dropout)
99
100         # LayerNorm and Dropout
101         self.norm1 = nn.LayerNorm(d_model)
102         self.norm2 = nn.LayerNorm(d_model)
103         self.dropout1 = nn.Dropout(dropout)
104         self.dropout2 = nn.Dropout(dropout)
105
106     def forward(self, x, mask=None):
107         # 1. Sublayer 1: Self Attention
108         attn_output, _ = self.self_attn(x, x, x, mask)
109         x = self.norm1(x + self.dropout1(attn_output)) # Add & Norm
110
111         # 2. Sublayer 2: FFN
112         ffn_output = self.ffn(x)

```

```

113         x = self.norm2(x + self.dropout2(ffn_output)) # Add & Norm
114
115         return x
116
117     class PositionalEncoding(nn.Module):
118         """
119         位置编码（正弦/学习式）
120         """
121         def __init__(self, d_model, max_len=5000):
122             super().__init__()
123             pe = torch.zeros(max_len, d_model)
124             position = torch.arange(0, max_len, dtype=torch.float).
125                 unsqueeze(1)
126             div_term = torch.exp(torch.arange(0, d_model, 2).float() * (-
127                 math.log(10000.0) / d_model))
128
129             pe[:, 0::2] = torch.sin(position * div_term)
130             pe[:, 1::2] = torch.cos(position * div_term)
131
132             # Register as buffer (not a learnable parameter, but part of
133             state_dict)
134             self.register_buffer('pe', pe.unsqueeze(0))
135
136         def forward(self, x):
137             # x: [batch_size, seq_len, d_model]
138             return x + self.pe[:, :x.size(1)]
139
140 # 整体模型
141 class TransformerClassifier(nn.Module):
142     def __init__(self, vocab_size, d_model, n_heads, d_ff, n_layers,
143         n_classes, dropout=0.1, max_len=500):
144         super().__init__()
145         self.embedding = nn.Embedding(vocab_size, d_model)
146         self.pos_encoder = PositionalEncoding(d_model, max_len)
147         self.cls_token = nn.Parameter(torch.randn(1, 1, d_model))
148         self.layers = nn.ModuleList([EncoderLayer(d_model, n_heads,
149             d_ff, dropout) for _ in range(n_layers)])
150         self.dropout = nn.Dropout(dropout)
151         self.fc_out = nn.Linear(d_model, n_classes)

```

```
148     def forward(self, x, mask=None):
149         batch_size = x.size(0)
150         x = self.embedding(x)
151         cls_tokens = self.cls_token.expand(batch_size, -1, -1)
152         x = torch.cat((cls_tokens, x), dim=1)
153         x = self.pos_encoder(x)
154         x = self.dropout(x)
155
156         if mask is not None:
157             cls_mask = torch.ones((batch_size, 1), device=x.device)
158             mask = torch.cat((cls_mask, mask), dim=1)
159             mask = mask.unsqueeze(1).unsqueeze(2)
160
161         for layer in self.layers:
162             x = layer(x, mask)
163
164         cls_output = x[:, 0, :]
165         return self.fc_out(cls_output)
```