

目 录

第一章 数据集与预处理	2
1.1 数据集介绍	2
1.2 数据预处理	4
1.2.1 数据清洗	5
1.2.2 词表构建	6
1.2.3 数据增强	6
第二章 模型介绍	8
2.1 模型概述	8
2.2 输入层	8
2.3 Transformer 编码器层.....	8
2.4 基于注意力的池化层.....	9
2.5 类别感知特征融合	9
2.6 带残差连接的双分支分类器.....	10
第三章 训练过程	11
3.1 损失函数	11
3.2 优化与正则化	11
3.3 渐进式训练与解冻	12
第四章 结果与分析	13
第五章 结论与改进方向	15

第一章 数据集与预处理

1.1 数据集介绍

ChnSentiCorp-Htl-all 是中文自然语言处理（NLP）领域中广泛应用的标准情感分析数据集之一，源自谭松波博士整理的 ChnSentiCorp 语料库。该数据集专门聚焦于酒店行业，通过采集国内主要在线预订平台的用户评论，并由人工进行情感标注，将评论划分为“正面”与“负面”两类，其中正面评论标签为“1”，负面评论标签为“0”。凭借其覆盖面广、标注质量高等特点，该数据集被广泛用于训练和评估中文情感分析模型，尤其在理解消费者反馈、优化服务流程和把握行业趋势等方面具有重要实用价值。

从数据分布的整体特征来看，ChnSentiCorp-Htl-all 存在较为显著的类别不均衡现象。根据图 1，数据集正面评论共计 5322 条，约为负面评论 2443 条的两倍以上。这一结构差异可能反映了现实情境中的用户表达倾向：多数消费者在服务尚可时不倾向于留言，而在体验极好或极坏时，更有动机表达意见。

此外，在评论长度方面，正负评论也呈现出如图 2 的差异性。统计数据显示，负面评论的平均长度为 164 个字符，显著高于正面评论的 112 个字符，表明用户在表达不满情绪时更倾向于详尽描述具体问题，为后续服务优化提供了精准的参考依据。

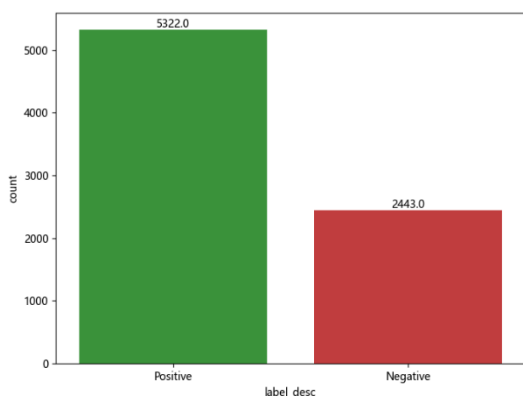


图 1 酒店评论情感倾向分布图

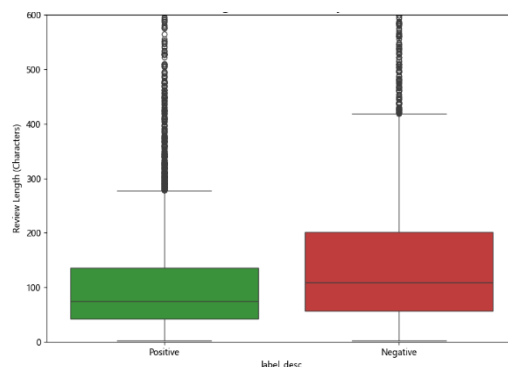


图 2 正负面评论文本长度箱形图

进一步对评论进行词频分析。首先是对数据集中所有评论做词频分析，如图 3 所示，“服务”、“位置”、“交通”、“早餐”、“设施”等词汇构成了所有顾客评价的共同基础。这些高频词汇代表了构成一次住宿体验的核心维度，是每位顾客无论满意与否都会关注的内容。然而，这些词语本身是中性的，例如，“服务”一词并不能告诉我们服务是好是坏。因此，总体词频分析的最大价值只在于帮助识别出顾客的核心关注点，但要真正洞察业务优劣，就必须进一步深入到正、负面评价的细节中去。



图 3 总体评论高频词词云图

从图 4 图 6 可以看出，在正面评论中，关键词如“服务”、“位置”、“交通”、“方便”、“早餐”、“干净”等频繁出现，描绘出消费者对地理便利性、服务质量及环境整洁度的高度认可。相比之下，从图 5 图 7 可以看出，负面评论中高频词汇则集中于“前台”、“设施”、“隔音”、“空调”、“毛巾”、“味道”等，凸显出设施老化、隔音差及服务态度等问题的普遍性与严重性。



图 4 正面评价核心关键词词云

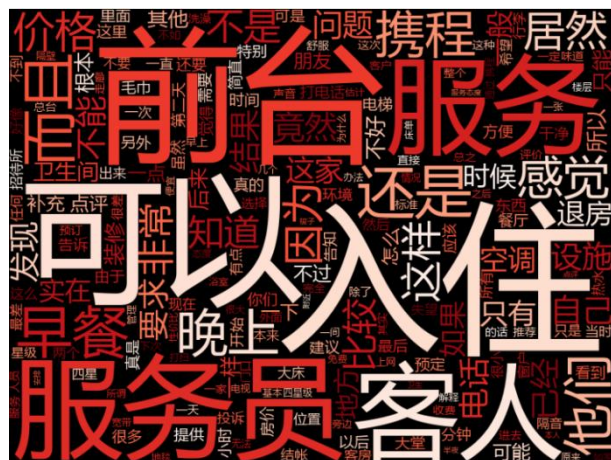


图 5 负面评价核心关键词词云

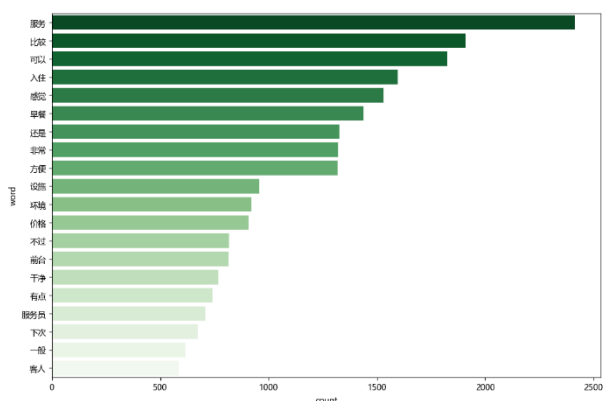


图 6 正面评价高频词

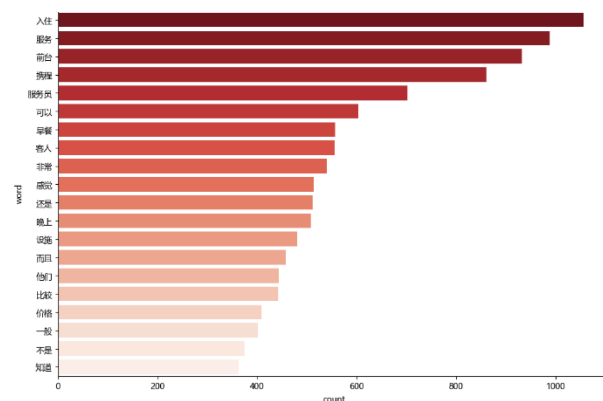


图 7 负面评价高频词

借助 Word2Vec 词向量模型与 t-SNE 降维算法，可将高维词向量投影至二维空间，以图

形方式展示词语间的语义关系。从图 8 可观察到，具有语义相关性的词语自然地聚集成多个语义簇。例如，“位置”、“地铁”、“交通”、“方便”形成了“地理位置”类聚，而“服务”、“前台”、“热情”、“态度”则构成了“服务质量”概念群。此类无监督学习所揭示的语义聚类结构，不仅印证了词频分析所得结论，也从潜在结构层面进一步界定了消费者评价的核心维度。

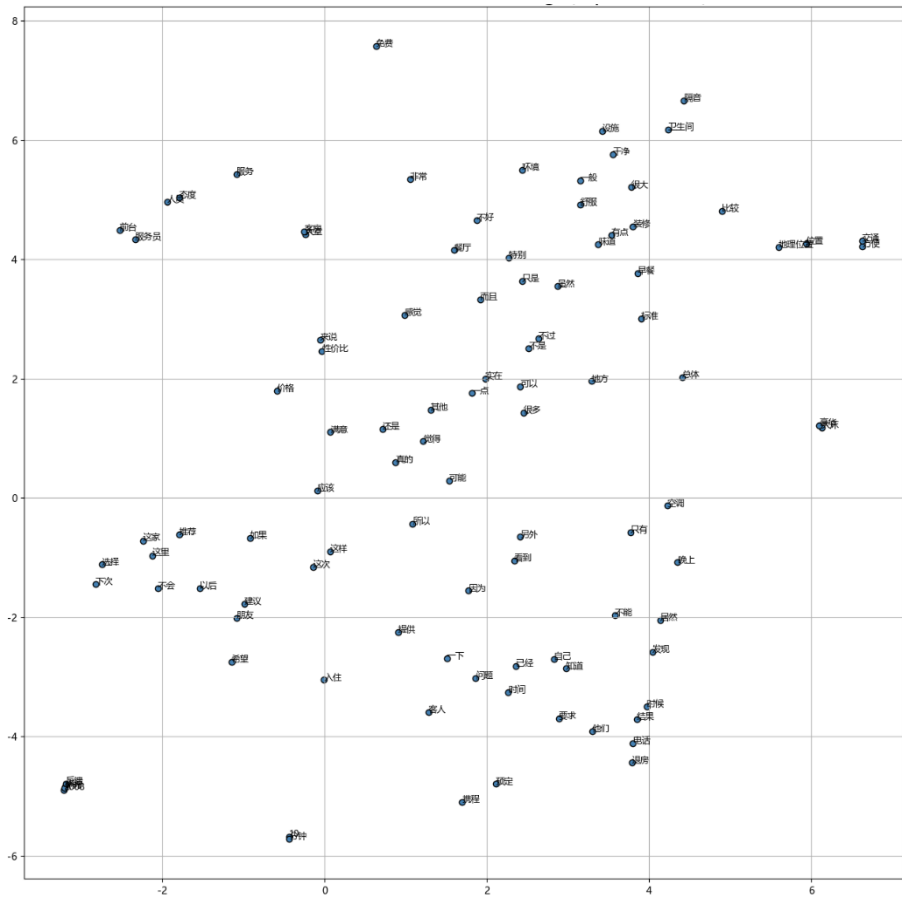


图 8 评论关键词语义聚类分析图

综上所述，ChnSentiCorp-Htl-all 数据集凭借其真实、丰富的用户评论内容，不仅为中文情感分析模型的构建与评估提供了高质量资源，也在宏观分布、微观表达与深层语义等多个层面上，为酒店行业的用户体验研究和服务优化策略提供了多维度的参考视角。

1.2 数据预处理

从上面分析可以看出，数据集有严重的类别不平衡问题，据极易导致模型在训练过程中产生偏见，即模型为了最小化整体误差，会过度关注多数类（正面评论），而忽视对少数类（负面评论）的学习。这最终会导致模型在实际应用中对真正需要关注的负面反馈“视而不见”，从而失去其实用价值。

为提升模型训练的效率与效果，在将数据输入模型前，对数据进行预处理操作，主要包括数据清洗、词表构建与数据增强三个核心环节，旨在规范化文本数据、构建高效的特征空

间，并缓解数据不均衡问题。

1.2.1 数据清洗

原始数据中包含格式不一、信息冗余或无关的元素，会对模型性能构成干扰。因此对原始数据进行清洗。

首先是数据删除。一方面，对数据集进行完整性检查，删除缺失样本。分析发现，数据集中第 6376 行存在 `review` 字段为空的记录，删除该数据，以免对后续产生影响。另一方面，删除信息无效的样本，长度过短的评论往往信息量极低，对模型学习复杂的语义模式贡献有限，甚至可能引入噪声，因此，移除原始文本字符数小于 7 的样本，移除掉的所有样本如表格 1 所示，确实未携带对正负情感判别有帮助的信息

表格 1 删除掉的无效样本

原始文本	标签
房间设施不	1
酒店房间	1
酒店地理位	1
经朋	1
酒店位于深圳	1
就在美美百	0
房间	0

接下来是文本正则化，去除样本噪声，分别进行标准化重复标点、重复字符处理、字符集筛选。首先是标准化重复标点，将连续的标点处理成单个标点，在保持情感信号的同时，减少对于有效文本输入长度的占用。处理结果如表格 2 所示。接下来是重复字符处理，这是对用户评论中常见的情感宣泄式表达进行处理，如将“太棒棒棒棒了”处理为“太棒棒了”，将多个重复字符保留为两重复字符，保留了情感强化的表达方式。最后是字符集筛选，移除了除`[\u4e00-\u9fa5a-zA-Z0-9s! ? ,. ; : "" ' () !?.,:;"""]`外的无关字符。

表格 2 标准化重复标点处理前后样本对比

原始文本	标准化重复标点后文本
还可以。。。。。。。。。	还可以。
总的来说很不错。。。。。。。。。	总的来说很不错。
点评点啥呢，不知应该说什么了。。。。。。。。。	点评点啥呢，不知应该说什么了。

1.2.2 词表构建

完成文本预处理后，进行词表构建以待后续使用，分别进行中文分词、词频统计与过滤、特殊标记添加。

首先是中文分词，采用 jieba 分词库对所有清理后的文本进行分词。其中 jieba 是一个成熟的中文分词工具，能够较好地处理口语化文本，分词结果如表格 3 所示，可见确实完成了较好的中文分词。

表格 3 jieba 中文分词示例

原始文本	地理位置不错,离溧阳汽车站确实很近.但是该酒店新开业没多久,还是有装修味的,最郁闷是另外两件事,根本不能刷 pos 机,别说信用卡了,什么卡都没有,网页上说能接受信用卡简直就是骗人的,很没信誉,不知道谁会带那么多的现金去住宾馆?另外餐厅和酒店是分开的,在餐厅吃饭不能签单的.
分词结果	地理位置 / 不错 / , / 离 / 溧阳 / 汽车站 / 确实 / 很近 / . / 但是 / 该 / 酒店 / 新开业 / 没多久 / , / 还是 / 有 / 装修 / 味 / 的 / , / 最 / 郁闷 / 是 / 另外 / 两件事 / , / 根本 / 不能 / 刷 / pos / 机 / , / 别说 / 信用卡 / 了 / , / 什么 / 卡 / 都 / 没有 / , / 网页 / 上 / 说 / 能 / 接受 / 信用卡 / 简直 / 就是 / 骗人 / 的 / , / 很 / 没 / 信誉 / , / 不 / 知道 / 谁 / 会 / 带 / 那么 / 多 / 的 / 现金 / 去 / 住 / 宾馆 / ? / 另外 / 餐厅 / 和 / 酒店 / 是 / 分开 / 的 / , / 在 / 餐厅 / 吃饭 / 不能 / 签单 / 的 / .

接下来是词频统计与过滤，统计所有词汇的出现频率，并设置 min_freq=2 的阈值。这意味着只将数据集中至少出现过两次的词汇纳入词表。通过过滤，可以完成对词表的降噪与词表规模的控制。其中，出现一次的词，很可能是拼写错误、罕见用语或无意义的字符组合，并且一个过于庞大的词表会增加模型的参数量，延长训练时间，并可能导致对低频词的嵌入学习不充分。

最后是特殊标记添加。在词汇表的最前端，添加<PAD>、<UNK>、<CLS>和<SEP>，其作用分别是将所有序列填充到相同长度、表示所有在词汇表中未出现的词、序列开头添加的分类标记、用于分隔两个不同的句子。

1.2.3 数据增强

在处理类别不平衡问题时，存在多种层面的解决方案，包括调整算法(如代价敏感学习)、修改模型架构以及在数据层面进行操作。其中，数据层面的方法，如重采样(Resampling)和数据增强(Data Augmentation)，通常被认为是最直接且最有效的策略之一。其核心思想是，与其让模型去被动适应有偏的数据，不如主动地去改善数据本身的分布，使其更有利于模型

的学习。

项目使用多种策略进行数据增强，分别是同义词替换、增强负面情感词、随机删除、随机交换、语调调整与组合增强，增强效果如所示表格 4。

表格 4 数据增强结果示例

原始样本	总体感觉比较差,房间家具老旧,商务套间的地毯居然到处是黑斑,让人看了心里非常不舒服。我们因为工作关系，每个月都有几天住在昆山，我的同事们也一致评价誉兴太脏。以后只要有其他的选择，坚决不住誉兴。
同义词替换	总体感觉比较 烂 ,房间家具老旧,商务套间的地毯居然到处是黑斑,让人看了心里非常不舒服。我们因为工作关系，每个月都有几天住在昆山，我的同事们也一致评价誉兴太脏。以后只要有其他的选择，坚决不住誉兴。
增强负面情感词	总体感觉比较 相当差 ,房间家具老旧,商务套间的地毯居然到处是黑斑,让人看了心里非常不舒服。我们因为工作关系，每个月都有几天住在昆山，我的同事们也一致评价誉兴太脏。以后只要有其他的选择，坚决不住誉兴。
随机删除	比较差房间商务套间的地毯居然到处是黑斑,让人了心里非常不。我们工作关系月几天住昆山，的同事一致评价誉兴太脏。以后只要的选择，坚决不住誉兴。
随机交换	总体感觉比较差,房间家具老旧,商务套间的地毯居然到处是黑斑,让人看了心里非常不舒服。我们因为工作关系，每个月都有几天住在昆山， 我也同事们的一 致评价誉兴太脏。以后只要有其他的选择，坚决不住誉兴。
语调调整	总体感觉比较差,房间 确实 家具老旧,商务套间的地毯居然到处是黑斑,让人看了心里非常不舒服。我们因为工作关系，每个月都有几天住在昆山，我的同事们也一致评价誉兴太脏。以后只要有其他的选择，坚决不住誉兴

第二章 模型介绍

2.1 模型概述

为应对酒店评论情感分类任务中的类别不平衡挑战，本项目构建了一个集成了多种先进模块的深度学习模型。该模型以增强 Transformer 编码器为核心，旨在深入捕捉文本的复杂语义特征。模型整体架构自下而上依次由输入层、Transformer 编码器层、基于注意力的池化层以及带残差连接的双分支分类器构成。

其中，输入层将文本转换为包含语义和顺序信息的数值向量，Transformer 编码器层作为模型的核心，深度提取文本的上下文依赖特征，池化层将 Transformer 编码器层输出的可变长度的序列信息聚合成一个固定长度的全局表示，带残差连接的双分支分类器使用独立的分类路径处理正面和负面评价，以应对数据不平衡带来的挑战。

这种设计不仅利用了 Transformer 在处理序列数据上的强大能力，还通过专门的模块设计来提升对关键信息，尤其是负面评论的敏感度，并通过双分支结构针对性地学习正、负面评论的独特表达模式，从而提高模型在不均衡数据集上的分类准确性和鲁棒性。

2.2 输入层

输入层将离散的、非结构化的文本评论转化为模型能够理解和处理的、信息丰富的数值化张量。这是后续所有复杂计算的基础，这部分包括词嵌入与位置嵌入。

首先是词嵌入，通过词嵌入层将输入文本中的每个词汇在词表中索引映射为一个高维、稠密的向量表示。

接下来是位置嵌入，Transformer 模型是“无序”的，它平等地看待所有词语，无法感知它们的先后顺序。然而，在自然语言中，语序至关重要。为了让模型理解词语的位置信息，引入了位置嵌入。该层为句子中的每一个位置生成一个唯一的位置向量。对于生成方式，本项目采用了可学习的位置编码，即位置编码本身也是模型的参数，通过反向传播进行学习，让模型自行学习最优的位置表示，可能比使用通用的、固定的函数具有更强的适应性和表达能力。

位置向量与对应位置的词嵌入向量进行相加。通过这种方式，每个词的最终表示都融合了其语义信息和位置信息，为 Transformer 编码器提供了完整的输入。

2.3 Transformer 编码器层

模型的主体是一个由 6 层 Transformer 编码器堆叠而成的深度网络。每一层编码器的具体配置如表格 5 所示。

模型堆叠 6 层编码器，使模型有足够的去捕捉文本中从局部到全局的深层语义依赖

关系；采用 8 头自注意力机制，允许模型在计算一个词的表示时，能够同时从 8 个不同的“表示子空间”关注序列中的其他部分，从而捕获更多样化的依赖关系；模型内部的隐藏层维度为 512，而前馈神经网络的中间层维度扩展到 2048，提供足够的容量进行非线性变换；将 Layer Normalization 操作置于自注意力/前馈网络之前，而非之后，研究表明，这种结构能够使梯度传播更加稳定，从而支持更深层次的网络训练，并加速模型收敛；选用 GELU 作为激活函数，它被证明在许多基于 Transformer 的模型中性能优于传统的 ReLU。

表格 5 超参数配置

参数	值	描述
Vocal_size	14000	词汇表大小，根据数据集构建
d_model	512	模型内部的隐藏层维度
nhead	8	多头注意力机制的头数
num_layers	6	Transformer 编码器的堆叠层数
dim_feedforward	2048	前馈神经网络的中间层维度
dropout	0.1	Dropout 比率
max_len	400	模型处理的最大序列长度
activation	Gelu	激活函数类型
norm_first	True	采用 Pre-LN 结构

2.4 基于注意力的池化层

在经过 6 层 Transformer 编码器处理后，得到序列中每个 token 的上下文感知表示。为了得到代表整个句子情感的单一向量，传统的做法是使用最大池化（Max Pooling）或平均池化（Mean Pooling）。然而，这两种方法都较为粗糙：平均池化会平等对待所有词，而最大池化则只关注最“突出”的词，两者都可能丢失关键信息。

本项目采用一个独立的多头注意力层作为池化层，将所有 token 表示的平均值作为 Query 向量，表示句子整体，将所有 token 的表示本身同时作为 Key 和 Value，表示单个 token。

通过这种方式，模型可以学习到一个对所有 token 的加权平均。权重的大小由 Query 与每个 Key 的相关性决定。这意味着，模型能够自动地为对判断整个句子情感最重要的词汇，如情感词、否定词，分配更高的权重，从而生成一个信息量更丰富、更具代表性的句子向量。

2.5 类别感知特征融合

传统的分类模型通常在网络的最后一层才引入类别信息，这意味着在整个特征提取过程中，模型本身对于最终的分类目标是“无知”的。本项目引入类别感知机制，创建一个可学习的类别嵌入层 class_embedding，它为每个类别生成一个独立的、与词嵌入维度相同的向量。

将类别嵌入向量与句子表示向量进行拼接。最后将这个拼接后的、更长的向量通过一个线性层 `global_feature` 进行投影和融合。

最终分类之前，向模型的特征表示中显式地注入了关于分类任务本身的信息。相当于给模型一个强烈的信号“请学习那些能够最好地区分这两个类别的特征”。这种机制引导模型在特征学习阶段就生成更具判别力的表示，而不是将所有压力都留给最后的分类器。

2.6 带残差连接的双分支分类器

在特征融合后设计双分支独立分类器，每个分支的输出都是一个标量（logit），分别代表样本属于负面和正面的置信度，以防止类别间的特征干扰。因为在类别不平衡的数据集上，一个共享的分类器很容易被多数类的特征“主导”，导致其为少数类学习到的决策边界被“挤压”或“扭曲”。

因此，设计正负类分支，使得负类分支可以专注于学习和识别那些细微的、独特的负面情感信号，而不用受到大量正面特征的干扰，正类分支同理。

并且每个分支本身并非一个简单的线性层，而是一个包含残差的深度网络。每个 `ResidualBlock` 由两个线性层、两个 `LayerNorm` 层、激活函数和 `Dropout` 组成，并带有一个残差连接，使得每个分类分支都具有很强的非线性建模能力，能够学习到比单一线性层复杂得多的决策边界，同时残差连接又能保证梯度的有效传播，防止网络退化。

第三章 训练过程

3.1 损失函数

在类别不平衡的分类任务中，标准的交叉熵损失函数存在明显缺陷：它平等对待所有样本，导致数量占优的多数类样本产生的总损失会淹没少数类样本的损失，使模型优化方向偏向多数类。为了克服这一点，对损失函数进行设计，使得损失为 `AdaptiveFocalLoss` 和 `LabelSmoothingLoss` 的加权结合，公式如下：

$$Loss = 0.8 \times AdaptiveFocalLoss + 0.2 \times LabelSmoothingLoss$$

首先是 `FocalLoss`，旨在解决难易样本不平衡问题，其核心思想是通过一个调制因子 $(1 - pt)^\gamma$ 来动态调整每个样本在总损失中的权重。其中， pt 是模型预测正确类别的概率， γ 是聚焦参数，本项目中设为 2.0。对于简单样本， pt 趋近于 1， $(1 - pt)^\gamma$ 接近于 0，其损失贡献被大大降低。对于困难样本， pt 趋近于 0， $(1 - pt)^\gamma$ 接近于 1，其损失贡献保持不变。通过这种方式，模型被迫将更多的“注意力”和计算资源用于学习那些难以区分的样本，从而提高了整体性能。

接下来是动态 `Alpha` 类别加权，标准的 `Focal Loss` 使用一个固定的超参数 `alpha` 来处理类别不平衡，而本项目实现了一种动态的 `alpha` 权重。首先根据类别不平衡的先验知识，为不同类别分配一个固定的基础权重。`alpha = torch.where(targets == 0, 0.75, 0.25)` 这行代码意味着，负面样本的基础权重被设为 0.75，而正面样本则为 0.25，直接在数值上弥补了其数量上的劣势。在此基础上，进一步用模型的预测置信度来动态调整这个权重：

$$alpha = alpha * (1 - target_probs)$$

这里的 `target_probs` 是模型对正确标签的预测概率，它将类别不平衡和难易不平衡两个问题联系了起来。一个样本的最终权重，不仅取决于它属于哪个类别，还取决于模型对它的分类有多“自信”。一个模型很不确定的负面样本，将获得最高的损失权重，从而得到模型最优先的“关照”。

最后，我们将上述自适应 `Focal Loss` 与标签平滑进行了加权融合。其中，标签平滑是一种强大的正则化技术。它将原本的硬标签替换为软标签，其中 `label_smoothing=0.`，防止模型在训练时做出过于绝对和自信的预测。通过鼓励模型输出更“柔和”的概率分布，标签平滑可以有效提升模型的泛化能力和校准度。

损失函数构建了一个既能宏观调控类别平衡，又能微观聚焦困难样本，同时还能抑制过拟合的强大训练目标，为模型收敛到更优的解提供了坚实的数学基础。

3.2 优化与正则化

首先在优化器的选择上，选用 `AdamW` 而非传统的 `Adam`。`AdamW` 优化器修复了 `Adam`

中权重衰减实现方式的缺陷，能够更有效地进行 L2 正则化，防止模型过拟合。

接下来是学习率调度器，采用 `CosineAnnealingWarmRestarts`，让学习率在一个周期内、按照余弦曲线从初始值平滑地下降到最小值。在周期结束后，学习率会“热重启”到初始值，并开始下一个、可能更长的周期。这种周期性的升降有助于优化器跳出可能陷入的局部最优“山谷”，去探索更广阔参数空间，从而找到更优、更泛化的解。

在训练过程中，每一个进 batch 后进行梯度裁剪，将所有参数的梯度范数裁剪到最大值 1.0，是训练深度网络尤其是 RNN 和 Transformer 时，防止梯度爆炸、确保训练稳定的操作。

在每个 epoch 结束后，在验证集上评估模型的宏平均 F1 分数。如果 F1 分数连续多个 epoch 没有超过历史最佳值，训练将自动停止。选择 F1 分数而非准确率作为监控指标，是因为在不平衡数据集上，F1 能更均衡地反映模型在所有类别上的表现。

3.3 渐进式训练与解冻

从零开始训练一个包含 6 层 Transformer 的深度模型是极具挑战性的，很容易因为不恰当的初始化和学习率设置导致训练不稳定或收敛缓慢。为了解决这个问题，设计了 `ProgressiveTrainer`，其核心是分阶段、渐进式的训练策略。

该策略将整个训练过程划分为两个主要阶段。

第一个阶段预训练分类头，在训练的初始阶段冻结模型中参数量巨大的主干部分，包括词嵌入层和全部 6 层 Transformer 编码器。此时，只有双分支分类器部分的参数是可训练的。这个阶段的目标是让随机初始化的分类器先快速学习到一个“像样”的决策边界。由于只训练一小部分参数，我们可以使用一个相对较大的学习率 ($lr=1e-3$) 来加速这个过程，而不用担心破坏整个模型的稳定性。这相当于先“教会”模型的“大脑”如何分类。

第二个阶段遵循渐进式解冻的时间表，进行全模型微调。训练进度 $< 30\%$ 时继续保持主干网络冻结，专注训练分类头； $30\% \leq \text{训练进度} < 60\%$ 时，解冻词嵌入层和 Transformer 的后 3 层。这是一种折衷，允许与输入和高层语义关系最密切的部分开始调整，同时保持底层特征提取器的稳定；训练进度 $\geq 60\%$ 时，解冻所有参数，对整个模型进行精细微调，提升了训练的稳定性和最终的性能。

第四章 结果与分析

本次实验中，构建的文本情感分类模型在测试集上取得了令人满意的性能。在训练集、验证集和测试集上的核心性能指标汇总如表格 6 所示。从上表可以看出，模型在训练集上表现优异，准确率超过 96%，F1 分数接近 0.96。在从未见过的验证集和测试集上，性能虽有小幅下降，但准确率依然保持在 91%–92% 的较高水平，宏平均 F1 分数也达到了 0.90 以上。并且模型对于正负样本的分类准确率相近，说明通过模型设计，很好地处理了正负样本不平衡问题。

这表明模型具备了良好的泛化能力。训练集与测试集之间约 4.8% 的准确率差距和约 5.6% 的 F1 分数差距，处于一个合理的范围内，说明模型在一定程度上有效避免了严重的过拟合现象。这得益于我们在训练过程中采用的 AdamW 优化器、余弦退火学习率、梯度裁剪以及标签平滑等正则化策略。

表格 6 模型输出结果

数据集	准确率	F1	TP	TN	FP	FN	负样本 准确率	正样本 准确率
训练集	0.97	0.96	4082	1908	47	175	0.98	0.96
验证集	0.93	0.92	499	222	22	34	0.91	0.94
测试集	0.92	0.90	500	212	33	32	0.87	0.94

查看误报数据如表表格 7，这类错误的核心在于褒中带贬或先抑后扬的复杂句式。模型被评论中的负面词汇或抱怨语气所吸引，而忽略了其最终的正面或中立立场。

表格 7 误报样本分析

误报样本	分析
"硬件设施没的说,房间豪华气派.唯一缺点是前台的退房服务不好...不时要停下来应付别的服务,一点都不尊重面前的客人."	这条评论以强烈的正面肯定“硬件没的说”、“豪华气派”开头，但后半段详细描述了服务上的缺点。尽管评论者可能整体上认可酒店的硬件，但大量的负面描述主导了文本的情感信号，模型给出负面判断
"作为酒店的老客户，恐怕以后要做另外的选择了服务水平在下降，价格却一升再升，再这样下去，下次不会再入住了。"	这条评论虽然标签为正面，但其内容完全是负面的抱怨和警告。这是一个典型的数据标注歧义或错误的例子。模型将其正确地识别为负面情绪，正面置信度仅 0.18，这反而证明了模型的判断力。

查看漏报数据如表格 8，可以看出误报的主要原因与漏报类似，多为贬中带褒的结构，或是客观陈述中夹杂了少量正面词汇，导致模型误判。

表格 8 漏报样本分析

漏报样本	原因分析
酒店地理位置较好...不过房间非常之小...餐饮部手艺不错，不过价格太高...最让人生气的是...收了 60 元的洗涤费！	此评论包含了多个褒贬交织的子句，“位置较好”但“房间小”、“手艺不错”但“价格高”。这种复杂的结构对模型构成了巨大挑战。模型似乎难以在多个对立的情感信号中权衡并找出主导情绪，最终被局部的正面词汇如“较好”、“不错”所误导
方便倒是真方便。可是晚上噪音太大了。正好是个十字路口。整个晚上全过车。鸣笛不止。睡觉轻的还是不要去住了	评论以“真方便”开头，这是一个强烈的正面信号。尽管后文紧跟着“可是”进行了转折，并详细描述了噪音问题，但模型显然给予了开头的正面词汇过高的权重，最终做出了错误的判断，负面标签的样本，被预测为正面的置信度为 0.20

第五章 结论与改进方向

综合来看,本模型在酒店评论情感分析任务上取得了良好的整体性能和泛化能力。然而,通过深入的错误分析,我们发现模型的主要挑战在于理解和处理包含复杂转折、褒贬交织的长句式。模型倾向于被局部出现的强情感词汇所影响,而对上下文的整体逻辑和核心情感的把握尚有不足。此外,数据集中也可能存在少量标注不一致或有歧义的样本,同样对模型训练和评估造成了干扰。

未来的改进可以从数据、模型、和训练策略三个层面系统地展开,以期构建一个更加鲁棒和智能的情感分析系统。

数据的质量和多样性是模型性能的基石。当前模型在处理褒贬不一的复杂评论时表现不佳,这首先提示我们应从数据源头寻找解决方案。分析中发现,部分样本的标签与其内容存在明显歧义。建议引入人工交叉验证或利用置信度学习等算法,系统性地筛选出训练集中可能存在标签错误的样本。对这些样本进行专业的重标注,可以净化训练数据,为模型提供更一致的学习信号,从根本上提升模型的可靠性。

并且鉴于模型在处理“先扬后抑”或“先抑后扬”句式上的困难,可以设计一套针对性的数据增强策略。例如,可以编写模板来生成大量包含“虽然…但是…”、“优点是…缺点是…”、“除了…之外都很好”等典型转折结构的伪样本。通过让模型在训练阶段接触更多此类复杂样本,可以显著增强其捕捉上下文逻辑和主要情感的能力。

在模型结构层面,可以引入更强大的预训练语言模型,将当前的 Transformer 编码器替换为更大规模、在中文语料上表现更优异的预训练模型,例如百度的 ERNIE 系列模型。这些模型通过融入知识图谱等外部信息,对中文的语义理解更为深刻,可能在处理复杂逻辑和隐晦情感时表现更佳。

在训练策略上,可以实施课程学习,模拟人类的学习过程,让模型“先易后难”。在训练初期,只使用情感明确、结构简单的样本;随着训练的进行,逐步增加包含转折、含蓄表达的复杂样本。这种由简到难的训练课程,可以帮助模型构建一个更稳定、更扎实的语义理解基础,避免在早期就被复杂样本误导。或者探索对抗式训练,通过在输入文本中添加微小的、人眼难以察觉的扰动来生成“对抗样本”,并让模型在训练中努力对这些样本做出正确分类。这种方法可以极大地提升模型的鲁棒性,使其在面对不规范或略有变化的输入时,表现得更加稳定,降低对局部关键词的过度敏感。