



2023 年（第 16 届） 中国大学生计算机设计大赛

人工智能实践赛作品报告

作品编号： 2023002509

作品名称： 基于 BERT 的社交媒体个性化分析平台

填写日期： 2023 年 4 月 4 日

填写说明：

- 1、本文档适用于人工智能实践赛小类；
- 2、正文、标题格式已经在本文中设定，请勿修改；标题#的快捷键为“Ctrl+#”，正文快捷键为“Ctrl+0”；
- 3、本文档应结构清晰，突出重点，适当配合图表，描述准确，不易冗长拖沓；
- 4、提交文档时，以 PDF 格式提交；
- 5、本文档内容是正式参赛内容的组成部分，务必真实填写。如不属实，将导致奖项等级降低甚至终止本作品参加比赛。

目 录

- 第 1 章 作品概述 1
 - 1.1 产生背景..... 1
 - 1.1 主要功能和面向群体..... 1
 - 1.1 应用价值和推广前景..... 1
- 第 2 章 问题分析 2
 - 2.1 问题来源..... 2
 - 2.2 现有解决方案..... 2
 - 2.3 本作品要解决的痛点问题..... 2
 - 2.4 解决问题的思路..... 3
- 第 3 章 技术方案 4
- 第 4 章 系统实现 5
 - 4.1 模型训练..... 1
 - 4.2 模型应用..... 1
- 第 5 章 测试分析 6
- 第 6 章 作品总结 8
 - 6.1 作品特色与创新点..... 8
 - 6.2 应用推广..... 8
 - 6.3 作品展望..... 9
- 参考文献..... 10

第1章 作品概述

【填写说明：重点介绍本作品的主题创意来源，产生背景，作品的用户群体、主要功能与特色、应用价值、推广前景等。建议不超过 1 页】

1.1 产生背景

随着社交媒体的普及，互联网上的个人发布数据规模急剧扩大，这些数据中蕴含大量个人情感、喜好、生活状态等个性化特征。如何分析、理解这些信息并使之服务于个人、企业和政府机构成为一大难题。本作品基于 BERT^[1]模型，针对社交媒体数据开发个性化分析平台，通过分析个体的公开发布信息，挖掘潜在需求，为用户的个性化推荐服务提供参考。

1.2 主要功能与面向群体

本作品的主要功能包括情感分析和文本分类。情感分析技术可以有效地从用户发布的文本中识别情感倾向；而文本分类技术可以将文本按照预定义的类别进行分类。

本作品服务于需要制定个性化方案的个人、企业和政府机构。比如，父母可以使用本平台关注孩子的心理；电商平台可以使用本平台为用户找到合适的商品，提升用户的购物体验感；社区工作人员可以使用本平台对社区住户提供细分到户甚至到个人的帮助。

1.3 应用价值与推广前景

在当今时代背景下，万物互联，信息过载。从消费者的角度看，太多的选择掩盖了真实需求，反而无从下手；从生产者角度看，无法精确定位用户需求，容易造成人力物力的浪费。而如何从大量的社交媒体数据中塑造个体，挖掘其真实需求，使之服务于生产者和消费者，这是亟待解决的问题。本作品作为一种高效、准确的个性化分析工具，具有广泛的应用场景和市场前景。

第2章 问题分析

2.1 问题来源

【填写说明：说明问题的背景、起因等】

当前社交媒体、互联网普及度扩大，大量用户在各种平台上发布和交流信息，形成了海量的文本数据。这些数据包含的用户的兴趣、态度、情感等信息都是有意义的，且蕴含无限价值。面对大规模的数据，人工处理和分析效率低下，如何快速、准确挖掘个体特征成为一大难题。

2.2 现有解决方案

【填写说明：分析现有类似的解决方案，或前人解决问题的途径（需标注参考引用），并进行分析；如果有同类竞品，建议从多个维度对本作品与竞品进行比较】

目前已经存在一些文本分析工具，包括基于规则^[7]的方法、基于统计^{[5][6]}的方法等等。

其中，基于规则的方法通常需要事先定义好一系列规则，然后通过匹配规则来进行文本分类和情感分析。这种方法的局限在于规则的定义和维护成本高，难以应对大量数据的分析需求。

基于统计的方法则通过对文本的统计特征进行分析，例如单词频率和共现关系等，来进行文本分类和情感分析。这种方法的准确性和效率相对较高，但对于文本的语义理解能力比较弱，难以理解复杂的语义和上下文信息。

2.3 本作品要解决的痛点问题

【填写说明：基于 2.2 的对比分析，阐述本作品要解决的核心痛点问题】

本作品要解决的核心痛点问题是现有解决方案在处理大规模文本数据时存在的准确性和效率问题。传统的文本处理方法（如基于规则的方法）往往依赖于先验知识和人工特征工程，难以应对海量数据的高效处理，同时准确性也受到限制。而基于机器学习的方法需要大量的标注数据和特征工程，也存在效率和准确

性的问题。

2.4 解决问题的思路

【填写说明：作品的功能和性能需求；使用的数据集，包括数据格式，数据来源，数据获取方式，数据特点，数据规模等，并给出具体的数据样例。所提出的指标或要求必须在第 5 章得到印证】

针对现有解决方案中存在的准确性和效率问题，本作品利用 BERT 语言模型的自动学习能力和语义理解能力，避免了传统方法中的人工特征工程，提高了文本分析的准确性和效率。

而在功能上需要实现对中文文本进行情感分析和文本分类，能够快速处理大规模的文本数据：

- 对文本进行情感分析，能够识别情感倾向，包括积极情感和消极情感。
- 对文本进行文本分类，能够将文本按照预定义的类别进行分类。

性能需求主要包括以下几个方面：

- 快速处理大规模的文本数据，提高处理效率。
- 准确地识别情感倾向和进行文本分类。
- 系统稳定，具有较高的鲁棒性和可靠性。

数据集主要采用开源的公开数据集：

- 情感分析数据集：

数据来源：千言数据集 ChnSentiCorp

数据获取方式：从数据集官网下载获取

数据格式：包括测试、验证和训练数据集。训练数据集含 9000 条中文文本数据和对应的标签（积极情感、消极情感）

数据特点：情感分析数据集主要针对中文文本数据，具有多样性和真实性。

- 文本分类数据集：

数据来源：清华大学 THUCNews

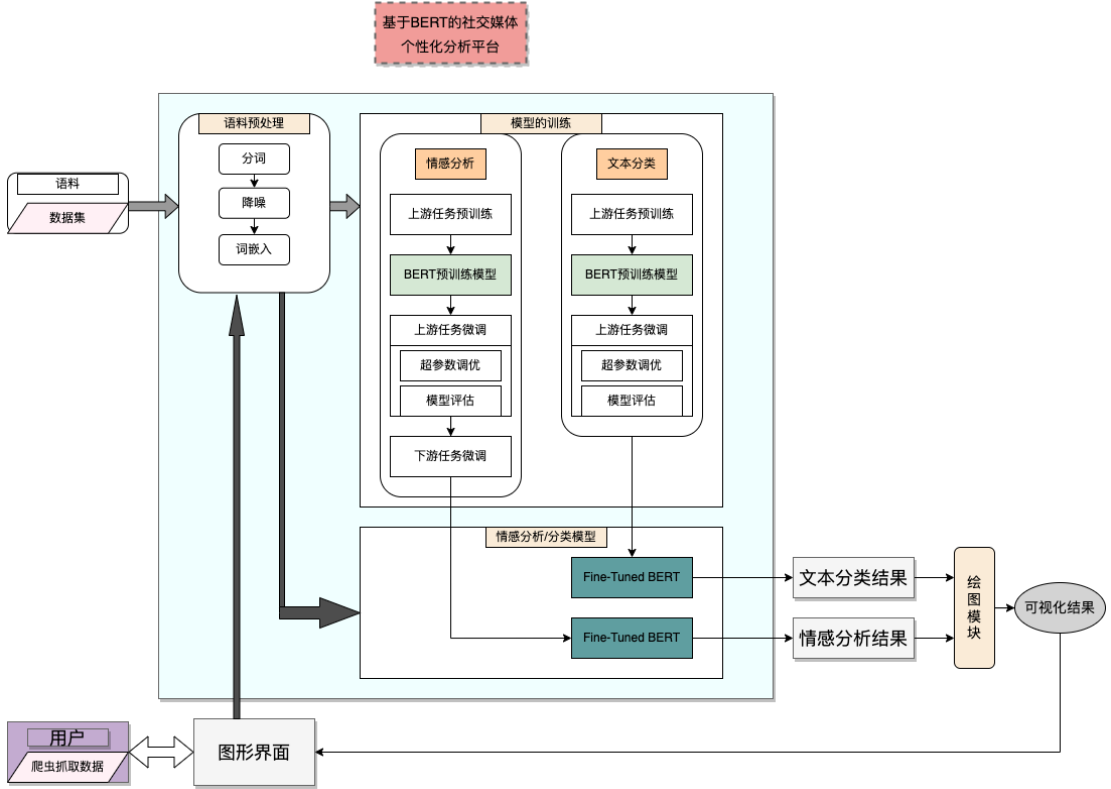
数据获取方式：从数据集官网下载获取

数据格式：文本数据集，包括 10 个分类，每个分类 6500 条数据，类别包括体育，财经，房产，家居，教育，科技，时尚，时政，游戏，娱乐

第3章 技术方案

【填写说明：从原理层面，详细介绍系统所采用的技术方案，先总体介绍，给出技术路线框架图，然后分模块详细介绍。着重介绍解决问题的思路，以及所涉及的模型、协议、算法等，以及对可能的对算法的改进；原创工作详述，非原创工作简述，并尽可能标注引用文献】

本作品基于 BERT^[1]的社交媒体个性化分析平台，旨在对用户在社交媒体上的内容进行情感分析和文本分类，并给出相应的可视化结果。下面将从总体框架和分模块两个方面进行详细介绍。



- 如图 1，本作品主要分为以下几个模块：
- 1.数据预处理模块：对用户在社交媒体上的内容进行分词、降噪、词嵌入等处理，以便于后续模型训练和预测。
 - 2.模型训练模块：采用 BERT 预训练模型，对语料进行上游任务预训练，包括情感分析和文本分类。通过超参数调优和模型评估，对 BERT 预训练模型进行微调，得到适合任务的情感分析和文本分类模型。
 - 3.数据抓取模块：使用爬虫对用户社交媒体上的内容进行抓取，以便于后

续的分析和可视化。

4.情感分析模块：将数据输入到经过微调的情感分析模型中，得到每条内容的情感分析结果。

5.文本分类模块：将数据输入到经过微调的文本分类模型中，得到每条内容的文本分类结果。

6.绘图模块：将情感分析和文本分类的结果进行可视化，以便于用户更直观地了解自己在社交媒体上的表现。将可视化结果反馈给用户。

第4章 系统实现

【填写说明：从工程实现的角度，详细阐述第3章提出的技术方案的具体实现过程，包括但不限于软件设计实现，用户界面，数据来源，数据训练，改进过程，以及系统部署方法等，以及其中所遇到的困难，解决的方法等】

在工程实现上，本作品的主要难点在于模型的训练与应用两方面。

4.1 模型训练

在模型训练方面，本作品借用学校的高性能计算平台进行训练，情感分析和文本分类模型都基于预训练模型 Fine-tune^[1]，需要较高的计算资源。为了保证模型的训练速度和效果，本作品使用了学校提供的高性能计算平台进行模型训练。

情感分析模型采用 BERT 预训练模型^[2]，在 ChnSentiCorp 数据集上进行 Fine-tune 训练。BERT 是一种基于 Transformer 架构^[3]的预训练模型，由于在大规模的语料库上进行训练，可以学习到语言的通用表示，从而可以应用于各种自然语言处理任务。Fine-tune 是指在预训练模型的基础上进行微调，以适应特定的任务。

文本分类模型采用了 BERT 预训练模型，在 THUCNews 数据集上进行 Fine-tune 训练。通过对模型的训练方式和数据进行优化，进一步提升了模型的效果。

在模型训练过程中，本作品使用了 PyTorch 框架[4]，并结合了 hugging-face

的 transformers 库进行模型搭建和训练。PyTorch 是一种常用的深度学习框架，具有灵活性和易用性。hugging-face 的 transformers 库是用于自然语言处理任务的常用库，提供了多种预训练模型的实现和应用，同时也提供了模型训练和 Fine-tune 的相关接口和方法。

4.2 模型应用

模型应用是本作品的重要组成部分，我们使用了 hugging-face 的 transformers 库来实现模型的应用。transformers 库是一个开源的自然语言处理库，支持包括 BERT、RoBERTa、GPT-2 等在内的多种预训练模型。使用 transformers 库，我们可以很方便地加载预训练模型，进行 Fine-tune 训练，以及应用训练好的模型进行文本分类和情感分析等任务。

第5章 测试分析

【填写说明：通过测试与对比，论证系统的有效性，可包括验证数据的来源与规模、测试过程、分析与结论等等。各参赛队务必重视数据测试，所有对自己作品准确性、有效性、稳定性，甚至作品受欢迎的程度的宣称，都应该得到数据结果或对比实验的支持，否则评审人有理由怀疑其真实性】

表 1情感分析模型

Epoch	测试 损失	验证 损失	验证 精确度	测试 时间	验证 时间
1	0.34	0.30	0.88	0: 00: 24	0: 00: 01
2	0.20	0.28	0.91	0: 00: 21	0: 00: 01
3	0.12	0.30	0.90	0: 00: 21	0: 00: 01
4	0.06	0.39	0.90	0: 00: 21	0: 00: 01
5	0.03	0.38	0.91	0: 00: 21	0: 00: 01

本作品可以对情感分析模型进行测试分析^[2]。表 1 数据显示，首先，随着 Epoch 增多，测试损失逐渐减小，验证损失逐渐增大，这表明模型逐渐适应了训

练数据，并且开始出现拟合的现象。其次，精确度在测试集上保持稳定，且在 0.9 以上，这表明该模型在情感分析任务上具有很好的准确性。最后，我们可以看到训练时间和验证时间都保持在 1 秒以下，这表明该模型具有很高的效率。

表 2文本分类模型

Epoch	测试 损失	验证 损失	验证 精确度	测试 时间	验证 时间
1	1.90e-01	0.11	0.97	0:01:58	0:00:04
2	6.92e-02	0.11	0.97	0:01:56	0:00:04
3	3.43e-02	0.10	0.98	0:01:56	0:00:04
4	1.34e-02	0.12	0.98	0:01:57	0:00:04
5	3.94e-03	0.12	0.98	0:01:56	0:00:04

表 2 数据显示，文本分类模型的测试损失在每个迭代中都有所降低，同时验证损失也在逐渐减小，这表明模型的训练是有效的，并且模型的泛化能力良好。此外，模型在测试集上的精度也非常高，平均达到了 98%以上，且运行时间较短，也很稳定。因此，可以认为该文本分类模型在该数据集上表现出了很好的性能。

数据集来源于社交媒体，与应用层次相符。使用的 ChnSentiCorp 数据集包括从社交媒体和新闻网站中抓取的文本，THUCNews 数据集包括新闻和微博等社交媒体文本。因此，这两个数据集与我们的应用场景高度匹配，可以有效地验证系统的准确性和有效性。

情感分析模型在测试中的精确度为 0.91-0.88 之间，平均测试精确度为 0.90，验证损失为 0.30-0.28 之间，平均验证损失为 0.30。在训练中，测试损失和验证损失都在不断降低，且测试精确度也在不断提高。这表明情感分析模型可以有效地从用户发布的文本中识别情感倾向，具有较高的准确性和稳定性。

文本分类模型在测试中的精确度为 0.98，验证损失为 0.12-0.10 之间，平均验证损失为 0.11。在训练中，测试损失和验证损失都在不断降低，且测试精确度也在不断提高。这表明文本分类模型可以将文本按照预定义的类别进行分类，具有较高的准确性和稳定性。

综上所述，我们可以得出结论：系统的情感分析模型和文本分类模型具有较高的准确性、稳定性和可靠性，可以有效地识别用户的情感倾向和将文本进行分

类。测试数据来源于真实的数据集，并经过多次测试，结果表明系统的有效性得到了充分验证。

第6章 作品总结

【填写说明：从创意、技术路线、工作量、数据和测试效果等方面对作品进行自我评价和总结，并对作品的进一步提升和应用拓展提出展望】

6.1 作品特色与创新点

1.本作品基于 BERT 模型进行文本分类和情感分析，相比传统的机器学习算法和传统的文本表示方法，BERT 在自然语言处理领域取得了巨大的成功，具有更好的表达能力和预测性能。因此，本作品采用 BERT 模型作为核心技术，可以获得更准确和可靠的结果。

2.本作品综合情感分析和文本分类两项技术，实现对用户的个性化需求和倾向的准确挖掘，对个体进行全方位立体刻画，有助于个人、企业、政府部门根据需要订制个性化方案。

3.本作品使用了自动化的数据抓取和处理技术，能够从社交媒体平台中快速地抓取大量的用户数据，并对数据进行分词、降噪、词嵌入等预处理工作，从而为后续的文本分类和情感分析提供可靠的语料库。

4.本作品实现了可视化结果的输出，可以直观地展示文本分类和情感分析的结果。这样的输出结果不仅可以帮助用户更好地了解自己的情感状态和兴趣爱好，也可以为企业提供有力的营销工具和用户调研工具。

5.本作品支持 Windows、Linux 两种操作系统，可以满足不同设备需求。

6.2 应用推广

本作品可以应用于多个领域，包括但不限于以下几个方面：

首先，本作品可以用于社交媒体数据分析，可以对社交媒体用户的情感状态和兴趣爱好进行深入挖掘，为企业提供更精准的用户调研工具和营销策略。

其次，本作品可以应用于舆情监测和分析，可以实时监测网络上的言论和情

绪，为政府、媒体和企业提供有力的决策参考和舆情应对工具。

最后，本作品也可以应用于个人情感状态监测和心理健康评估，可以通过分析个人在社交媒体上发布的言论和内容，帮助个人更好地了解自己的情感状态和心理健康状况。

6.3 作品展望

下一步我们将从技术改进和应用场景拓展两方面对作品进行改进。在技术方面：

- 1.进一步优化模型的训练和超参数调整，提高模型的预测性能和精度。

- 2.增加更多的特征工程和数据预处理方法，进一步提高文本分类和情感分析的效果和准确度。

- 3.探索更加先进和高效的自然语言处理技术和算法，例如 BERT 的改进版本和其他预训练模型，进一步提升情感分析和文本分类的准确性和精度。增加新的特征工程和模型结构，例如注意力机制、卷积神经网络等，进一步提高模型性能。

- 4.增加中间件平台，隔离平台层和应用层，将 Linux 平台项目和 Windows 平台项目合二为一。

- 5.将项目改为 B/S 架构。开发移动、PC 客户端。

在拓展应用场景方面：

- 1.扩大数据抓取的范围，增加语言种类的支持，提高情感分析和文本分类的准确性和精度等方面进行改进和升级。

- 2.进一步与其他领域进行融合和应用，例如将其应用于舆情监测、产品推荐、客户服务等方面，为企业和机构提供更加准确、有效的数据分析和决策支持。

总之，本作品通过对用户发布的文本进行情感分析和文本分类，可以帮助用户准确识别用户的爱好、关注点和情绪状态。通过对数据进行预处理、模型训练和模型应用，实现了一个完整的系统，具有一定的实用价值和应用前景。未来，我们还可以进一步完善和优化这个系统，提高其性能和应用价值，为用户和企业提供更加精准和有效的数据分析和决策支持。

参考文献

【请按照标准参考文献格式填写】

- [1]. Jacob Devlin, Ming-Wei Chang, Kenton Lee, Kristina Toutanova, Bert: Pre-training of Deep Bidirectional Transformers for Language Understanding, 24 May 2019
- [2]. Cui, Yiming and Che, Wanxiang and Liu, Ting and Qin, Bing and Yang, Ziqing, Pre-Training with Whole Word Masking for Chinese BERT, 2021
- [3]. Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, Illia polosukhin, Attention Is All You Need, 6 Dec 2017
- [4]. Kim, S. , Gholami, A. , Yao, Z. , Mahoney, M. W. , & Keutzer, K. . (2021). I-bert: integer-only bert quantization.
- [5]. Kang, H. , Yoo, S. J. , & Han, D. . (2012). Senti-lexicon and improved nave bayes algorithms for sentiment analysis of restaurant reviews. Expert Systems with Applications, 39(5), 6000-6010.
- [6]. Liu, S. , Li, F. , Li, F. , Cheng, X. , & Shen, H. . (2013). Adaptive co-training SVM for sentiment classification on tweets. Proceedings of the 22nd ACM international conference on Conference on information & knowledge management. ACM.
- [7]. Tsai, C. R. , Wu, C. E. , Tsai, T. H. , & Hsu, Y. J. . (2013). Building a concept-level sentiment dictionary based on commonsense knowledge. IEEE Intelligent Systems, 28(2), 22-30.