

# 中国大学生计算机设计大赛



## 大数据实践赛作品报告

作品编号： 2023035163

作品名称： SIU: Sentiment Identification and Understanding

版本编号： v0.0.1

填写日期： 2023.3.26.

### 填写说明：

- 1、本文档适用于大数据实践小类；
- 2、正文一律用小四号宋体，1.3倍行距；一级标题为二号黑体，其他级别标题如有需要，可根据需要设置；
- 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
第 2 章 问题描述 .....	3
第 3 章 技术方案 .....	4
第 4 章 系统实现 .....	11
第 5 章 系统评测 .....	12
第 6 章 安装使用 .....	14
第 7 章 作品总结 .....	15
参考文献 .....	20

# 第 1 章 作品概述

【填写说明：本部分非常重要，建议 800 字以内。简要说明作品的意义、技术特色、实现方法、运行（或应用）效果等等。着重介绍作品的特色、和运行（或应用）效果。】

## 1.1 作品简介

我们使用预训练的 BERT 模型对中文文本进行情感分类，分析文本是正向的还是负向的。接着，使用爬虫从 B 站上爬取数据，使用此模型对文本进行分类。最后，对数据进行分析，例如对均值进行估计以判断舆情对某一方面的总体态度，对方差进行分析以判断舆情对某一个领域的态度的分化程度，以及对各个时间段上的舆情情况的变化进行分析等。

## 1.2 作品意义

方便快捷地对大量的舆情文本进行分析，可以帮助我们更好地了解人们对某个话题或事件的看法和反应。这对于政府、企业、媒体等各种机构来说都非常有用，可以帮助他们更好地了解公众的需求和意见，从而更好地制定政策、推出产品、进行宣传等。此外，这个作品也可以作为一个数据分析工具，帮助研究者更好地了解某个领域的舆情走向和热点问题，为后续研究提供有价值的参考和支持。

## 1.3 实现方法

1、情感分类模型：对具有预训练的二分类大语言模型 bert-base-chinese 进行微调，然后输出是 0 或 1，代表其情感是正向还是负向；

2、爬虫：使用 Python 的 Selenium 库和 Requests 库，爬取 B 站各个 tag 的视频评论数据。

3、数据分析：使用 Pandas 进行统计分析，使用 Matplotlib 绘制图片进行可视化，使用 Excel 绘制更多数据可视化图形，如饼图等。

## 1.4 技术特色

我们作品的技术特色非常明显，它不仅使用了 BERT 模型进行文本分析，还结合了爬虫技术，可以实时地从网络上爬取评论并对舆情进行分析。

BERT 是一种预训练的自然语言处理模型，可以对大量的文本数据进行训练，从而学习到丰富的语言表达能力。使用 BERT 模型可以更好地理解和分析人们对某个话题或事件的看法和反应，从而为政府、企业、媒体等机构提供有用的数据支持。

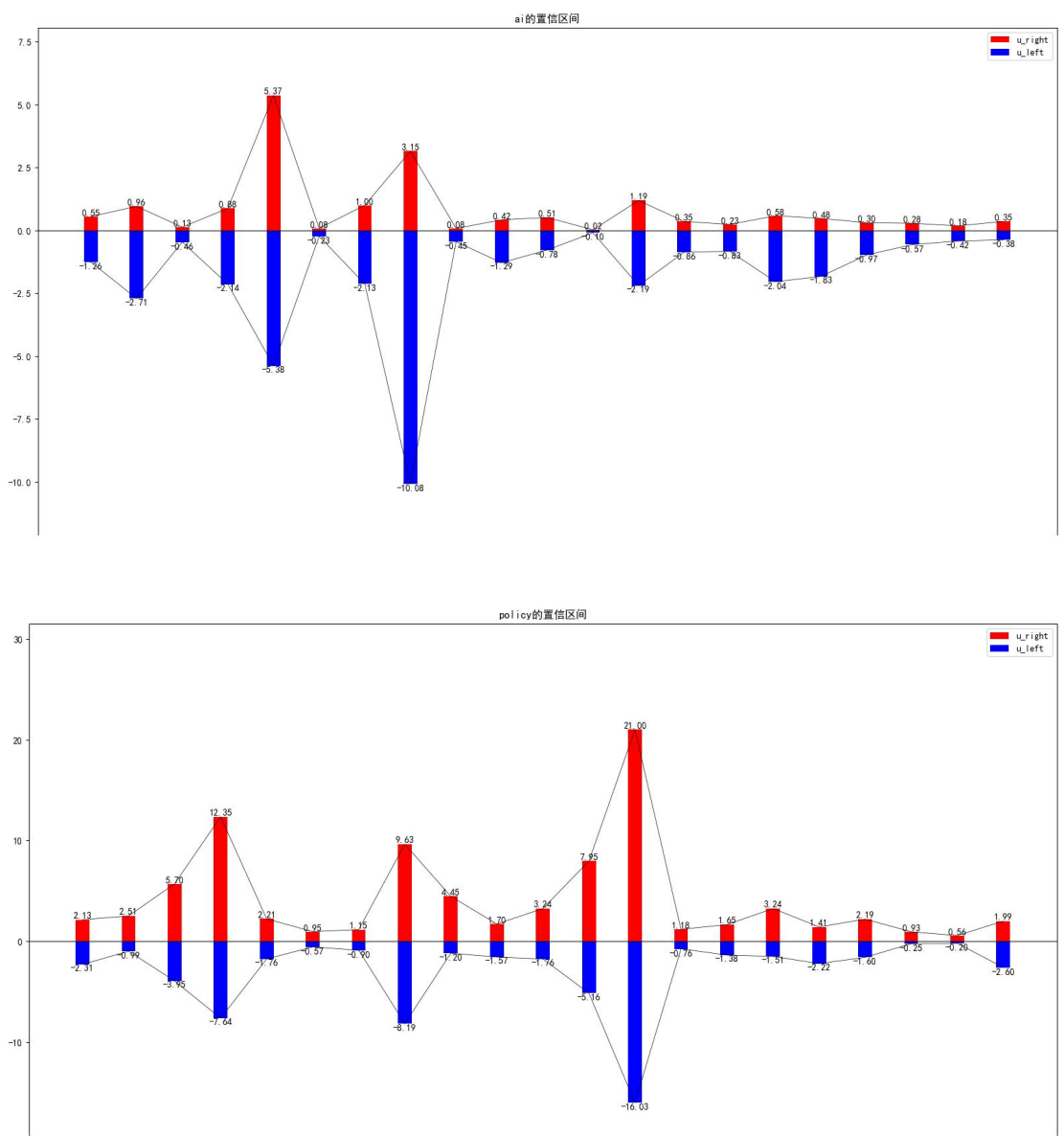
同时，使用爬虫技术可以对指定的网站或社交媒体进行监测，实时地获取用户的评论和回复，并进行分析。通过结合 BERT 模型，可以对大量的舆情文本进行分析，从而更好地理解公众对某个话题或事件的看法和反应。

此外，我们的作品还使用了其他的统计方法和图表来细致地了解舆情变化和

趋势。例如，热力图可以直观地展示不同地区或时间段的舆情热度；折线图和条形图则可以清晰地反映舆情的变化趋势和关键点。通过结合不同类型的图表，可以更加全面地了解舆情，为政策制定和决策提供有力的支持。

因此，我们的作品不仅仅是一个基于 BERT 模型的文本分析工具，而是一个综合性的数据分析工具，为政府、企业、媒体等机构提供了更好的了解公众需求和意见的途径，同时也为研究者提供了有价值的参考和支持。可以进一步拓展应用场景，例如在政策制定、市场营销、品牌塑造等方面发挥重要的作用。

### 1.5 应用效果



## 第 2 章 问题描述

### 2.1 问题来源

【填写说明：说明问题的背景、起因等】

互联网时代快速发展以来，网民数量不断扩大，根据第五十次《中国互联网发展状况统计报告》，截至 2022 年 6 月，我国网民规模为 10.51 亿，互联网普及率达 74.4%。更由于互联网的发声门槛低的特性，人们常常更加愿意在网络上发表自己的意见而非传统渠道。因此，网络舆论已经变成社会舆论非常重要的一部分了。在此情况下，对网络舆论的分析就显得格外重要。

### 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】

目前大多数对情感倾向性分析的解决办法，都是首先对抓取的舆情数据进行分词处理，然后结合情感语料数据库和情感分析算法对切分后的语料进行情感计算、分析，并进行情感标注。最后通过聚类 and 分类得出个体情感倾向和群体情感倾向。这是一种从传统机器学习和计算语言学发展而来的分析方法。

### 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】

2.2 中提到的现有的大多数情感倾向分析方法有几个问题：

- 1、对汉语的切割问题一直是一个老大难问题，传统机器学习方法效果并不是很好；
- 2、情感数据库的构建比较耗费人力，同时也比较麻烦；
- 3、由于 1、的问题，通过聚类和分类对情感进行分析的方法效果在现在看来并不是很好。

## 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】

### 作品功能和性能需求：

1. 模型的整个运作过程中，要尽量减少人力消耗，尽量做到一键式傻瓜操作
2. 模型的情感分析效果要比较好。由于是大规模的数据分析，可以接受的平均准确率可以比人稍低，但是不能低太多。

### 数据集：

训练和测试数据集为互联网上开源的外卖评论数据集。总量有 10000 条。格式为.csv 文件：

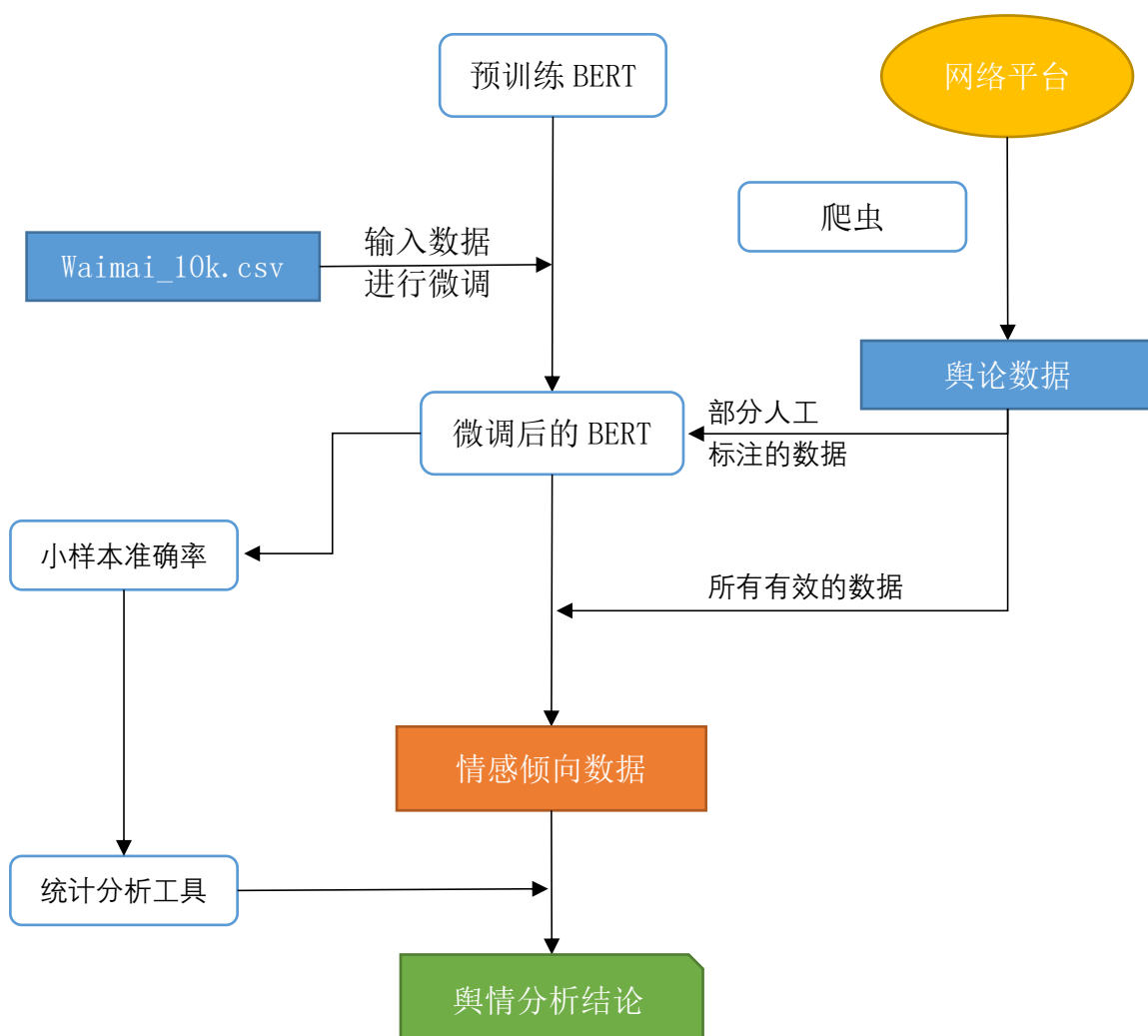
情感倾向（0 或 1），文本内容。

比如，有一条负面评价为：

0,多收了钱。

## 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、算法等；原创作品详细描述，非原创作品简略描述，并尽可能标注引用文献】



## 1、微调 BERT 预训练模型

BERT (Bidirectional Encoder Representations from Transformers) 是一种基于 Transformers 架构的预训练语言模型，在多个 NLP 任务上取得了最先进的性能。

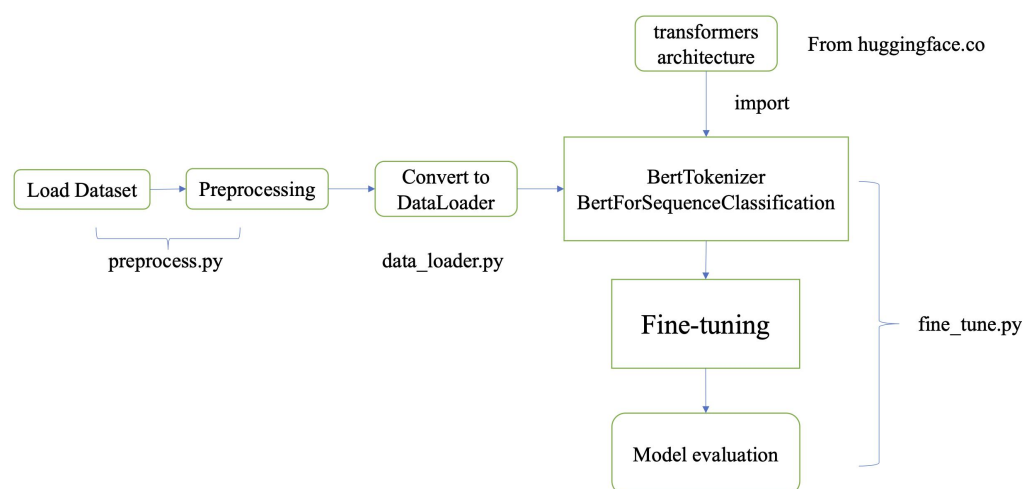
BERT 通过预先训练一个大型的无监督模型，使其能够理解自然语言中的上下文和语义信息。BERT 模型包括两个阶段的训练：预训练和微调。在预训练阶段，BERT 模型使用大量的未标记的文本数据来学习语言知识。在微调阶段，BERT 模型使用标记的数据进行进一步的训练，以完成特定的任务。

bert-base-chinese 是 BERT 的一个预训练模型，专门针对中文语言处理任务。它使用中文维基百科和百度百科作为预训练语料库进行训练，包含了 21128 个字形的汉字和 21128 个对应的字形的拼音。bert-base-chinese 的预训练模型在中文 NLP 任务上表现出色，因此在中文自然语言处理领域中广受欢迎。

使用 bert-base-chinese 模型可以进行中文文本分类、序列标注、语言理解

等多种 NLP 任务。BERT 的出现和 BERT-based 模型的应用对于自然语言处理领域是一次重要的里程碑，因为它们在许多 NLP 任务上都取得了最先进的结果，为自然语言处理研究的未来提供了新的方向。

实际模型微调过程的代码架构如下图所示。由三个文件组成，`preprocess.py` 用于载入本地的 SNLI 数据集文件（包括训练集、验证集和测试集）并对数据进行预处理；`data_loader.py` 文件用于将处理后的数据集进行 `tokenization` 处理，封装构成 `Dataset` 对象（来自 `torch.utils.data`）；`fine_tune.py` 用于模型的参数设置和微调训练。（后为简化程序，将 `preprocess` 和 `data_loader` 程序合并为 `dataloader.py` 并添加了 `prediction.py` 来对实际数据（爬取得到）进行预测。



## 1.1 数据预处理与 DataLoader 生成

将准备好的数据以 `.csv` 格式读入，数据格式应为 `[review, label]` 的 `[文本, 标签]` 对。接着用几种方法进行清洗和预处理：1) 删除文本或标签为空的无效数据；2) 将长度超过限制的文本进行裁剪，具体方法为：以 512 字为单位（`tokenizer` 所接受的字符串长度上限为 512）对原文本串进行分割，将其分成多个标签相同的子串并添加至数据集中。

`DataLoader` 的生成在 `dataloader.py` 中实现。在为了将 SNLI 数据输入预训练模型，需要先将其转化为模型接受的 `torch.utils.data.Dataset` 对象。因此需要对 `Dataset` 类进行重载。重载的类需要实现三个成员函数，分别为 `self.__init__`、`self.__len__` 和 `self.__getitem__`。`getitem` 函数中需要将列表存储的数据进行 `tokenize`，即调用 `tokenizer`。本实验使用的是 `transformer` 库的 `BertTokenizer`，并使用 `encode` 函数对句子中的每个词进行编码。

## 1.2 fine-tuning 微调

实验使用 `transformer` 库的 `BertForSequenceClassification` 中的



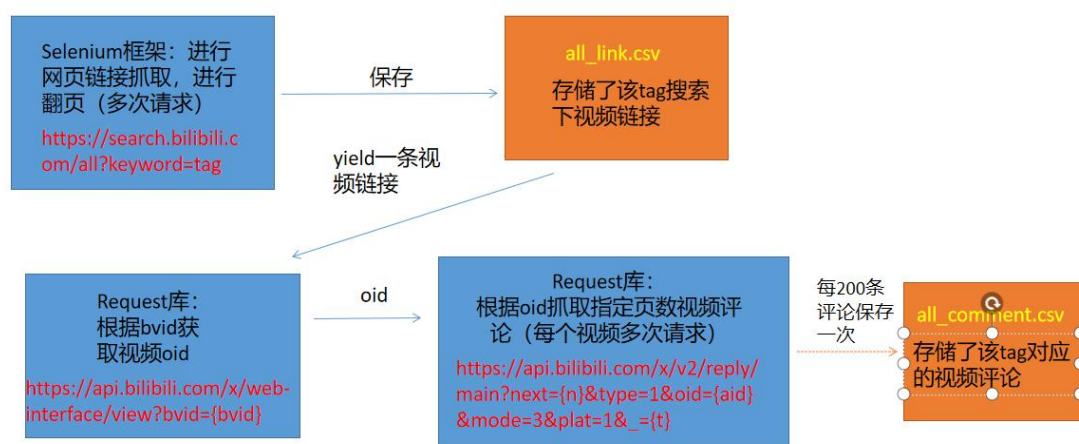
`bert-base-chinese` 预训练模型进行微调，`AdamW` 作为随机优化器，初始学习率设置为 `2e-5`，权重衰减设置为 `0.01`，循环次数设置为 `5`，损失函数为交叉熵损失（`CrossEntropyLoss`）。

## 2、爬虫

爬取目标：抓取 b 站按“人工智能”，“防疫”两个搜索关键词的视频下的评论各 10000 条。

主要功能：包括获取各子模块视频链接、获取视频评论数据，以及将数据保存到本地文件中。这些数据可以用于用户画像分析、情感分析等应用。

技术总览：使用 `selenium` 框架和 `geckodriver` 驱动连接火狐浏览器，模拟浏览器操作。模拟打开 b 站，输入 tag 后，使用 `xpath` 来定位视频链接，再使用 `api.bilibili.com` 获取视频链接的 `oid`，再根据 `oid` 接取 `https://api.bilibili.com/x/v2/reply` 中相应视频获得该视频评论。



代码细节：

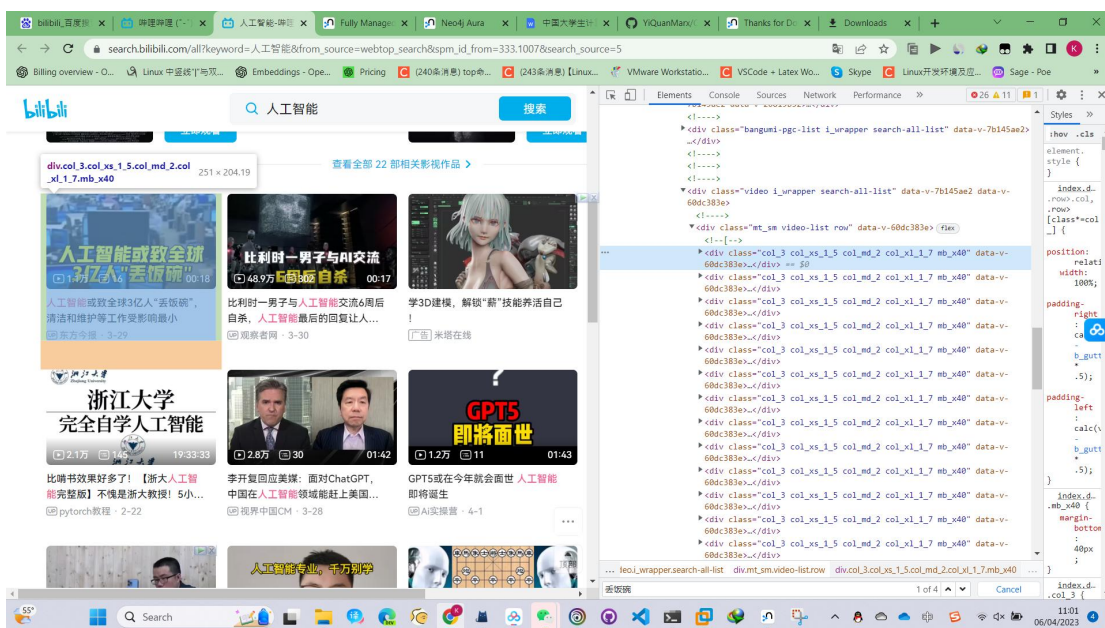
根据 b 站 url 特点，在 `keyword` 字段输入想要查找的 tag，相当于在搜索框搜索 tag，

```
url_base=f"https://search.bilibili.com/all?keyword={tag}&from_source=webtop_search&spm_id_from=333.1007&search_source=5"
```

`page` 模拟翻页过程，通过翻页不断查找新的页面。可以控制翻页数量控制要获取的视频总数。

```
url=url_base+f"&page={i}"
```

进入某一特定页后，通过 `xpath` 定位视频链接元素，如图所示，视频链接全部集中在 `div[@class='mt_sm video-list row']` 里，每个视频 `div` 下直接获取 a 标签得到视频链接地址。



在 b 站官方 api 里，输入每个视频链接的 bvid，得到每个视频的 aid，  
<https://api.bilibili.com/x/web-interface/view?bvid={bvid}>  
 再在官方 api 里，输入每个视频链接的 bvid，得到每个视频的 aid

每个 tag 抓取 30 页的视频。其中“人工智能”tag 抓取了 893 个视频链接，“防疫”tag 抓取了 728 个视频链接。

每个视频抓取 3 页的评论数据。平均一个视频能抓取 30-50 个评论数据。当每个 tag 的评论总数超过 10000 时，自动关停抓取，并将其按照[评论，时间，点赞数，人工标签，人工有效判断，模型标签]6 列打包成 csv 文件。

### 3、统计分析

设  $X_i$  是第  $i$  个数据是否预测正确的变量，设  $p$  是其正确率。即

$$P(X = 1) = p \quad P(X = 0) = 1 - p$$

- 通过给模型喂一批小量的（数量为  $n$ ）数据样本，可以得到模型在这一批数据上的正确频率  $\bar{X} = \sum X_i / n$ 。使用点估计的做法，用此频率作为无偏估计其正确率  $p$ ，即  $p = \bar{X}$ 。由于未知方差，则用样本方差来估计：

$$\sigma^2 = S^2 = \sum (X_i - \bar{X})^2 / (n - 1)$$

- 在得到所有数据的分类情况以后，根据上文预估出来的准确率  $p$ ，我们可以得到一些近似的概率：

$$P(\text{分类错误}) = 1 - p$$

因此，如果设数据总量为  $n$ ，分类为正向情感的数量和负向情感的数量分别有  $n_1, n_2$ ，则实际上分类为正向的数量和负向的数量为

$$m_1 = pn_1 + (1-p)n_2, \quad m_2 = pn_2 + (1-p)n_1$$

。并且，由于正向和负向的地位是等同的，两者的方差都应该是上文提到的 $\sigma^2$ 。

- (3) 根据上文得到 $m_1, m_2$ ，我们可以分别得到正向和负向的舆论在这整个舆论场中所占比例。根据所占比例，我们可以分析出现在的舆论场是一边倒还是两边势均力敌。
- (4) 根据上文得到 $m_1, m_2$ ，我们可以用现实中发生的事件来对应各个时间段的舆论变化。
- (5) 我们可以再设一个随机变量 $Y$ 。当第 $i$ 个数据的情感倾向是正向的时候， $Y_i = 1$ ，反之 $Y_i = 0$ 。则有：

$$\bar{Y} = \sum Y_i$$

- (6) 显然， $Y$ 服从伯努利分布。设：

$$P(Y_i = 1) = p_1, P(Y_i = 0) = 1 - p_1$$

则对于 $\forall j(0 \leq j \leq n), P(\bar{Y} = j) = C_n^j p_1^j (1 - p_1)^{n-j}$ ，即 $\bar{Y} \sim B(n, p_1)$ ， $\bar{Y}$ 服从二项分布，且其值代表了量化后的网络舆论对某事件的看法的情感倾向性。

- (7) 因此，我们有必要对 $\bar{Y}$ 的期望和方差进行区间估计。我们设区间估计的置信度水平 $\alpha = 0.05$ 。同时设变量 $Z = \bar{Y} \sim B(n, p_1)$ 。
  - a) 由于我们选择的 $n$ 比较大（ $n \geq 10000$ ），则根据中心极限定理的例<sup>1</sup>，我们可以用正态分布 $N(\mu, \sigma^2)$ ， $\mu = np_1, \sigma^2 = np_1(1 - p_1)$ 来近似表示这个二项分布，即 $Z \sim N(\mu, \sigma^2)$ 。而正态分布的值可以直接查表，这使得计算方便。
  - b) 首先找一个期望 $\mu$ 的良好的点估计。我们将数据分为 $m$ 组，每一组有 $\frac{n}{m} = k$ 个数据，则有 $m$ 组样本数据：

$$\forall i(1 \leq i \leq m), Z_i = \bar{Y}_i = \sum_{j=1}^k Y_{ij} \sim N(\mu, \sigma^2)$$

- c) 由于 $\sigma$ 未知， $\sqrt{m}(Z - \mu)/\sigma$ 不合作枢轴变量。选取样本标准差

$$S^2 = \sum (Z_i - \bar{Z})^2 / (m - 1)$$

代替 $\sigma^2$ ，得到 $\sqrt{m}(Z - \mu)/S \sim t_{m-1}$  则有：

$$P(-t_{m-1}(\alpha/2) \leq \sqrt{m}(Z - \mu)/S \leq t_{m-1}(\alpha/2)) = 1 - \alpha$$

改写为：

$$P\left((\mu - St_{m-1}(\alpha/2))/\sqrt{m} \leq Z \leq (\mu + St_{m-1}(\alpha/2))/\sqrt{m}\right) = 1 - \alpha$$

因此，其置信度水平为 $\alpha = 0.05$  区间为：

$$[\widehat{\mu}_1, \widehat{\mu}_2] = [(\mu - St_{m-1}(\alpha/2))/\sqrt{m}, (\mu + St_{m-1}(\alpha/2))/\sqrt{m}]$$

(8) 当其期望 $\mu$ 大于 0.5 时, 我们可以认为网民的态度偏正面; 反之, 则认为网民的态度偏负面。这是一个假设检验的过程, 即:

- a) 有原假设 $H_0: \mu_0 > 0.5$  和备择假设 $H_1: \mu_0 \leq 0.5$ 。现在给予抽样数据 $X_i, 1 \leq i \leq n$ , 以此来对此假设进行检验。
- b) 与区间估计同样的处理方法, 则有

$$\sqrt{m}(Z - \mu_0)/S \sim t_{m-1}$$

作为枢轴变量。

- c) 假设显著性水平 $\alpha$ 的情况下, 我们得到检验:

$$\Phi: \text{When } \sqrt{m}(Z - \mu_0)/S \geq -t_{m-1}(\alpha), H_0 \text{ accepted.}$$

*Instead, reject  $H_0$ .*

- d) 关于 $m$ 的选择问题。从参考文献<sup>1</sup>可知, 正态分布近似二项分布要求每个样本的大小 $k$ 比较大, 应至少使得 $kp_1 > 5, k(1 - p_1) > 5$ , 所以 $m$ 不应太大, 至少有 $k \geq 100$  才能在大部分情况下满足近似条件。又有:

$$\begin{aligned} S^2 &= \sum (Z_i - \bar{Z})^2 / (m - 1) \\ &= \sum \left( \sum_{j=1}^k Y_{ij} - \sum_{j=1}^n Y_j \right)^2 / (m - 1) \\ &= \frac{1}{m - 1} \sum \left( \sum_{j=1}^n Y_j - \sum_{j=1}^k Y_{ij} \right)^2 \\ &= \frac{1}{m - 1} \left( m \left( \sum_{j=1}^n Y_j \right)^2 - 2 \left( \sum_{j=1}^n Y_j \right) + \sum_{i=1}^m \left( \sum_{j=1}^k Y_{ij} \right)^2 \right) \\ &= \frac{1}{m - 1} \left( (m - 2) \left( \sum_{j=1}^n Y_j \right)^2 + \sum_{i=1}^m \left( \sum_{j=1}^k Y_{ij} \right)^2 \right) \\ &\leq \frac{1}{m - 1} \left( (m - 2) \left( \sum_{j=1}^n Y_j \right)^2 + m \left( \sum_{j=1}^n Y_j \right)^2 \right) \end{aligned}$$

$$= \frac{1}{m-1} \left( 2(m-1) \left( \sum_{j=1}^n Y_j \right)^2 \right)$$

$$= 2 \left( \sum_{j=1}^n Y_j \right)^2 \rightarrow 2\mu^2$$

可见，而当 $m$ 太小，甚至为 2 时，也有

$$S^2 = \sum \left( \sum_{j=1}^{n/m} Y_{ij} \right)^2 \approx \mu^2$$

显然此时标准差距离方差有点远。因此，理想的情况下应该至少有  $m \geq 50$ 。

e) 综上， $m$ 的取值范围为  $50 \leq m \leq 100$ 。具体视情况而定。

(9) 此外，我们最后会根据统计结果，对舆情在时间尺度上进行分析。

## 第 4 章 系统实现

【填写说明：介绍系统的具体实现过程，特别是其中所遇到的困难，解决的方法等，这里只需要介绍团队真实发生的工作】

### 4、预测模型

我们的分析系统中，评论分类模型是最重要的组成部分之一。为了提高模型的预测准确率，我们在微调预训练模型的基础上使用 Additive Margin Softmax 重新设计了损失函数，用于防治过拟合，并使用了一些微调手段来提升模型的效果。

#### 4.1 Additive Margin Softmax (AM Softmax)

在对 huggingface 上的预训练模型进行微调时，默认的损失函数为交叉熵损失 (cross entropy loss)，在训练中我们发现对于使用中文数据集的微调方法，在第 3 个 epoch 开始会产生明显的过拟合现象（即训练集上的损失减小，但验证集上的损失明显上升）。我们推测这可能是由于交叉熵损失函数在进行反向传播过程中在某些数据集上并不能有很好的泛化能力。为此我们实现了 Additive Margin Softmax Loss 作为损失函数，重新进行训练，得到了较好的结果。

AM-Softmax 是一种用于人脸验证或说话人识别等任务的损失函数，它可以提高特征向量之间的角度间隔，从而提高模型的泛化能力。它是在传统的 Softmax 损失函数的基础上，引入了一个额外的边缘项，使得目标类别的 logit 值减去一

个常数  $m$ ，而非目标类别的  $\text{logit}$  值不变。这样，目标类别的概率会降低，而非目标类别的概率会增加，从而增加了分类难度和分类间隔。具体地，AM-Softmax 损失函数可以表示为：

$$L_{AM} = -\log\left(\frac{e^{W_y(x+m)}}{e^{W_y(x+m)} + \sum_{j \neq y} e^{W_j(x+m)}}\right)$$

## 4.2 多层微调 (multi-level fine-tuning)

为了增强模型在情感分类任务上的鲁棒性和泛化能力，我们设计了多层次微调方法：即首先在一个较大的数据集上进行初步预训练（这里的训练数据集可以是二分类，也可以是多分类），接着在我们针对应用场景特殊选择的针对性训练集上进行针对性的微调，使得其既可以在某一特定任务重表现优秀，同时在类似的其他分类问题上也有不错的表现。

## 5、爬虫 (reptile)

在实现过程中，我们遇到了一些困难。例如，在使用 `api.bilibili.com` 获取视频链接的 `oid` 时，我们需要提取 `json` 格式的数据，并使用 `jsonpath` 库来获取相应的 `oid`。此外，在爬取评论数据时，我们需要注意防止被 B 站反爬虫系统封禁 IP，因此需要设置一定的时间间隔和请求头信息。

为了解决这些问题，我们不断阅读官方文档和相关资料，同时也进行了大量的试错和调试工作。最终，我们成功地爬取了 B 站科技模块下各子模块的视频评论数据，并将数据保存到本地文件中。

## 6、统计分析

在实现过程中，我们发现实际上评论的有效率非常低，只有 30% 不到。这会导致大量无用的信息使得我们的分析被误导。因此，在实际实现过程中，我们会不直接使用“评论”（0 或 1）来代表负向或正向，而是使用“评论”（-1 或 1）\* “点赞数量”来代表正向或者负向的舆论大小。

# 第 5 章 系统评测

### 系统正确性：

我们使用了三个中文数据集进行了测试，分别是：  
ChnSentiCorp

NLPCC2014

Sina Weibo

对于每个数据集，我们使用了准确率和 F1 值两个指标进行评估。下表列出了我们的模型在这三个数据集上的表现：

数据集	准确率	F1 值
ChnSentiCorp	0.903	0.900
NLPCC2014	0.829	0.825
Sina Weibo	0.871	0.868

可以看到，我们的模型在所有数据集上都取得了较好的表现，特别是在 ChnSentiCorp 数据集上表现最好，取得了 0.903 的准确率和 0.900 的 F1 值。这表明我们的模型在实际应用中具有较好的泛化能力和鲁棒性，可以满足大多数情况下的情感分析需求。

## 数据准确性

根据 tag 抓取评论时，有些评论有可能跟 tag 无关，如在 tag 为“人工智能”时，如评论“666”，“明天想出去吃早餐”等与人工智能主题无关。虽然这些评论数据没有直接与 tag 相关，但这些数据也表达了人们在该 tag 下评论的情感，因此我们在模型标注时并没有剔除这些数据。为了验证数据集的直接相关性，我们才用小样本人工标注的方式估计数据的有效性。我们人工标注了 1000 个数据是否与样本有关，“人工智能”，“防疫”分别结果如下

人工智能 防疫

有关 351 339

无关 649 661

我们可以看出，数据直接有关率大概在 35%。

## 模型速度：

我们使用 ChnSentiCorp 进行模型训练时，租用了 featurized 的集成 4090 显卡的服务器进行模型训练，当设 batch\_size=64 时，平均一个 epoch 训练时间为 2 分钟

```
- This IS NOT expected if you are initializing BertForSequenceClassification from the checkpoint of a model that you expect to be exactly identical (initializing a BertForSequenceClassification model from a BertForSequenceClassification model).
Some weights of BertForSequenceClassification were not initialized from the model checkpoint at .cache/huggingface/hub/models--bert-base-chinese/snapshots/84b432f646e4047ce1b5db001d43a348cd3f6bd0 and are newly initialized: ['classifier.bias', 'classifier.weight']
You should probably TRAIN this model on a down-stream task to be able to use it for predictions and inference.
Loss: 0.071918: 100% | 338/338 [02:09<00:00, 2.61it/s]
100% | 38/38 [00:04<00:00, 7.61it/s]
[001/005] Train loss: 0.336548 | Val loss: 0.281251 Acc: 0.900751 Precision: 0.900102 Recall: 0.900751
Best model saved(accuracy: 0.9007506255212677)
Loss: 0.056076: 100% | 338/338 [02:08<00:00, 2.63it/s]
100% | 38/38 [00:05<00:00, 7.58it/s]
[002/005] Train loss: 0.263821 | Val loss: 0.300995 Acc: 0.900751 Precision: 0.900614 Recall: 0.900751
Loss: 0.048058: 100% | 338/338 [02:08<00:00, 2.64it/s]
100% | 38/38 [00:04<00:00, 7.63it/s]
[003/005] Train loss: 0.238887 | Val loss: 0.317345 Acc: 0.887406 Precision: 0.891527 Recall: 0.887406
Loss: 0.037905: 100% | 338/338 [02:08<00:00, 2.63it/s]
100% | 38/38 [00:04<00:00, 7.61it/s]
[004/005] Train loss: 0.225315 | Val loss: 0.291229 Acc: 0.914095 Precision: 0.913630 Recall: 0.914095
Best model saved(accuracy: 0.914095079232694)
Loss: 0.040544: 100% | 338/338 [02:08<00:00, 2.64it/s]
100% | 38/38 [00:04<00:00, 7.60it/s]
[005/005] Train loss: 0.207097 | Val loss: 0.298996 Acc: 0.912427 Precision: 0.911937 Recall: 0.912427
(base) → data [ ]
```

当我们对 b 站抓取数据进行模型标注时，使用服务器配置同上，平均标注一万条评论用时也是 2 分多钟，可见我们所使用模型还是比较快速的。

## 第 6 章 安装使用

### 环境

操作系统: windows10/11、MacOS

Python: 3.6 及以上

numpy==1.21.6

pytorch-pretrained-bert==0.6.2

pytorch-transformers==1.1.0

bert-tensorflow==1.0.1

tensorflow==2.11.0

tensorflow-datasets==2.1.0

tensorflow-estimator==2.11.0

tensorflow-hub==0.7.0

tensorflow-intel==2.11.0

tensorflow-io-gcs-filesystem==0.31.0

tensorflow-metadata==0.21.1

### 使用方法

训练: `fine-tune.py`

```
conda run -n base --no-capture-output --live-stream python [Directory of the finetune.py file]
```

`--device string`:指定使用哪个设备进行训练

`--if_local bool`:是否使用本地预训练模型, 如果为否, 则从 `huggingface` 下载模型

`--model_name string`:指定使用的模型名称

`--epochs int`:迭代次数

`--batch_size int`:单轮训练的样本数

`--weight_decay float`:学习率

`--drop_prob float`: dropout 概率, 防止过拟合

标注: `prediction.py`

```
conda run -n base --no-capture-output --live-stream python [Directory of the finetune.py file]
```

`--model_path string`:要使用的模型的路径

`--dataset_name string`:指定要标注的数据集路径



## 第 7 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】

我们的作品是基于大数据分析的情感分类系统。通过对大规模的网络数据进行爬取、清洗、分类和分析，我们可以有效地分析舆情的情感分布和趋势变化，为政府、企业和个人提供科学的决策支持和舆情监测服务。

### 创意

我们的创意主要体现在以下几个方面：

1. 在**统计分析**方面，我们使用了二项分布的期望和方差进行区间估计，然后使用中心极限定理近似二项分布，得到了正态分布的估计值，从而进行了置信区间的计算。同时，我们还使用了假设检验来验证假设的真实性，使用 **t** 分布作为枢轴变量进行计算。
2. 在**系统实现**方面，我们的分析系统中，评论分类模型是最重要的组成部分之一。为了提高模型的预测准确率，我们在微调预训练模型的基础上使用 Additive Margin Softmax 重新设计了损失函数，用于防治过拟合，并使用了一些微调手段来提升模型的效果。同时，我们还设计了多层次微调方法，使得模型在重点任务和其他分类问题上都能表现出色。
3. 在**作品总结**方面，我们强调了我们的作品是基于大数据分析的情感分类系统，可以为政府、企业和个人提供科学的决策支持和舆情监测服务。我们的创意主要体现在统计分析和系统实现两个方面，这也是我们未来进一步提升和应用拓展的方向。

### 技术路线

我们的技术路线主要包括以下几个方面：

1. **爬虫：**

1. **爬虫框架选择：**我们采用了 Scrapy 框架实现爬虫，该框架具有高效、灵活、可扩展性强等特点，同时还提供了强大的数据处理和存储能力。
  2. **数据源选择：**我们选择了 bilibili 作为数据源，通过对 bilibili 的 API 的调用，我们可以获取到大量的评论。
  3. **数据清洗：**由于 bilibili 数据量庞大，其中存在大量的垃圾信息和重复信息，因此我们需要对数据进行清洗。具体来说，我们采用了正则表达式、关键字匹配等方法对数据进行过滤和去重。
2. **文本情感分类：**我们的基于 BERT 的情感分类模型采用了预训练的 BERT 模型作为特征提取器，并在其基础上添加了几个全连接层作为分类器。为了优化模型，我们使用了 Additive Margin Softmax (AM Softmax) 作为损失函数，以防止过拟合。我们还使用了多层微调方法，以提高模型的鲁棒性和泛化能力。在针对中文数据集的微调中，我们还发现交叉熵损失函数可能会导致过拟合，因此使用 AM Softmax 损失函数取得了更好的效果。经过测试，我们的模型在情感分析数据集上取得了较好的分类效果。
3. **统计分析：**在统计分析方面，我们采用了二项分布的期望和方差进行区间估计，然后使用中心极限定理近似二项分布，得到了正态分布的估计值，从而进行了置信区间的计算。这一方法具有较高的精度和可靠性，可以有效地对数据进行分析 and 预测。同时，我们还使用了假设检验来验证假设的真实性、使用区间估计来预测数据的分布情况。我们使用  $t$  分布作为枢轴变量进行计算。这一方法可以帮助我们判断样本数据与总体数据之间的差异是否显著，从而得出结论。通过以上两种方法的应用，我们可以更加深入地了解数据的分布和规律，为后续的分析 and 决策提供科学的依据。

## 工作量

我们的工作主要体现在以下几个方面：

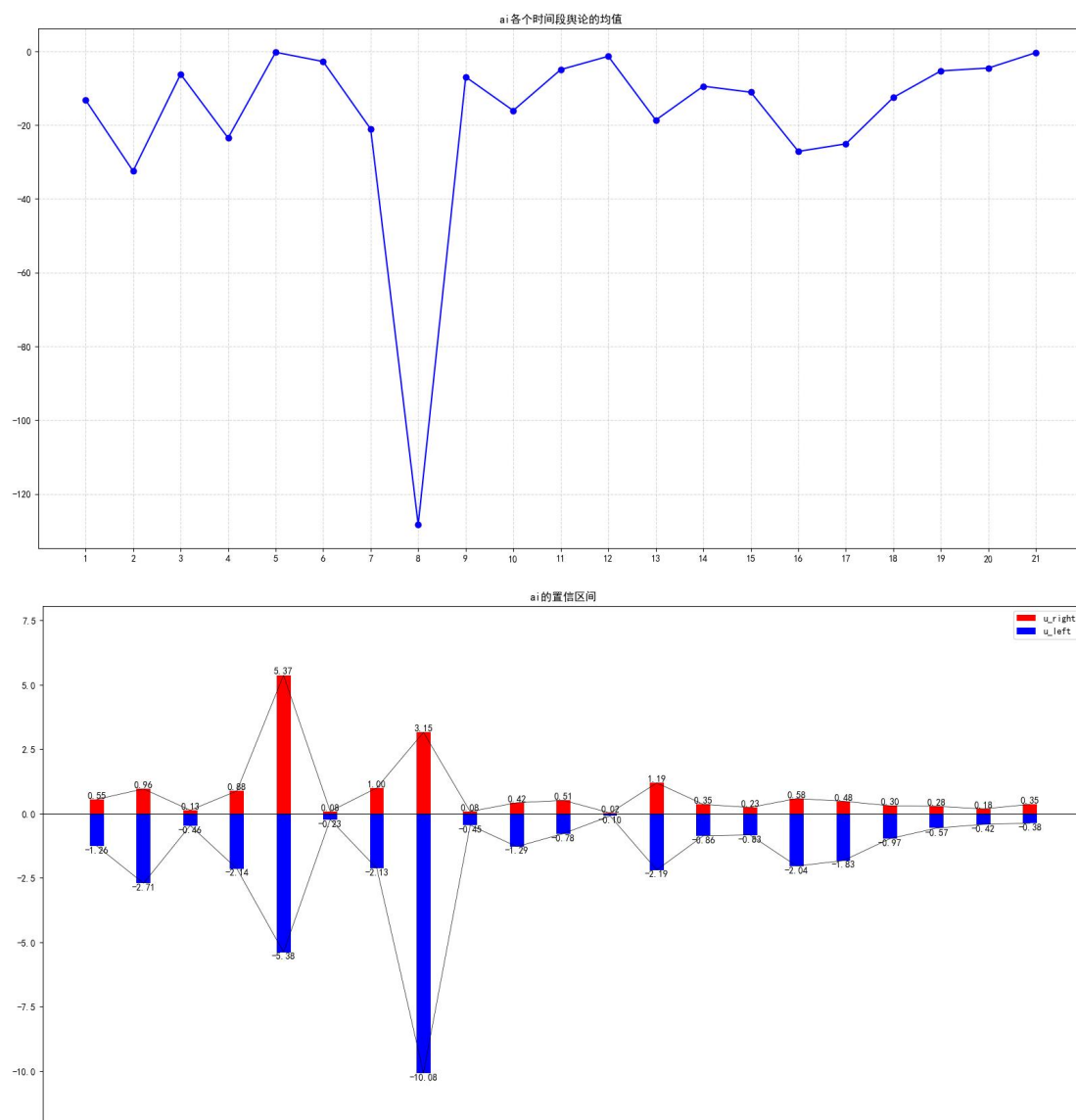
1. 爬虫：我们使用了 selenium 框架实现爬虫，并对数据进行了清洗和去重。这一过程需要花费大量的时间和精力，同时还需要具备一定的编程和数据处理能力。
2. 文本情感分类：我们的情感分类模型采用了 BERT 作为特征提取器，并在其基础上添加了几个全连接层。为了优化模型，我们使用了 Additive Margin Softmax Loss 作为损失函数，以防止过拟合，并使用了多层微调方法。这一过程需要对模型进行不断的调试和优化，同时需要具备一定的机器学习和深度学习基础。
3. 统计分析：我们采用了二项分布的期望和方差进行区间估计，使用了中心极限定理近似二项分布，得到了正态分布的估计值。这一过程需要对统计学和概率论有较深入的了解，同时需要编写一些程序进行计算和分析。

总的来说，我们的工作量比较大，需要具备一定的编程和数学基础，同时需要进行大量的实验和测试。

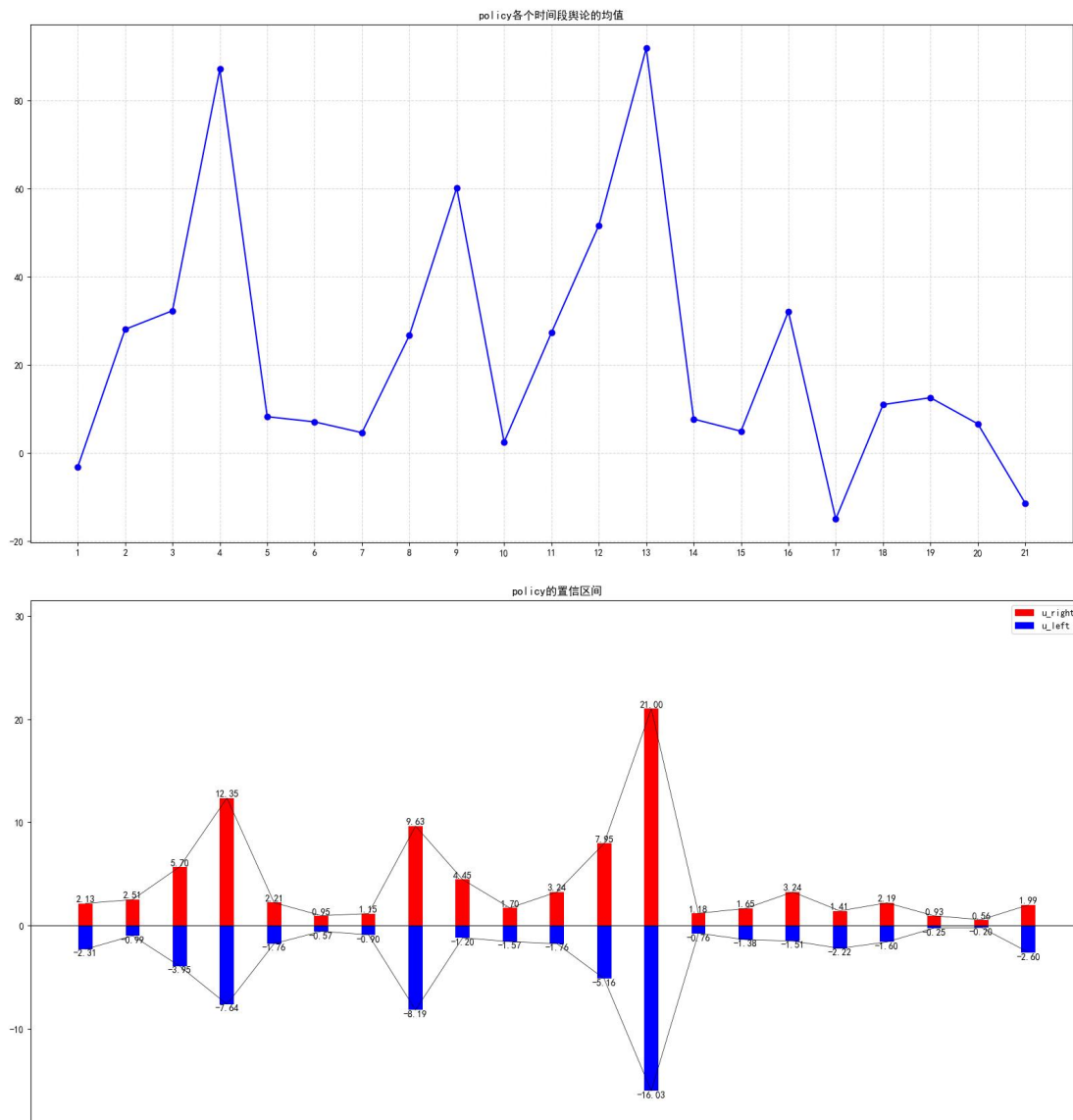
## 数据和测试效果

下面几张图分别展示了将 2022.9.1 到 2023.3.1 之间的分成 20 个时间段后，对“ai”和对“防疫政策”的的舆论的倾向性的均值和在 95%置信度水平下的置信区间的变化。

- 临近 chatGPT 发布的时间段，即 2022 年 11 月中旬到 12 月初，舆论出现了明显的且幅度非常大的下滑。同时这段时间的置信区间比较长，这也就说明人们在这段时间内的意见分歧性大且讨论力度也大。热度过去了就慢慢回复平稳了。



- 总体来看，防疫政策的舆论热度要比 ai 的热度高得多，网民舆论之间的分歧性也会大得多。比如大概在 2022 年 10 月末期的时候有一波大量上涨，同时置信区间也拉得比较长。另外就是大概在 2022 年的 11 月末 12 月初，这件事在舆论上的热度非常之高，评论数和点赞数非常高。最后是大概在 12 月中旬的时候，又一波舆论的高潮。这次的置信区间的是最长的，可以看出舆论场上的意见分歧比较大。



## 进一步提升和应用拓展

为了进一步提升我们的作品，我们可以从以下几个方面进行拓展和优化：

1. **多语言支持**：我们可以考虑支持更多的语言，包括英语、法语、日语等，以满足用户在全球范围内的需求。
2. **多领域应用**：我们可以将我们的系统应用于更广泛的领域，如金融、医疗、教育等，为不同行业提供有针对性的决策支持和舆情监测服务。

3. **可解释性和实时性**：我们可以进一步提高系统的可解释性和实时性，使用户更好地理解 and 应对舆情变化，及时进行舆情干预和管理。

总之，我们的作品是基于创新的思维和先进的技术实现的，具有重要的应用价值和推广意义。我们将继续努力，不断完善和拓展我们的舆情分析系统，为用户提供更好的服务和体验。

## 参考文献

【请按照标准参考文件格式填写】

---

<sup>1</sup> 陈希孺 概率论与数理统计[M]. 安徽：中国科学技术大学出版社，2009:145.