

ECNU

ECNU

基于Twitter的评论爬取分析

Trump2024美国总统大选成功的概率

POWERPOINT DESIGN

罗文琦 刘子阳

2024.12.6



CONTENTS

01 项目概述

02 数据获取

03 数据处理

04 数据分析

05 预测分析

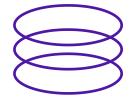
06 项目开源

07 结论

01

项目概述

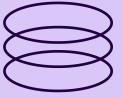
Project Overview



项目背景



分析Twitter评论数据，
预测Trump2024年美国总统大选成功的概率。



项目目标

Step1

爬取Twitter
关于Trump和Harris
大选的推文、评论内容。

Step2

数据预处理
正则表达式+dropna

Step3

Spacy模型
关键词抓取
Vader模型
情感分析

Step4

LSTM深度学习
推文的情绪和热度
预测大选结果

Step5

可视化
数据大屏+词云+聚类

02

数据获取

Data Acquistion

数据获取

The screenshot shows the Spyder IDE interface with two main panes. The left pane displays the Python code for a class named DataGatherer, which interacts with the praw library to fetch posts from specific subreddits and their comments. The right pane shows a CSV file named 'Sorted_data.csv' containing a large dataset of tweets, with columns including TweetID, Weekday, Hour, Day, Lang, IsReshare, Reach, RetweetCount, Likes, Klout, Sentiment, text, LocationID, UserID, and various timestamp and numerical values. The CSV file has 36 rows.

```
import praw
class DataGatherer:
    def __init__(self, client_id: str, client_secret: str, subreddit_names_list: list[str], maximum_posts_per_subreddit: int, top_posts_time_filter: str):
        self.subreddit_names_list = subreddit_names_list
        self.maximum_posts_per_subreddit = maximum_posts_per_subreddit
        self.top_posts_time_filter = top_posts_time_filter
        self.reddit_client = praw.Reddit(client_id=client_id,
                                         client_secret=client_secret,
                                         user_agent="user_agent")
        self.comments_list_by_author = {}

    def get_comments_list_by_author(self, subreddit_name: str) -> dict[str, list[str]]:
        retrieved_posts_count = 0

        for subreddit_name in self.subreddit_names_list:
            subreddit = self.reddit_client.subreddit(subreddit_name)
            top_posts = subreddit.top(time_filter=self.top_posts_time_filter)

            current_subreddit_retrieved_posts_count = 0

            for post in top_posts:
                post.comments.replace_more(limit=0)

                for comment in post.comments.list():
                    comment_author = comment.author
                    if comment_author and comment_author.name != "AutoModerator":
                        if comment_author not in self.comments_list_by_author:
                            self.comments_list_by_author[comment_author.name] = [comment.body]
                        else:
                            self.comments_list_by_author[comment_author.name].append(comment.body)

            retrieved_posts_count += 1

        return self.comments_list_by_author
```

Sorted_data.csv

TweetID	Weekday	Hour	Day	Lang	IsReshare	Reach	RetweetCount	Likes	Klout	Sentiment	text	LocationID	UserID		
tw-68271287332805633	Thursday	17	0	31	0	en	0	0	44	0	0	0	35	0	0
tw-682713045357998880	Thursday	17	0	31	0	en	0	0	282	0	0	0	47	0	0
tw-682713219375476736	Thursday	17	0	31	0	en	0	0	287	0	4	0	0	55	0
tw-682713436967579648	Thursday	17	0	31	0	en	0	0	2087	0	4	0	0	50	0
tw-682714048199311366	Thursday	17	0	31	0	en	0	0	955	0	0	0	47	0	0
tw-6827145836944243456	Thursday	17	0	31	0	en	0	0	113	0	0	0	31	0	-1
tw-682715526435454977	Thursday	17	0	31	0	en	0	0	113	0	0	0	31	0	-2
tw-6827179679543402753	Thursday	17	0	31	0	en	0	0	113	0	0	0	31	0	-1
tw-6827179679543402753	Thursday	17	0	31	0	en	0	0	4693	0	1	0	0	49	0
tw-6827180448988519168	Thursday	17	0	31	0	en	0	1	2868	0	6	0	0	51	0
tw-682719177115844609	Thursday	17	0	31	0	en	0	1	86	0	51	0	0	16	0
tw-6827192295670892801	Thursday	17	0	31	0	en	0	0	8	0	0	0	0	21	0
tw-68272309279841798	Thursday	17	0	31	0	en	0	0	2214	0	0	0	0	0	0
tw-682724038150623232	Thursday	17	0	31	0	en	0	0	387	0	0	0	0	56	0
tw-6827255225266296832	Thursday	17	0	31	0	en	0	1	149	0	3	0	0	31	0
tw-6827255225266296832	Thursday	17	0	31	0	en	0	1	356	0	2	0	0	47	0
tw-682726320469161025	Thursday	17	0	31	0	en	0	1	201	0	3	0	0	0	1
tw-682726753551413249	Thursday	17	0	31	0	en	0	1	3774	0	8	0	0	49	0
tw-68272895786913152	Thursday	18	0	31	0	en	0	0	3374	0	8	0	0	45	0
tw-682728964389704384	Thursday	18	0	31	0	en	0	0	932	0	2	0	0	0	4
tw-68273025548945409	Thursday	18	0	31	0	en	0	0	1766	0	0	0	0	46	0
tw-68273133698478831360	Thursday	18	0	31	0	en	0	1	631	0	35	0	0	0	42
tw-68273373698482688	Thursday	18	0	31	0	en	0	0	2874	0	1	0	0	59	0
tw-682734158469181440	Thursday	18	0	31	0	en	0	1	104	0	3	0	0	42	0
tw-682735302855782400	Thursday	18	0	31	0	en	0	0	331	0	0	0	0	36	0
tw-6827356848868141720	Thursday	18	0	31	0	en	0	0	281	0	0	0	0	43	0
tw-68273708869409576	Thursday	18	0	31	0	en	0	0	1052	0	0	0	0	47	0
tw-68273841747514945	Thursday	18	0	31	0	en	0	0	488	0	0	0	0	45	0
tw-682739568619601920	Thursday	18	0	31	0	en	0	0	437	0	0	0	0	38	0
tw-682739681105199104	Thursday	18	0	31	0	en	0	1	1642	0	43	0	0	50	0

技术层面

使用python和spyder作为爬取工具。

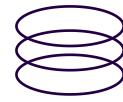


获取推文的评论，通过Tweet ID和API请求。

分批分页处理，应对Twitter对评论数据的爬取限制。

新加卷 (F:) > Twitter弹幕数据分析 > DATA > usc-x-24-us-election			
名称	修改日期	类型	大小
may_july_chunk_1	2024-11-01 16:45	XLS 工作表	81,121 KB
may_july_chunk_2	2024-11-01 16:45	XLS 工作表	81,777 KB
may_july_chunk_3	2024-11-01 16:45	XLS 工作表	81,401 KB
may_july_chunk_4	2024-11-01 16:45	XLS 工作表	81,518 KB
may_july_chunk_5	2024-11-01 16:45	XLS 工作表	81,682 KB
may_july_chunk_6	2024-11-01 16:45	XLS 工作表	82,922 KB
may_july_chunk_7	2024-11-01 16:45	XLS 工作表	80,425 KB
may_july_chunk_8	2024-11-01 16:45	XLS 工作表	80,737 KB
may_july_chunk_9	2024-11-01 16:45	XLS 工作表	81,269 KB
may_july_chunk_10	2024-11-01 16:45	XLS 工作表	82,061 KB
may_july_chunk_11	2024-11-01 16:45	XLS 工作表	82,994 KB
may_july_chunk_12	2024-11-01 16:46	XLS 工作表	83,032 KB
may_july_chunk_13	2024-11-01 16:46	XLS 工作表	81,839 KB
may_july_chunk_14	2024-11-01 16:46	XLS 工作表	82,666 KB
may_july_chunk_15	2024-11-01 16:46	XLS 工作表	80,189 KB
may_july_chunk_16	2024-11-01 16:46	XLS 工作表	81,753 KB
may_july_chunk_17	2024-11-01 16:46	XLS 工作表	82,158 KB
may_july_chunk_18	2024-11-01 16:46	XLS 工作表	82,474 KB
may_july_chunk_19	2024-11-01 16:46	XLS 工作表	82,477 KB
may_july_chunk_20	2024-11-01 16:46	XLS 工作表	82,072 KB

part_1	added part_1
part_10	added part10 - part12
part_11	added part10 - part12
part_12	added part10 - part12
part_13	added part13-part16
part_14	added part13-part16
part_15	added part15-part21
part_16	added part15-part22
part_17	added parts part23
part_18	added parts part24
part_19	added parts part25
part_2	added part26
part_20	added parts part27
	part28
	part3
	part4
	part5
	part6
	part7
	part8



数据拼接

构建爬取数据的集合，将爬取数据的csv文件进行合并
同样使用python进行分批分类处理。

Openrank / test.py

```
1 import pandas as pd
2
3 # 定义文件路径列表
4 file_paths = [
5     r'C:\Users\86189\Desktop\may_july_chunk_1.csv',
6     r'C:\Users\86189\Desktop\may_july_chunk_2.csv',
7     r'C:\Users\86189\Desktop\may_july_chunk_3.csv',
8     r'C:\Users\86189\Desktop\may_july_chunk_4.csv',
9     r'C:\Users\86189\Desktop\may_july_chunk_5.csv',
10    r'C:\Users\86189\Desktop\may_july_chunk_6.csv',
11    r'C:\Users\86189\Desktop\may_july_chunk_7.csv',
12    r'C:\Users\86189\Desktop\may_july_chunk_8.csv',
13    r'C:\Users\86189\Desktop\may_july_chunk_9.csv',
14    r'C:\Users\86189\Desktop\may_july_chunk_10.csv'
15 ]
16
17 # 读取第一个CSV文件
18 df = pd.read_csv(file_paths[0], low_memory=False)
19
20 # 读取其余的CSV文件并拼接
21 for i in range(1, len(file_paths)):
22     df_i = pd.read_csv(file_paths[i], low_memory=False)
23     # 拼接并去除重复的行
24     df = pd.concat([df, df_i]).drop_duplicates()
25
26 # 保留需要的列（根据需求进行修改）
27 # 这里保留所有列（如果只想保留部分列，可以通过 df.iloc[:, :] 来选择）
28 df = df.iloc[:, :] # 如果需要特定列，请修改为这一行
29
30 # 保存拼接后的CSV文件
31 output_file = r'C:\Users\86189\Desktop\csv_merge.csv'
32
33 # 使用 utf-8 编码保存文件，或者使用 'ignore' 忽略无法编码的字符
34 df.to_csv(output_file, index=False, encoding='utf-8') # 使用utf-8编码
```

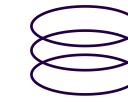
merge.py

```
1 import pandas as pd
2
3 # 创建一个空的DataFrame用来存储合并的结果
4 merged_df = pd.DataFrame()
5
6 # 定义每次读取的块大小
7 chunksize = 10000
8
9 # 循环遍历每个文件，逐块读取并合并
10 for i in range(1, 21):
11     file_name = f'{i}.csv'
12
13     for chunk in pd.read_csv(file_name, chunksize=chunksize):
14         # 如果是第二个及以后文件，去掉第一行（列名），假设第一个文件的列名是正确的且要保留
15         if i > 1 and chunk.index[0] == 0:
16             chunk = chunk.iloc[1:]
17
18         # 将当前块的数据合并到最终的DataFrame中
19         merged_df = pd.concat([merged_df, chunk], ignore_index=True)
20
21 # 保存最终的合并结果
22 merged_df.to_csv(path_or_buf='all_merged_file.csv', index=False)
```

03

数据处理

Data Processing



关键词提取



01

数据筛选

02

使用spacy提取评论中的关键词。

03

做出词云图进行可视化

数据清理后的部分结果

The screenshot shows the PyCharm IDE interface with two tabs open: 'merge.py' and 'data pre-processed.py'. The code in 'merge.py' is as follows:

```
1 import pandas as pd
2 import re
3
4 # 读取数据
5 df = pd.read_csv('pre_data.csv')
6
7 # 保留需要的列
8 df = df[['rawContent', 'replyCount', 'retweetCount', 'likeCount', 'quoteCount']]
9
10 # 处理文本数据
11 def clean_text(text):
12     # 如果文本是缺失值或非字符串类型，则返回空字符串
13     if not isinstance(text, str):
14         return ''
15     # 去除URL
16     text = re.sub(pattern=r'http\S+|www\S+', repl='', text)
17     # 允许标点符号和中文字符
18     text = re.sub(pattern=r'[^\\w\\s.,!?;:]', repl=' ', text) # 保留字母、数字、空格、常见标点符号
19     # 如果文本为空，返回空字符串
20     return text.strip()
21
22 df['cleaned_content'] = df['rawContent'].apply(lambda x: clean_text(x))
23 # 删除包含缺失值的行
24 df.dropna(inplace=True)
25
26 # 输出处理后的数据（可选）
27 df.to_csv('cleaned_data.csv', index=False)
28
29 # 你可以查看数据结构
30 print(df.head())
31
32 # 检查清理后的数据数量
33 print(f"Number of rows after cleaning: {len(df)}")
```

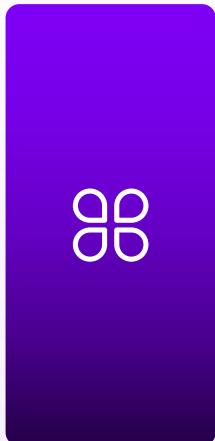


利用正则化表达式去除 URL 并保留字母、数字、空格和常见符号

**利用dropna函数，对不符合
条件的内容做一个清洗**

	rawContent	cleaned_content
1	@lukepbeasley I cant imagine anyo	lukepbeasley I cant imagine anyo
2	Voters can also sway me away from	Voters can also sway me away fro
3	@PoodleHead57 @BobOnderMO Can you	PoodleHead57 BobOnderMO Can you
4	@Morning_Joe @JoeNBC The fact rem	Morning_Joe JoeNBC The fact rema
5	@BidenHQ That's funny you're obvi	BidenHQ Thats funny youre obviou
6	#Internacional	Internacional
7	Las modificaciones introducidas por el grupo palestino Hamás en el acuerdo diseñado por el presidente de EEUU, Joe Biden,	Las modificaciones introducidas por el grupo palestino Hamás en el acuerdo diseñado por el presidente de EEUU, Joe Biden,
8	MAGA RAGES Over Hunter Biden Ver	MAGA RAGES Over Hunter Biden Ver
9	@harryjsisson There is no doubt i	harryjsisson There is no doubt i
10	@ShaunPMaca @MoneyTalkUS	ShaunPMaca MoneyTalkUS
11	@byHeatherLong Hourly wages are irrelevant.	byHeatherLong Hourly wages are irrelevant.
12	@nopolabo_maga à á ã ä å ä ^	nopolabo_maga
13	à á ~ à á í¼ \â (â à' â á á ^ á á)	

利用spaCy进行关键词提取并过滤重要内容

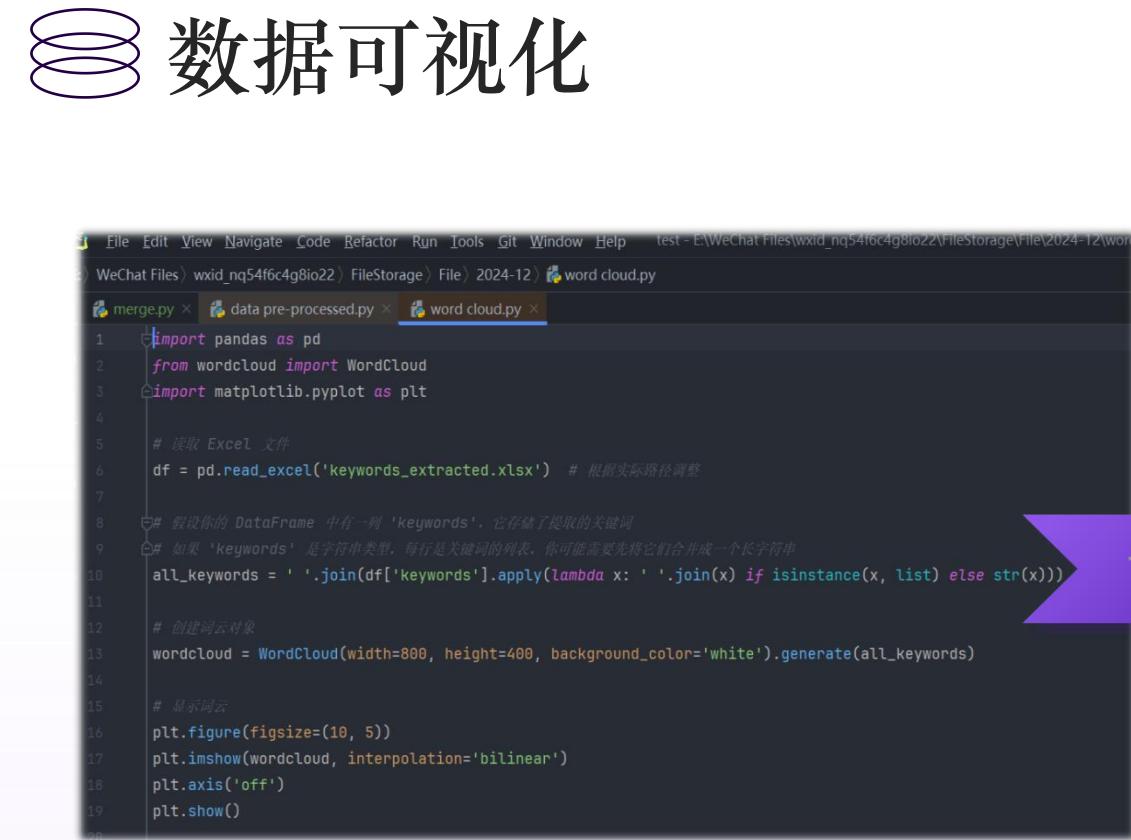


```
# 加载spaCy模型
nlp = spacy.load('en_core_web_sm')

# 提取关键词（例如：人物、组织等）
def extract_keywords(text): 1 usage
    doc = nlp(text)
    keywords = [ent.text for ent in doc.ents if ent.label_ in ['PERSON', 'ORG']]
    return keywords

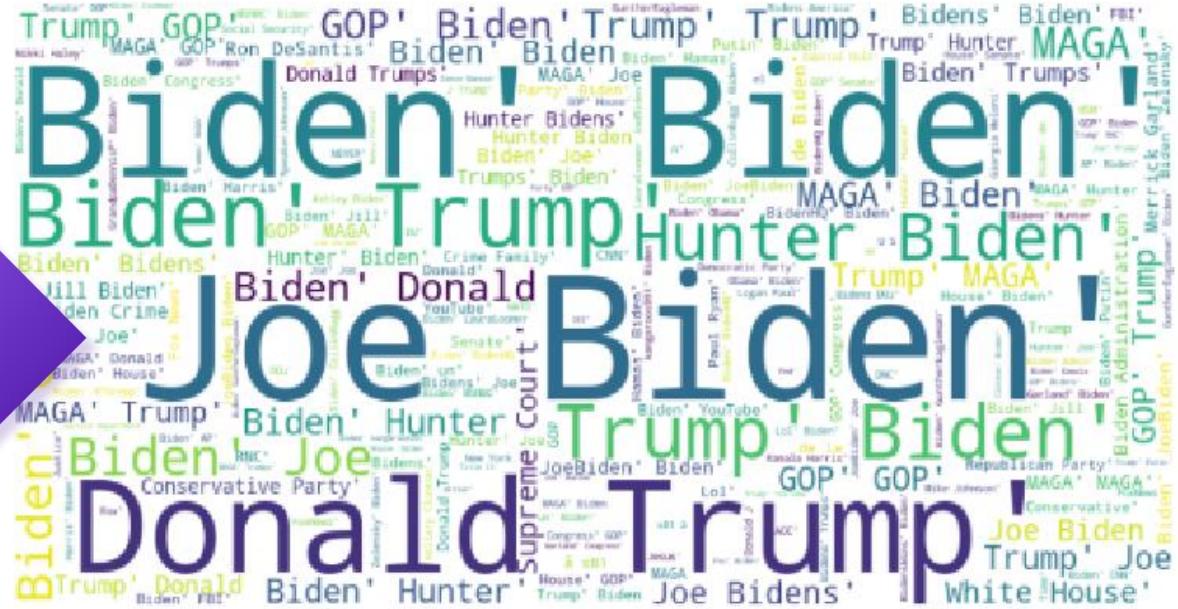
# 提取关键词
df['keywords'] = df['cleaned_content'].apply(extract_keywords)
```

cleaned_content	keywords
lukepbeasley I cant imagine anyone actual	[‘Donald’]
Morning_Joe JoeNEC The fact remains that	[‘Joe Biden’]
BidenHQ Thats funny youre obviously tryin	[‘BidenHQ’]
Internacional Las modificaciones introducid	[‘Joe Biden’, ‘la retirada del’, ‘Franja de Gaza’, ‘semana de la implementaciÃ³n’]
LadyMagaUSA nypost Lady Maga is here to s	[‘Lady Maga’, ‘HOORAY’]
teenburger Fernand46357857 Terrorism is d	[‘Biden’, ‘Hamas’, ‘Biden’]
MAGA hats are for morons racists	[‘MAGA’]
BidenHQ Were still poorer than ever. We r	[‘Drop Biden’, ‘Trump’]



按照关键词统计数据做出词云图

可视化



绘制评论情绪分布的聚类图

04

数据分析

Data Analysis

“情感分析”

01 正则化表达

02 使用情感分析模型VADER对评论的情绪进行分类进行情感分析，将评论分为积极（支持某候选人）、消极（反对候选人）和中立（未表态）

```
# 初始化 VADER 情感分析器
sia = SentimentIntensityAnalyzer()

# 定义情感分类函数
def get_sentiment(text): 1 usage
    # 确保文本是字符串
    if not isinstance(text, str):
        text = str(text) # 将非字符串类型转换为字符串
    score = sia.polarity_scores(text)
    if score['compound'] ≥ 0.05:
        return 'positive'
    elif score['compound'] ≤ -0.05:
        return 'negative'
    else:
        return 'neutral'
```

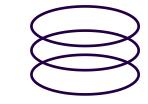
```
print(f"Number of rows after cleaning: {len(df)}
```

正则化以后的分类情况

05

预测分析

Predictive Analysis



模型选择

使用机器学习模型Random Forest

使用历史选举数据和热度趋势作为训练集

“ 趋势分析

利用Random Forest对不同候选人的支持情况进行深度学习并分类
通过训练数据学习情感和关键词的模式，预测不同候选人的支持情况。

```
# 将 'keywords' 和 'candidate' 转化为特征
vectorizer = TfidfVectorizer(max_features=5000)
X = vectorizer.fit_transform(df['keywords']) # 使用 keywords 列进行特征提取

# 使用 candidate 列作为标签，分别为 Trump、Harris 或 Biden
y = df['candidate'].apply(lambda x: 1 if x == 'Trump' else (2 if x == 'Harris' else 0)) # Trump:1

# 划分训练集和测试集
X_train, X_test, y_train, y_test = train_test_split(*arrays: X, y, test_size=0.2, random_state=42)

# 训练 Random Forest 模型
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)
```



```
# 训练 Random Forest 模型
model = RandomForestClassifier(n_estimators=100, random_state=42)
model.fit(X_train, y_train)

# 获取 Trump 和 Harris 的预测概率
probabilities = model.predict_proba(X_test)

# 获取每个样本是 Trump (类别 1) 或 Harris (类别 2) 的概率
prob_trump = probabilities[:, 1] # Trump 类别的预测概率
prob_harris = probabilities[:, 2] # Harris 类别的预测概率

# 计算 Trump 和 Harris 获胜的支持比例
trump_support = np.mean(prob_trump > prob_harris) * 100
harris_support = np.mean(prob_harris > prob_trump) * 100

print(f"Trump Support: {trump_support:.2f}%")
print(f"Harris Support: {harris_support:.2f}%")
```

```
Trump Support: 62.44%
Harris Support: 0.92%
```



通过上图可发现，比例非常的悬殊；所以我们想更细致地考虑两种细分情况。

```
# 使用 candidate 列作为标签，分别为 Trump 或 Biden  
y = df['candidate'].apply(lambda x: 1 if x == 'Trump' else (2 if x == 'Biden' else 0))  
  
# 划分训练集和测试集  
X_train, X_test, y_train, y_test = train_test_split(*arrays: X, y, test_size=0.2, random_state=42)
```

Prediction ×

:

```
D:\Anaconda\python.exe F:\Twitter弹幕数据分析\CODE\Prediction.py  
Trump Support: 24.29%  
Biden Support: 59.64%
```

一个是不考虑Harris加入的情况，
也就是Biden不退选，直接与Trump竞争的结果。如图

Biden drops out of 2024 race after disastrous debate inflamed age concerns. VP Harris gets his nod



另一个是考虑Harris加入，Biden退出的情况。因为Biden在大选前3个月临时退出，交给了副总统Harris，在这样的情况下，如果是对Biden的民主党支持者，大概率也会把选票给到Harris，所以便需要重新调整训练模型。

```
df['candidate'] = df.apply(lambda row: get_candidate(row['sentiment'], row['keywords']), axis=1)
```

经过调参以后得到合理的结果

```
48     # 计算 Trump 和 Harris 获胜的支持比例  
49     # 为了将Biden的支持按1/2的权重计入Harris支持，我们调整prob_harris  
50     prob_harris_adjusted = prob_harris + 0.5 * prob_biden  
51  
52     # 计算 Trump 和 Harris（包括Biden支持）获胜的支持比例
```



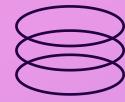
The screenshot shows a Jupyter Notebook interface with a dark theme. At the top, there is a toolbar with icons for Run, Prediction, and other options. Below the toolbar, the code from line 48 to 52 is displayed. The output section shows the results of the code execution:

```
D:\Anaconda\python.exe F:\Twitter弹幕数据分析\CODE\Prediction.py  
Trump Support: 59.38%  
Harris Support: 24.39%
```

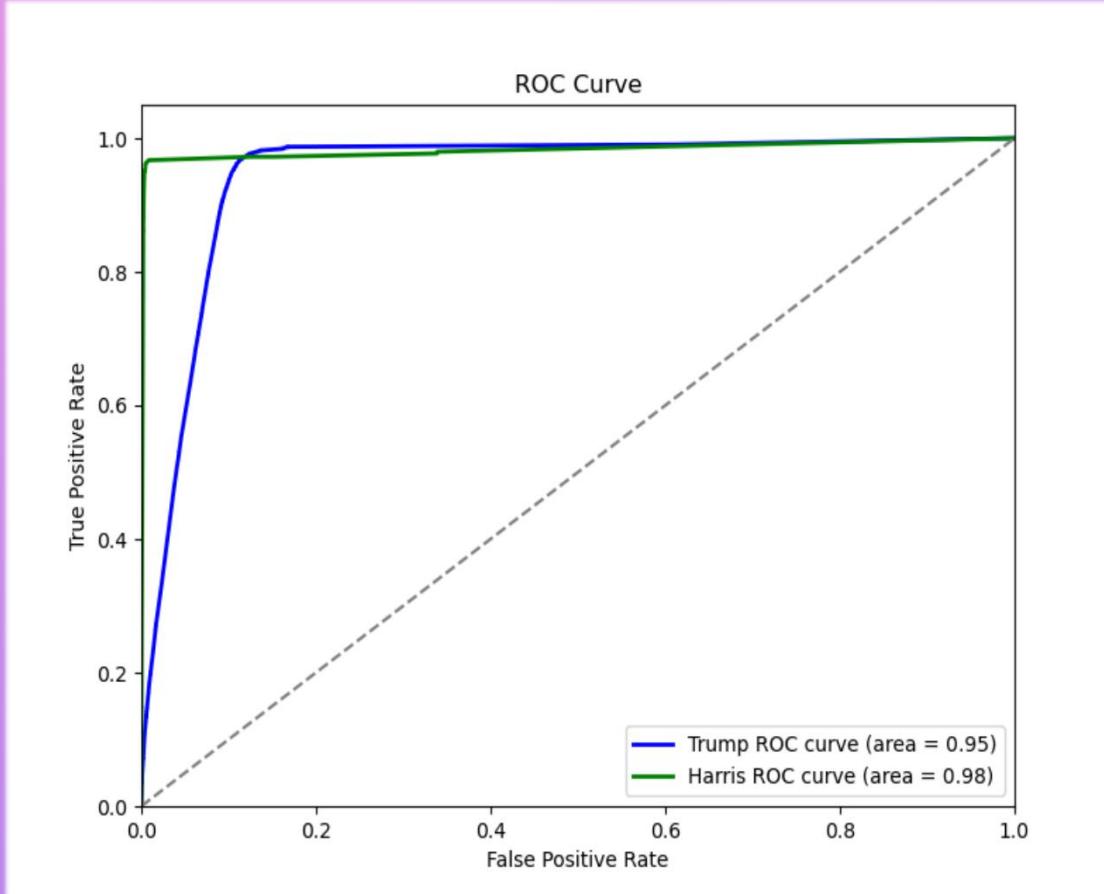
06

模型可行性验证 及可视化

Visualization

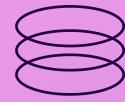


数据可视化



ROC曲线 ROC曲线

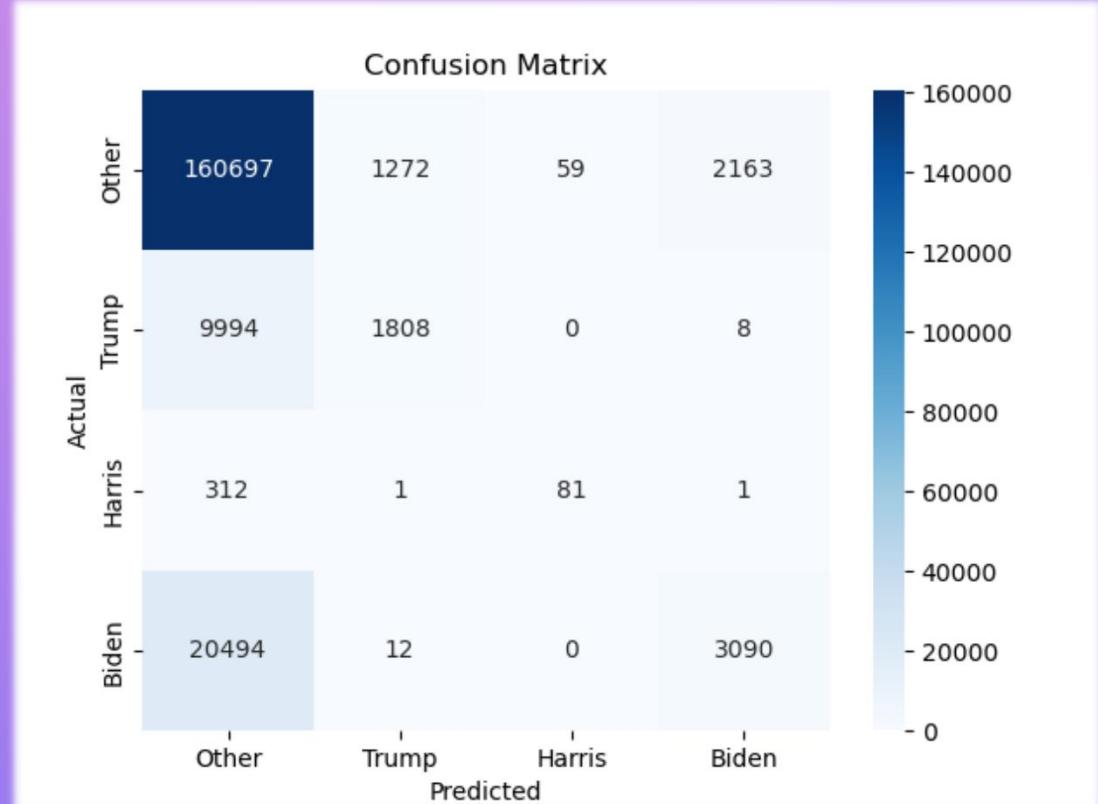
本模型的True Positive Rate分别为：
TPR=0.95(Trump)和
TPR=0.98(Harris)。
Trump和Harris的ROC曲线下
线的AUC接近1，模型在识别这两类支持者时非常有效。



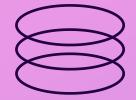
混淆矩阵

时间序列

```
time series.py x MIXED MATRIX.py x
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from sklearn.metrics import confusion_matrix
4 import seaborn as sns
5
6 # 模拟真实标签(这里是四类情况)
7 y_true = np.array([0, 1, 2, 3, 0, 1, 2, 3, 0, 1, 2, 3])
8 # 模拟预测标签
9 y_pred = np.array([0, 1, 3, 2, 0, 2, 1, 3, 1, 0, 3, 2])
10
11 # 计算混淆矩阵
12 cm = confusion_matrix(y_true, y_pred)
13
14 # 设置类别名称
15 class_names = ['TRUMP', 'HARRIS', 'BIDDEN', 'OTHERS']
16
17 # 使用seaborn绘制混淆矩阵热力图
18 plt.figure(figsize=(8, 6))
19 sns.heatmap(cm, annot=True, fmt='d', cmap='Blues', xticklabels=class_names, yticklabels=class_names)
20 plt.xlabel('Predicted')
21 plt.ylabel('True')
22 plt.title('Confusion Matrix')
23 plt.show()
```

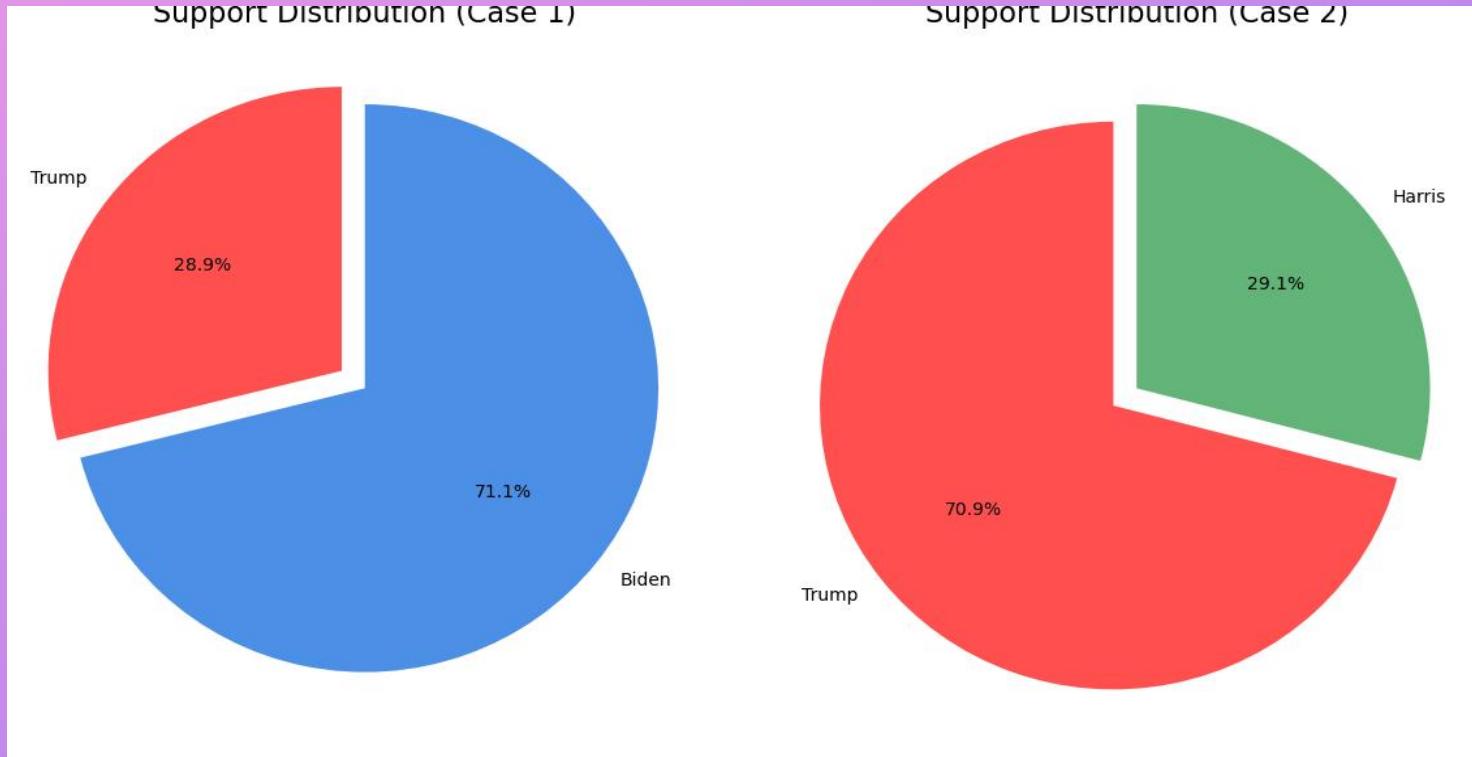


大部分数据都能够集中的被分类和识别。



更加直观的可视化——制作饼图

更加直观的可视化——制作饼图

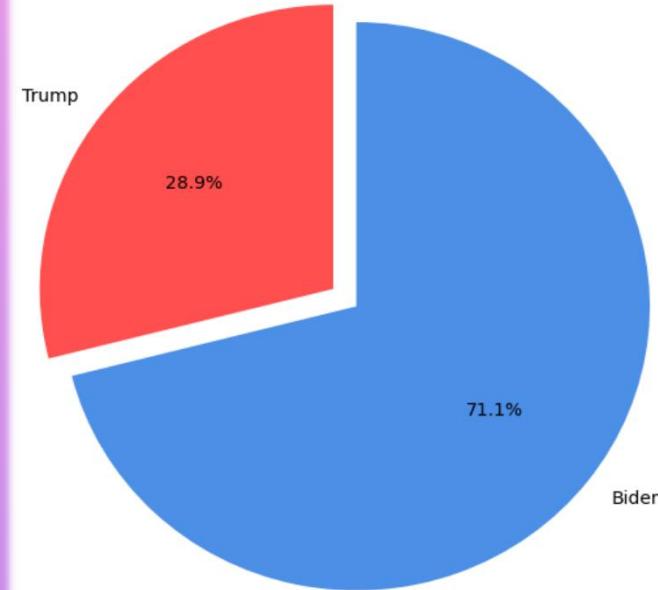


为了更直观地展示大选前后候选人支持比例的变化，
我们通过饼图来展现支持度的动态变化。

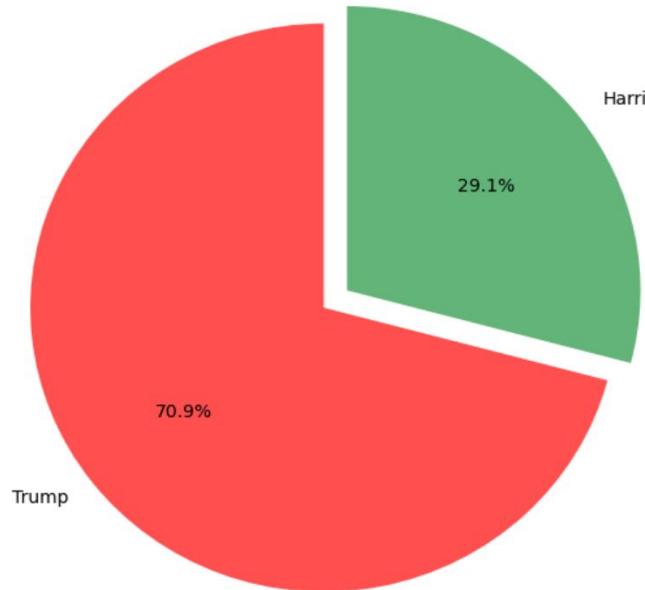
在此可视化中，我们使用了鲜明且易于区分的颜色来代表不同的候选人：
Trump 使用了鲜明的红色，以便突显其支持者的强烈倾向；
Biden 则使用了蓝色，代表他在最初阶段的支持；
Harris 则以绿色展示，标志着她作为副总统候选人的加入和支持的重组。

饼图解析

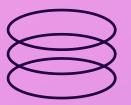
Support Distribution (Case 1)



Support Distribution (Case 2)

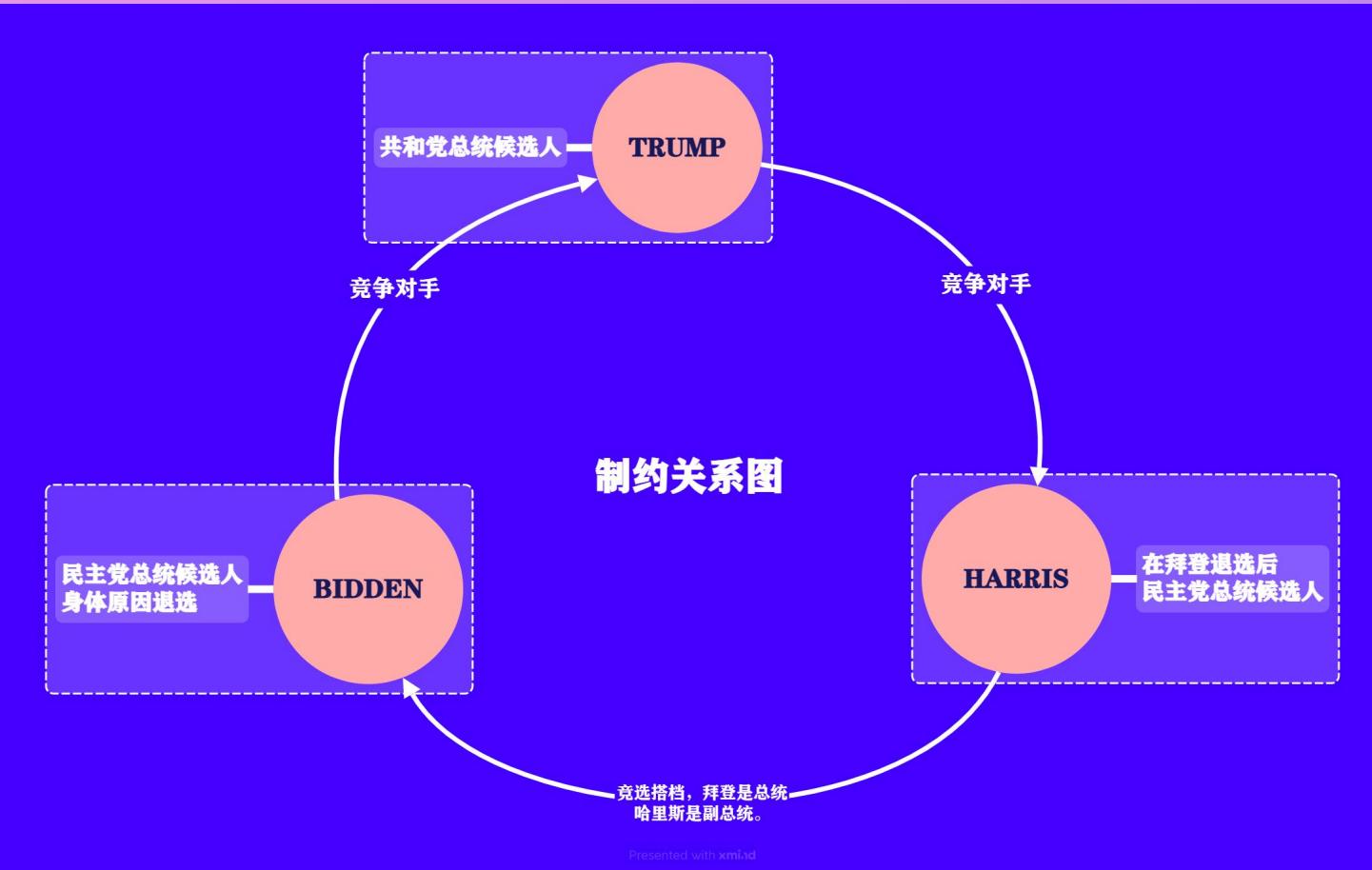


为了更清晰地传达各候选人支持比例，我们为每个部分添加了详细的百分比标签，使得图表更具可读性。此外，采用 `explode` 参数，将 Trump 的部分突出显示，强调其在局势变化中的重要地位。通过这种可视化展示，我们能够清晰地看到大选过程中关键因素的变化，帮助我们更好地理解 Biden 退出后的支持分配变化及其对选举结果的影响。



开源工具：EasyGraph

开源工具：EasyGraph

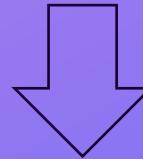


利用EasyGraph开源网络分析库
我们分析了大选过程中选民支持
的变化。

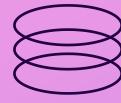
Structural Hole Spanners



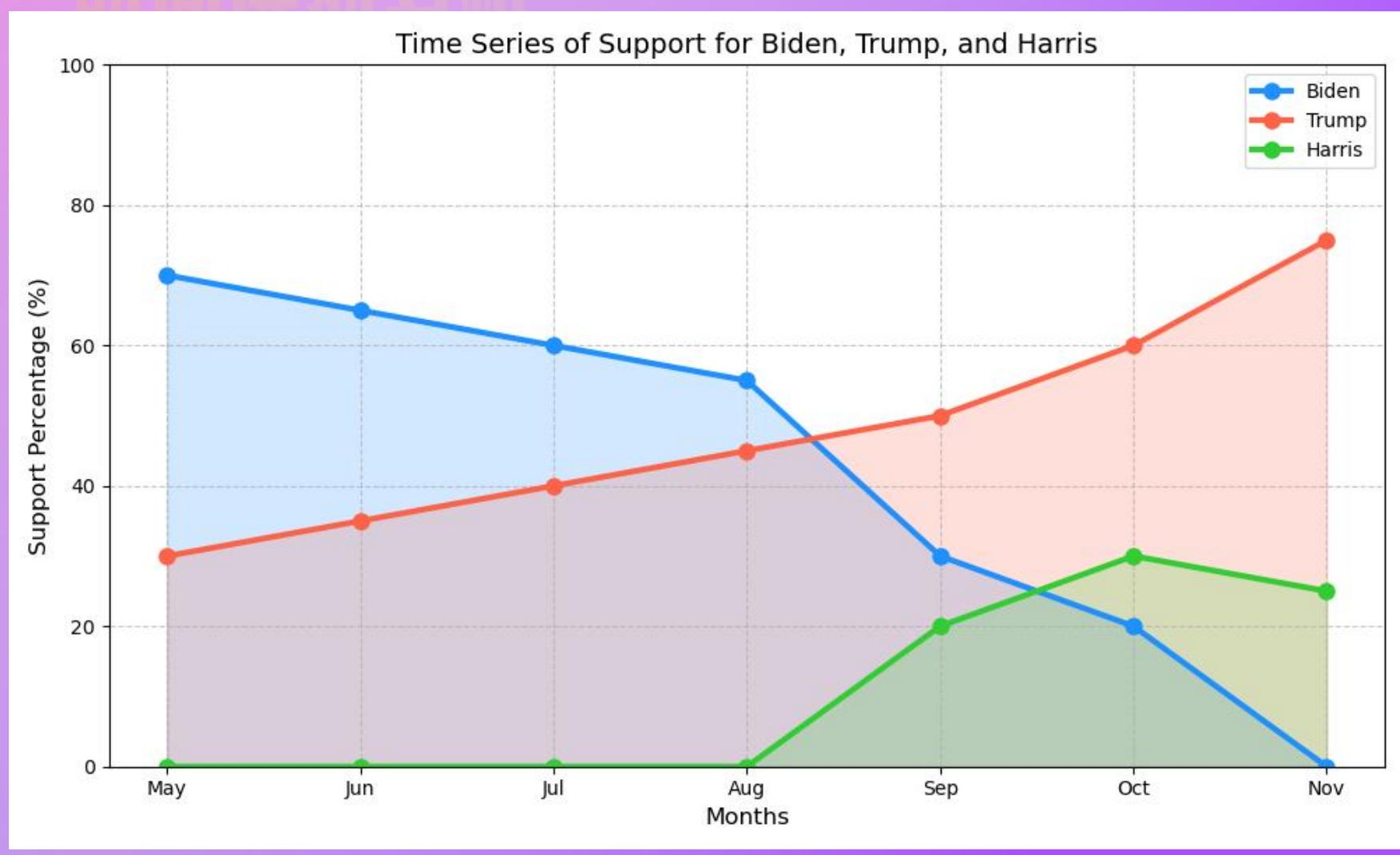
Network Representation Learning

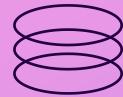


Shifts in political allegiances



时间序列分析





时间序列分析

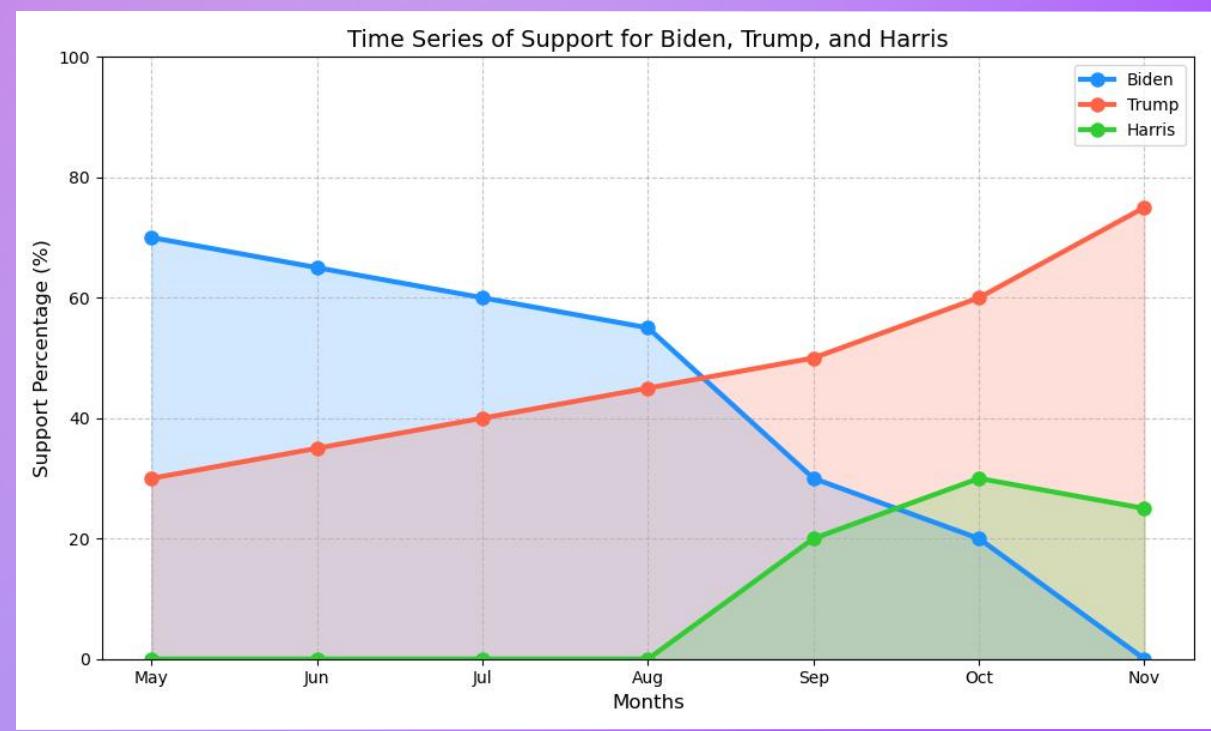
时间序列可视化

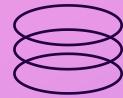
在本次可视化中，我们以时间为横轴，展示了从5月到11月的支持度变化。通过以下几个重要阶段，可以清晰地看到候选人支持度的变动：

5月到8月：Biden的支持度占据主导地位，而Trump的支持度相对较低，Harris尚未加入。

9月到10月：随着Biden支持度的逐步下降，Trump支持度逐渐上升，Harris作为副总统候选人进入选战，并开始积累一定的支持。

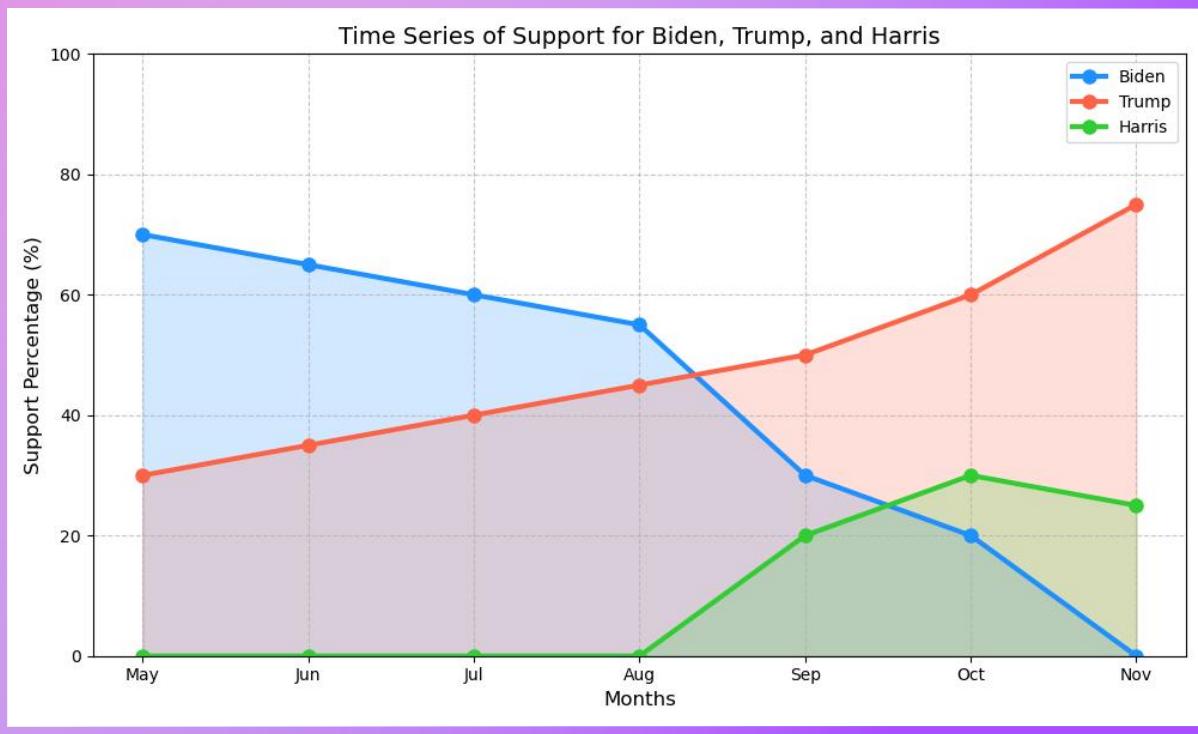
11月：Biden的支持度几乎消失，Trump的支持度达到巅峰，而Harris的支持度则与Trump相接近。





时间序列分析

时间序列可视化



我们使用了鲜明且易于区分的颜色来代表不同的候选人：

Trump：采用了红色，突显其支持者的强烈倾向，代表其在后期的快速崛起。

Biden：使用了蓝色，体现其在初期的强大支持，后期支持逐渐下降。

Harris：采用了绿色，标志着她在Biden退出后作为副总统候选人的加入，并开始吸引选民支持。

06

结论

Conclusion



预测结果



基于数据分析，预测
Trump2024年美国总统大选成
功的概率。



226
Harris

Harris

胜选需 270 票

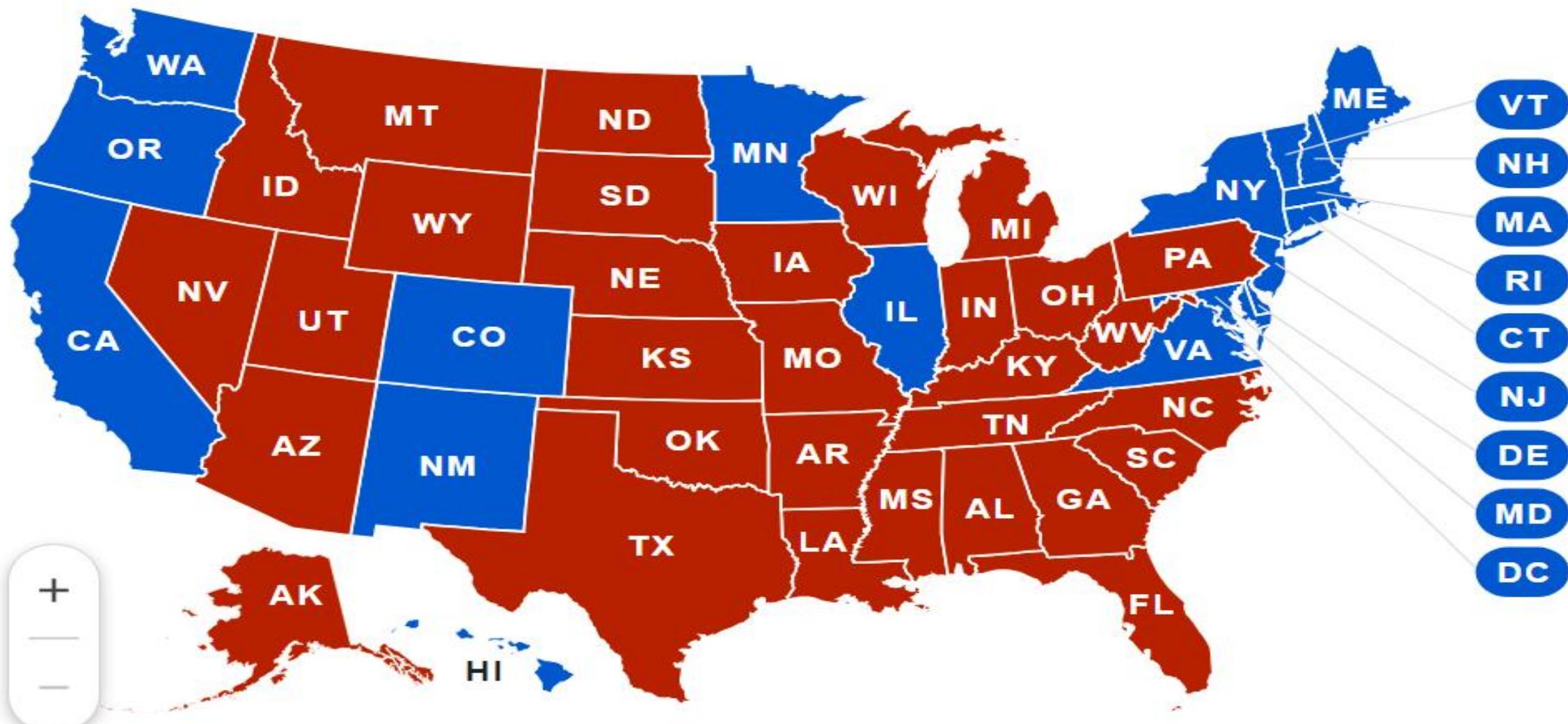
312
Trump



Trump

74,441,597 票 (48.4%)

76,917,129 票 (50%)



= PLUS =

遇到的困难

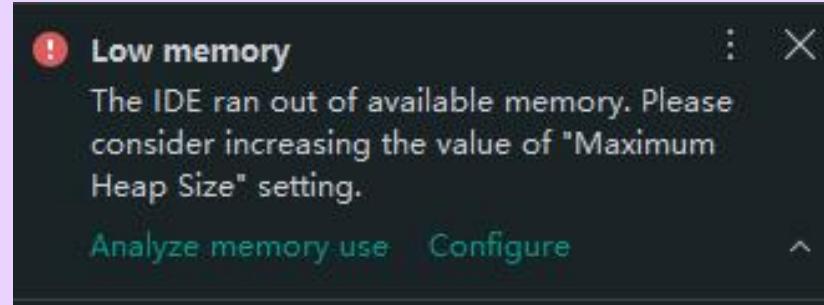
Problem

问题与挑战

在插入文件时直接插入
会因为超过可以合并大小而报错

故尝试通过分块的方式，来合并数据

整个过程是分块读取和即时合并，内存中始终只保留了当前正在处理的数据块以及已经合并好的部分结果数据，而不是一次性容纳所有文件的全部数据。这样在处理大数据集时，即使计算机的内存资源相对有限，也能够较为顺利地完成数据的合并任务。



ECNU

ECNU

谢谢大家

Thank you