



天津大学
TianJin University

疫情全回顾—— 疫情期间舆情演化分析与热点事件模式挖掘

李泽宇 (队长)

陈诺

李思思

王腾

贾世超

张加万 (指导老师)



ChinaVis 2020 Data Challenge

利用大数据分析技术和可视分析方法协助发现病毒传染源、监测疫情发展、调配救援物资，为疫情防控中的分析、指挥和决策提供有效依据和指南。

- 疫情时空态势分析
- 疫情传播模式分析
- 疫情预测与舆情监测
- 新冠病毒病理研究
- 疫情潜在影响与次生灾害分析



ChinaVis 2020 Data Challenge

利用大数据分析技术和可视分析方法协助发现病毒传染源、监测疫情发展、调配救援物资，为疫情防控中的分析、指挥和决策提供有效依据和指南。

- **疫情预测与舆情监测：**利用可视分析技术，预测疫情发展趋势与关键节点、分析社交媒体话题与情感的动态演变、对社会舆情进行态势感知。

疫情全回顾—— 疫情期间舆情演化分析与热点事件模式挖掘



天津大学
Tianjin University





数据、分析任务、分析流程



天津大学
Tianjin University





数据 2019-12-23 ~ 2020-05-10

| | 热文数据 | 热点事件数据 | 疫情新闻列表 |
|------|-------------------------------|------------------------------------------|--------------|
| 数量 | $300 * 140 \text{ 天} = 42000$ | 1781 | 690 |
| 来源 | 榜单：清博大数据 文章：今日头条 | 清博大数据 | 知微数据 |
| 核心字段 | 文章标题、内容 发布日期 阅读数 | 事件标题、概述、热词 发生日期 14天情感趋势 14天热度趋势 | 新闻标题 发生日期 |
| 类型 | 时序文本 | 时序文本，时间序列 | 时序文本 |






分析任务

1. 明确疫情发展阶段，识别重大热点事件
2. 从多个角度（文本内容/实体）和多个尺度（话题级/词级）探究舆论关注点的转变
3. 找出在情感相关维度上特殊的热点事件
4. 挖掘热点事件的热度演化模式



分析流程



天津大学
Tianjin University





数据模型与计算




天津大学
Tianjin University





统一投影空间（主题地图）的构造






与事件情感相关的度量准则的计算

| 度量准则 | 计算公式 | 描述 |
|------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|------------------------------------------------------------------------------------------------|
| 总体情感 | $\frac{\sum_i (p_i + m_i / 2)}{\sum_i (p_i + m_i + n_i)}$ | 以“修正”后的正面情感比例来度量，值越大，总体情感越偏正面。“修正”是为了避免将“中性占比很大，正面占比很小”这种事件识别为负面事件 |
| 争议程度 | $\frac{(P + N) - abs(P - N)}{P + M + N}$ | 正负情感值占总情感值的比例越大，且两者差异越小的事件认为争议程度越大 |
| 波动程度 | $[\sigma(\tilde{p}) + \sigma(\tilde{m}) + \sigma(\tilde{n})]/3$ | 以三种情感比例序列的标准差的平均值来度量。值越大，情感波动程度越大 |
| 反转情况 | $s = sign(indexOf(max(d)) - indexOf(min(d)))$ $\left\{ \begin{array}{ll} s * abs(max(d) - min(d)), & \text{if } max(d) * min(d) < 0; \\ undefined, & \text{else} \end{array} \right.$ | 对某事件，若情感比例差异序列d即有正值也有负值，则认为该事件存在情感反转，此时最大值和最小值之间的差异越大，意味着反转强度越高。最大值后于最小值出现，认为是负面转正面，反之认为是正面转负面 |





事件热度序列投影空间的构造



天津大学
Tianjin University





案例分析



天津大学
Tianjin University





案例1：探究舆论热点演化



天津大学
Tianjin University





疫情发展的四个阶段









疫情发展的四个阶段



天津大学
Tianjin University




疫情发展的四个阶段





理解“疫情通报”话题在第二和第四阶段的差异



天津大学
Tianjin University





“病毒研究”话题的演化




天津大学
Tianjin University





“韩国/欧洲/南美疫情”话题的演化




天津大学
Tianjin University





“美国/日本疫情”话题的演化




天津大学
Tianjin University





“武汉/医院/医护”话题的演化




天津大学
Tianjin University





“防疫物资”话题的演化



天津大学
Tianjin University






“防疫物资”话题的演化

01-20 - 02-15

02-16 - 03-12

03-13 - 04-08

04-09 - 05-10



武汉

口罩

悉尼

澳洲

车站

变态

全淹

泳池

风暴

新州

肺炎

警报

医用

湖北

药店

描述词

韩国

青田

伊朗

华联

巴基斯坦

青田县

意大利

蝗灾

巴林

空运

瑞士

齐齐哈尔

橘子

文文

动物

美国

留学生

呼吸机

英国

法国

回国

王女士

入境

北京

乘客

国际航班

天津

特朗普

目的地

祖国

渔船

货机

越南

出口

尼泊尔

海域

扬中

古巴

索马里

启迪

顺水

阿根廷

巴厘岛

格力

上海

01-20 - 02-15

02-16 - 03-12

03-13 - 04-08

04-09 - 05-10



卫健委

红十字会

国务院

外交部

京东

央视

波音

世界卫生组织

环球网

淘宝

武汉大学

中国政府

世卫组织

阿里

人民政府

机构

外交部

中国政府

国务院

卫生部

世卫组织

环球网

韩国政府

商务部

中国外交部

淘宝

世界卫生组织

民航局

央视

中新社

新华社

外交部

国务院

民航局

中国政府

卫生部

海关总署

环球网

新华社

世卫组织

卫健委

白宫

霍普金斯大学

央视

商务部

美国政府

外交部

商务部

海关总署

卫生部

环球网

中国政府

民航局

新华社

路透社

波音

卫健委

国务院

世界卫生组织

上海浦东国际机场

中国外交部



天津大学

Tianjin University





案例2：探究事件在话题和情感上的分布 以及事件的演化模式





天津大学
Tianjin University





全部事件和疫情相关事件在时间维度上的分布











天津大学
Tianjin University





事件在三个情感属性下的分布




天津大学
Tianjin University





事件在各个主题的分布情况





天津大学
Tianjin University





理解事件投影空间



天津大学
Tianjin University





讨论与总结



天津大学
Tianjin University





伸缩性

计算的伸缩性：只有描述词提取需要在线计算，对于 G^2 statistics 方法来说，采样是一个既能大幅减少计算量又不显著降低描述词质量的方法。

可视化的伸缩性：Canvas 或者 WebGL 绘制技术。当然，采样也是一种缓解办法。

流数据

基于模型的transform，新近的文章和事件可以实时、无缝地融入已有图中。两点需要注意：

- 新发生的事件和新闻由于尚未积累起足够的影响力从而可能会被忽略
- 新数据可能不会对已训练的模型产生实质影响

本质上并没有解决流数据的问题，仍然需要定期手动更新。

总结

我们的系统可以被广泛应用在舆情状态感知、关键节点识别、话题演变分析、事件演化范式探索上。





天津大学
TianJin University

谢谢!

李泽宇 (队长)

陈诺

李思思

王腾

贾世超

张加万 (指导老师)

lzytianda@tju.edu.cn