



同济大学交通运输工程学院

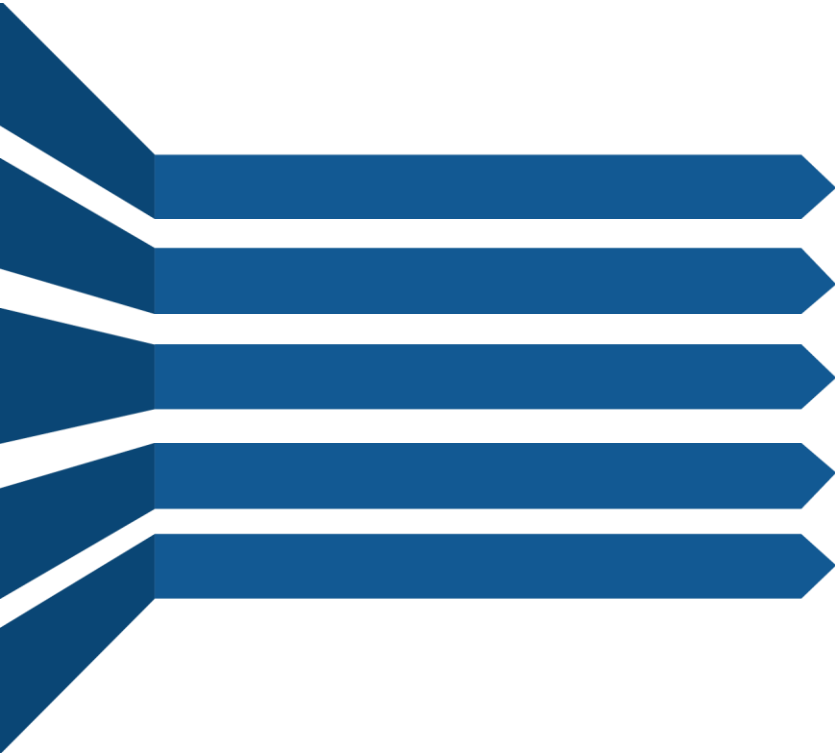
COLLEGE OF TRANSPORTATION ENGINEERING
TONGJI UNIVERSITY

基于文本挖掘的高铁列车 晚点影响因素分析

指导老师：邢莹莹

答辩人：赵冠华



- 
- A decorative graphic on the left side of the slide consisting of five horizontal blue arrows pointing to the right, stacked vertically. The arrows are of varying lengths and are set against a background of blue diagonal stripes.
- 第一节 研究背景与意义**
 - 第二节 研究思路与方法**
 - 第三节 研究成果与分析**
 - 第四节 研究结论与展望**
 - 第五节 参考文献与致谢**



第一节 研究背景与意义

1.1 研究背景

1.2 研究问题

1.3 研究意义



1.1 研究背景

高铁VS航空

工具	日运输量 (人次)	日运输量占比	票价 (元)	速度
铁路	90628	99.18%	98元 (火车)；高铁二等、一等、商务舱票价分别为314元、504元、995元	8~10小时 (火车)；2小时40分钟 (高铁)
航空	745	0.82%	票价不定：远期：240元，270元，290元；近期：880元，950元，1280元	1小时20分钟

表1.1 广州-长沙路段高铁、航空日运营情况对比表（20190812）

时速差距小
准点率高



选择高铁出行

高铁列车晚点





第二节 研究思路和方法

2.1 基于先验LDA的高铁列车晚点影响因素特征提取

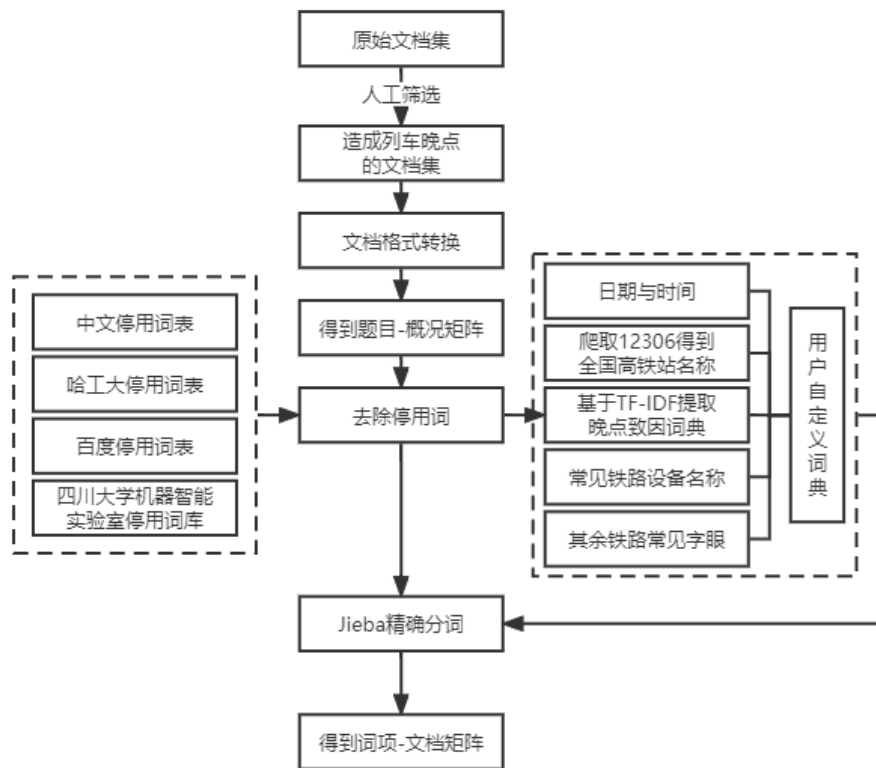
2.2 基于模糊事故树的高铁晚点影响因素分析

2.3 技术路线



2.1 基于先验LDA的高铁列车晚点影响因素特征提取

2.1.1 数据集预处理



原始文档集 → 词项-文档矩阵

- 目的：方便下一步的特征提取
- 分词：JIEBA分词精确模式

基于概率语言模型，生成句中所有可能词情况的有向无环图，动态规划查找最大概率路径。

- ### 词频-逆文档频率方法:
- ①用户自定义词典; ②晚点致因词典 (10类)





2.1.2 基于相关性的先验知识提取

“词项特征 w 与晚点致因模式 c ”



“词项特征 w 与潜在主题特征 z 的相关性”

$$\chi^2(w, c_i) = N * \frac{[P(w, c_i)P(w', c'_i) - P(w', c_i)P(w, c'_i)]^2}{P(w)P(w')P(c_i)P(c'_i)} \quad (1)$$

$$\bar{T}(w_i, c_j) = \frac{T(i, j)^2}{\sum_{i=1}^n T(i, j) * \sum_{j=1}^m T(i, j)} \quad (2)$$

卡方值量化相关性

- ◆ 假设：
 w 与 c 为一维自由度卡方分布
- ◆ 计算：式(1)
- ◆ 结果：
获得 w 与 c 的相关性矩阵 M

相关性矩阵初始化

- ◆ 目的：
方便对比
- ◆ 计算：式(2)
- ◆ 结果：
获得 w 与 c 初始化后的相关性矩阵 T

聚类判断相关性强度

- ◆ 方法：K-Means，并降序排列
- ◆ 分类：
强关联 —— 0,1,2
复杂关联 —— 3~8
弱关联 —— 9,10,11

相关性数值转换

- ◆ 方法：预分配10*10个潜在主题特征 z
- ◆ Γ 值：
强关联 —— 2
复杂 —— $(\theta_i - t)/t'$
弱关联 —— 10^{-12}



2.1.3 先验知识与LDA模型的整合

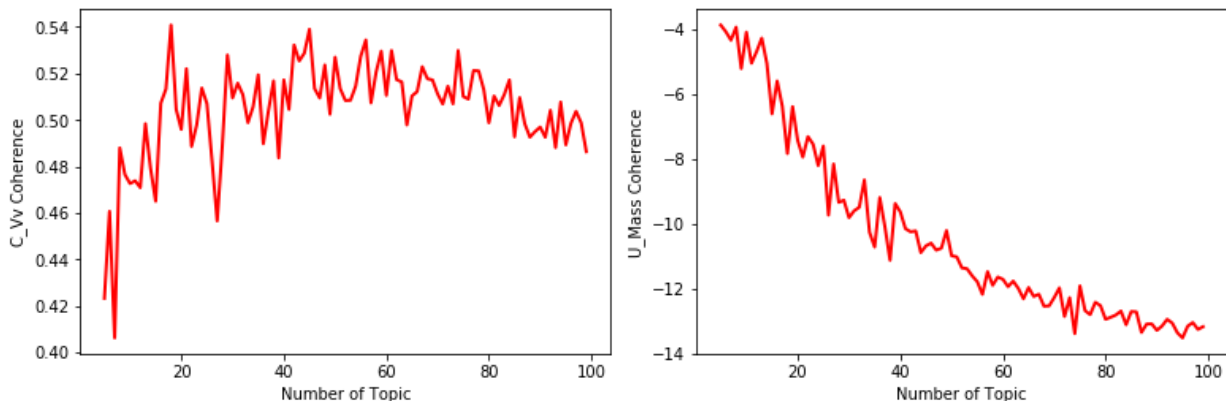
最优主题个数

方法：主题连贯性

U_{mass} 最低点

C_v 最高点

最优个数为74



主题更新公式

Γ 值	P(w,z)
2	加大
10^{-12}	减小
其它值	修改

φ 代表主题的词语分布, Θ 代表主题的多项式分布, α 、 β 为超参数

$$P(z_i = j | z_{-i}, w, \alpha, \beta) = \left(\frac{n_{-i,j}^{w_i} + \beta}{\sum_w n_{-i,j}^{w_i} + W\beta} \right) \left(\frac{n_{-i,j}^{d_i} + \alpha}{\sum_j n_{-i,j}^{d_i} + T\alpha} \right) * \Gamma(w_i' z_j)$$

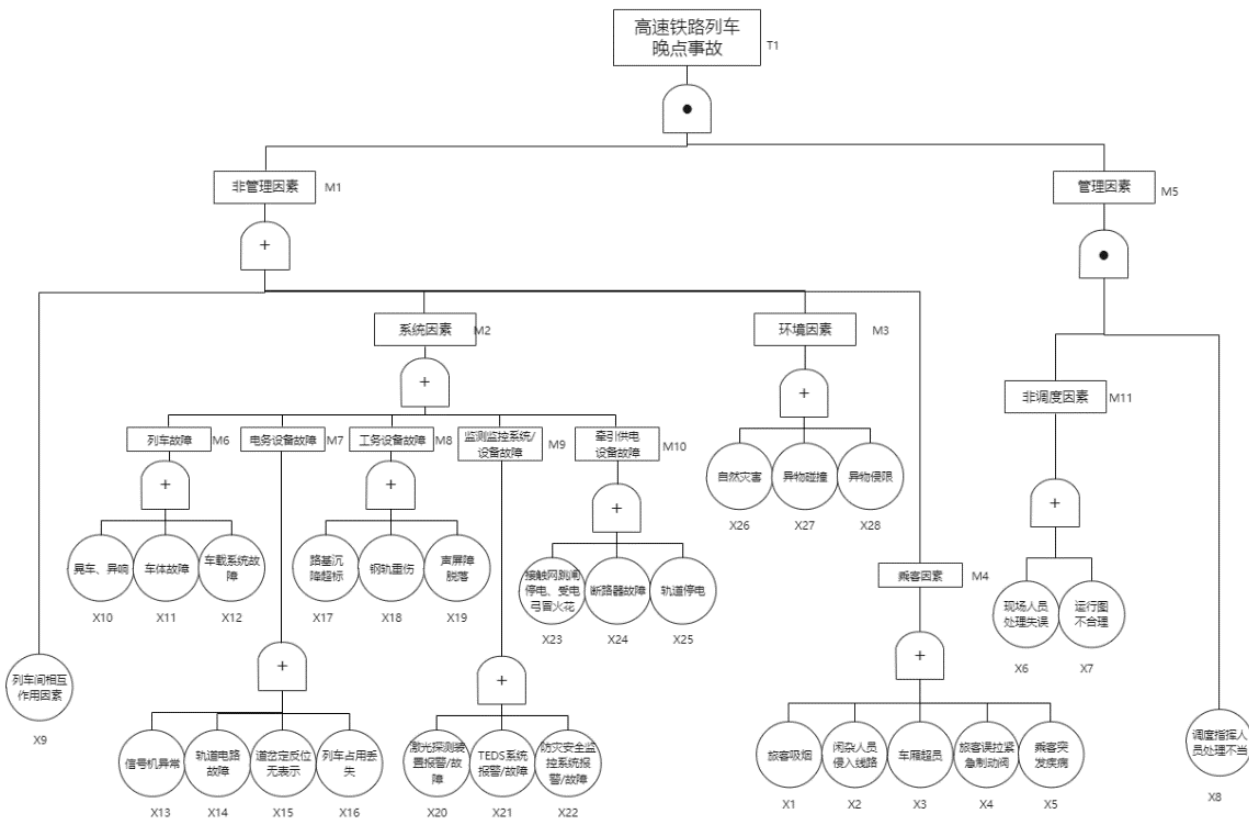
$$P(z_i = j | z_{-i}, w, \alpha, \beta) = \varphi' \left(\frac{n_{-i,j}^{d_i} + \alpha}{\sum_j n_{-i,j}^{d_i} + T\alpha} \right) \varphi' = \frac{(1 + \Gamma(w_i' z_j)) n_{-i,j}^{d_i} + \alpha}{\sum_j (1 + \Gamma(w_i' z_j)) n_{-i,j}^{d_i} + T\alpha}$$

» 得到高铁列车晚点致因主题-词分布, 利于理解相互关系, 利于事故树建立与进一步分析。



2.2 基于模糊事故树的高铁晚点影响因素分析

2.2.1 模糊事故树构建：顶上事件——高铁晚点



确定基本事件

- 共28个基本事件
- 共11个中间事件

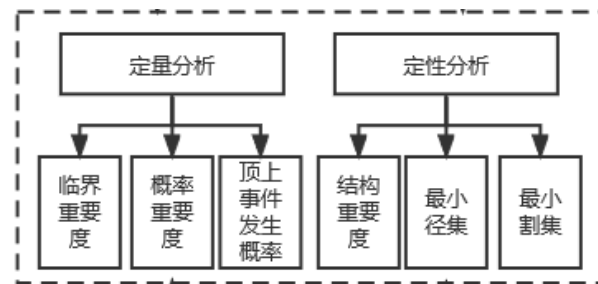
三角模糊处理

- 可统计基本事件：
 $m = P$, $\alpha = 0.95m$, $\beta = 1.05m$
- 不可统计基本事件：
采用专家评价 3σ 表征法
 $m = AVE$, $\alpha = m - 3\sigma$, $\beta = m + 3\sigma$



2.2.2 定性与定量分析

- 弄清系统内顶事件发生的可能性
- 了解高铁晚点的形成途径与控制途径
- 认识各基本事件在结构上对顶上事件的影响程度



定性分析

最小割集:

布尔代数法

- ◆ 导致顶事件发生的基本事件的集合
- ◆ 表明系统危险性
- ◆ 方便掌握事故发生规律

最小径集:

对偶树法

- ◆ 导致顶事件不发生的基本事件的集合
- ◆ 表明系统安全性
- ◆ 确定最经济有效的控制事故发生的方案

结构重要度I:

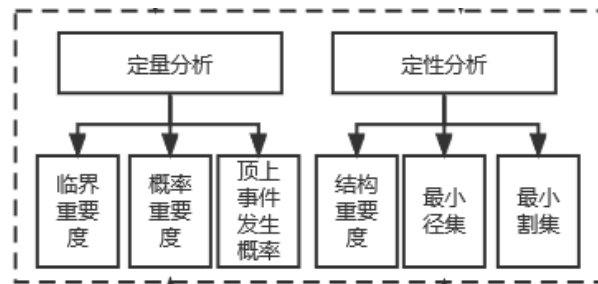
$$I(i) = \frac{1}{k} \sum_{j=1}^n \frac{1}{n_j} (j \in k_j)$$

- ◆ 各基本事件在结构上对顶上事件的影响程度
- ◆ 仅仅通过事故树结构, 不考虑发生概率



2.2.2 定性与定量分析

- 估算系统可靠性，验证模型适用性
- 得到降低晚点率的有效措施
- 得到治理高铁晚点的优先级顺序



定量分析

顶事件发生概率：最小径集逼近法

概率重要度Q： $Q(i) = \frac{\partial P(T)}{\partial q_i}$

临界重要度C： $C_i = \frac{P(q_i)}{P(T)} Q(i)$

- ◆ 估算系统的可靠性
- ◆ 与实际数据对比，验证模型适用性
- ◆ 基本事件概率变化引起顶事件概率变化的程度
- ◆ 通过中值法求解模糊概率重要度
- ◆ 得到减少事故发生的有效措施
- ◆ 基本事件发生概率的相对变化率与顶上事件发生概率的相对变化率之比
- ◆ 得到治理的优先级顺序



2.3 技术路线

① 问题提出：基于文本挖掘的高铁列车晚点影响因素分析

② 模型数据获取与处理

基于先验LDA的高速铁路列车晚点影响因素特征提取

高速铁路故障写实表与处置过程

数据预处理

筛选晚点记录

中文分词

引入停用词

提取相关性先验知识

计算卡方值

确立晚点因素词项特征与晚点类型的相关关系

先验知识与LDA模型整合

强相关性关系

弱相关性关系

复杂关联关系

③ 模型建立与分析

基于模糊事故树的影响因素分析

构建模糊事故树

确定基本事件

构筑树结构

节点模糊处理

定性分析

最小割集

最小径集

结构重要度

定量分析

顶上事件
模糊概率

模糊概率重要度

临界重要度

高铁列车晚点风险应对

降低晚点率
的有效措施

风险因素
治理优先级

④ 结论与展望



第三节 研究成果与分析

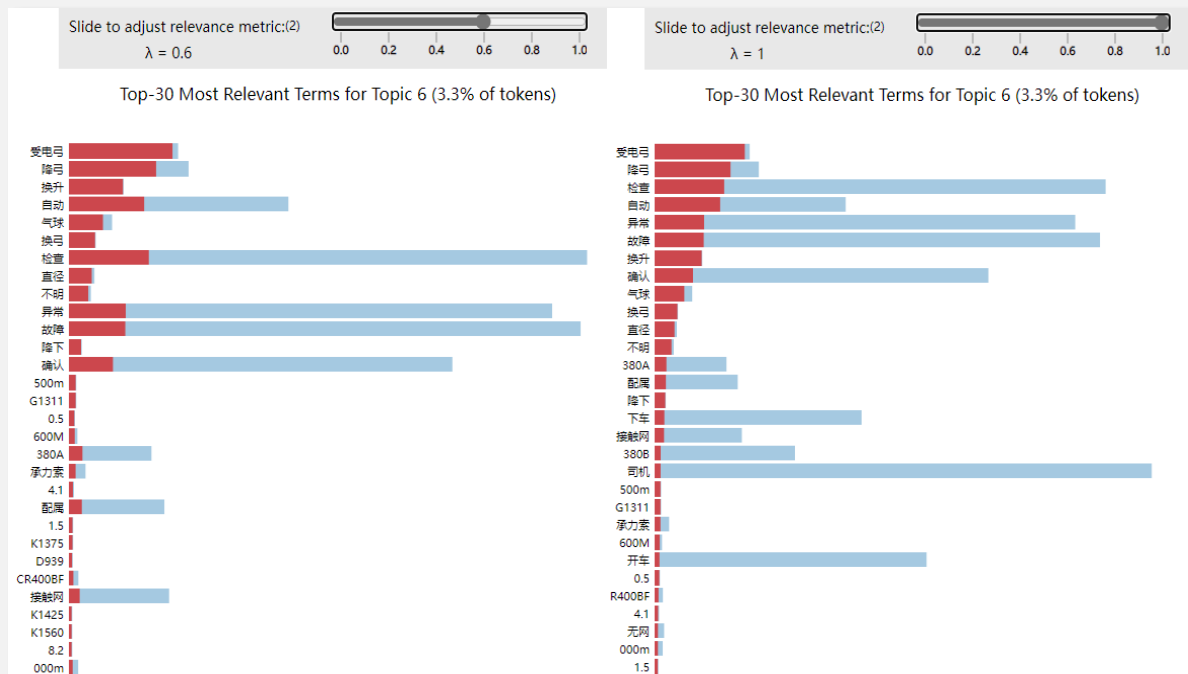
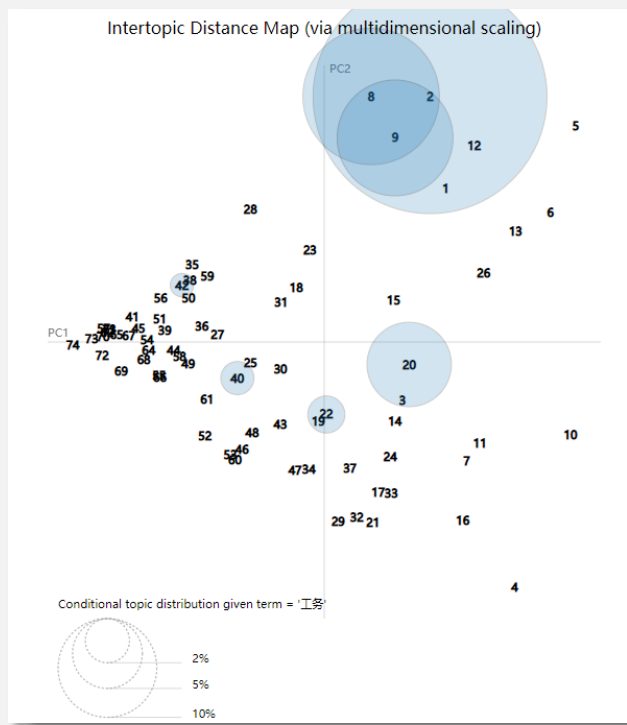
3.1 LDA建模可视化

3.2 模糊事故树分析 (1) (2) (3)



3.1 LDA建模可视化

利用 pyLDAvis 包可视化主题模型



气泡大小



气泡中心距离



最相关词汇



参数 λ



红条



蓝条



3.2 模糊事故树分析 (1)

最小割集与最小径集

$$\begin{aligned} T &= M_1 * M_5 \\ &= X_6 X_8 X_1 + X_6 X_8 X_2 + \cdots + X_6 X_8 X_{28} \\ &\quad + X_7 X_8 X_1 + X_7 X_8 X_2 + \cdots + X_7 X_8 X_{28} \end{aligned} \quad \left. \vphantom{\begin{aligned} T &= M_1 * M_5 \\ &= X_6 X_8 X_1 + X_6 X_8 X_2 + \cdots + X_6 X_8 X_{28} \\ &\quad + X_7 X_8 X_1 + X_7 X_8 X_2 + \cdots + X_7 X_8 X_{28} \end{aligned}} \right\} \text{共50个}$$

- 50个最小割集——晚点事件极易发生且发生路径多
- 基本事件数量少——引发路径较简单，风险性较大

$$\begin{aligned} T' &= M'_1 + M'_5 \\ &= X'_1 * X'_2 * X'_3 * X'_4 * X'_5 * X'_9 * X'_{10} * X'_{11} * X'_{12} * X'_{13} \\ &\quad * X'_{14} * X'_{15} * X'_{16} * X'_{17} * X'_{18} * X'_{19} * X'_{20} * X'_{21} * X'_{22} \\ &\quad * X'_{23} * X'_{24} * X'_{25} * X'_{26} * X'_{27} * X'_{28} + X'_6 X'_7 + X'_8 \end{aligned}$$

- 3个最小径集——存在3种方案使晚点事件不发生

顶上事件模糊概率

□ (0.02438 , 0.03289 , 0.04275)

即：高铁准点率 **96.7%**，
波动范围 **95.7% ~ 97.6%**

- ✓ 根据人民网报道，2015年全国
始发正点率98.8%
到达正点率95.4%
- ✓ 根据《中国的高速铁路发展报告》

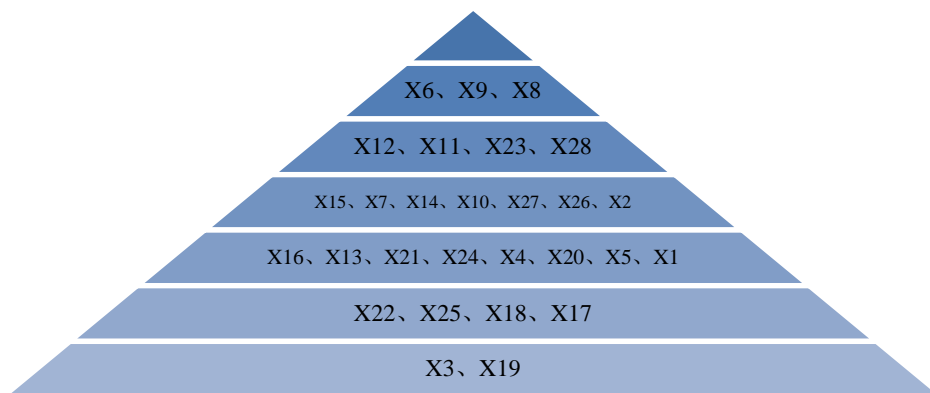
截至 2019.7.8	始发 准点率	到达 准点率
高铁	>98%	>95%
复兴号	99%	98%

» 验证了模糊事故树模型
在高铁晚点研究上的适用性



3.2 模糊事故树分析 (2)

- **结构重要度 I:** $I[X_8] > I[X_6] = I[X_7] > I[X_1] = I[X_2] = I[X_3] = I[X_4] = I[X_5]$
 $= I[X_9] = I[X_{10}] = I[X_{11}] = I[X_{12}] = I[X_{13}] = I[X_{14}] = I[X_{15}] = I[X_{16}]$
 $= I[X_{17}] = I[X_{18}] = I[X_{19}] = I[X_{20}] = I[X_{21}] = I[X_{22}] = I[X_{23}] = I[X_{24}]$
 $= I[X_{25}] = I[X_{26}] = I[X_{27}] = I[X_{28}]$
- ◆ X8最重要；X6、X7次之；除X8、X6、X7以外的剩余25个基本事件处于同等地位，最不重要
- **临界重要度 C:**



- ◆ C越大，基本事件越重要，其所在阶梯越在上方。
- ◆ 应更加关注第一、二阶梯内基本事件的治理

X_6, X_9, X_8

$X_{12}, X_{11}, X_{23}, X_{28}$



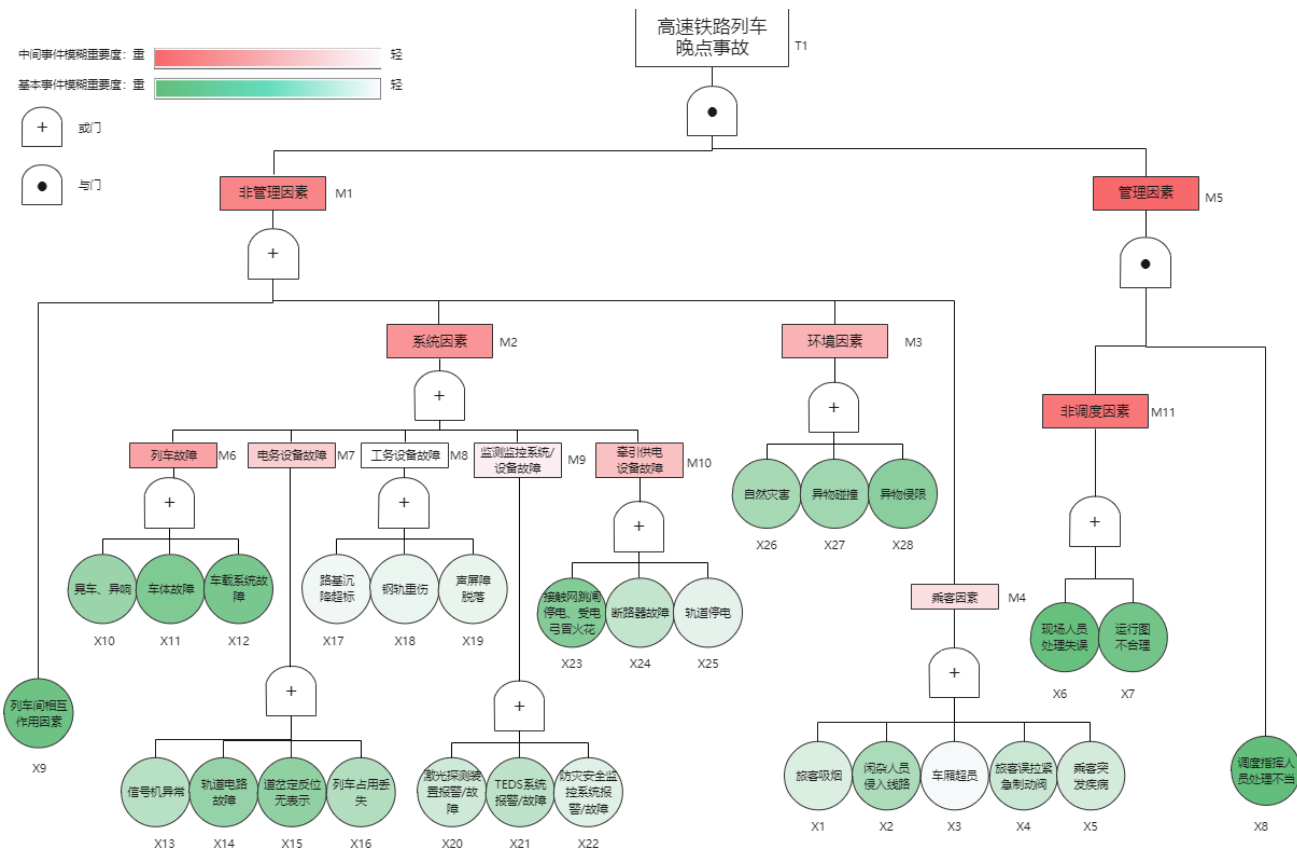
3.2 模糊事故树分析 (3)

➤ 模糊概率重要度 Q:

- ◆ 基本事件
绿色阶
中间事件
红色阶
- ◆ Q越大,
越重要
- ◆ 降低晚点发生
概率的有效措施

M_{11} 、 M_6 、 M_4

X_9 、 X_{23} 、 X_8





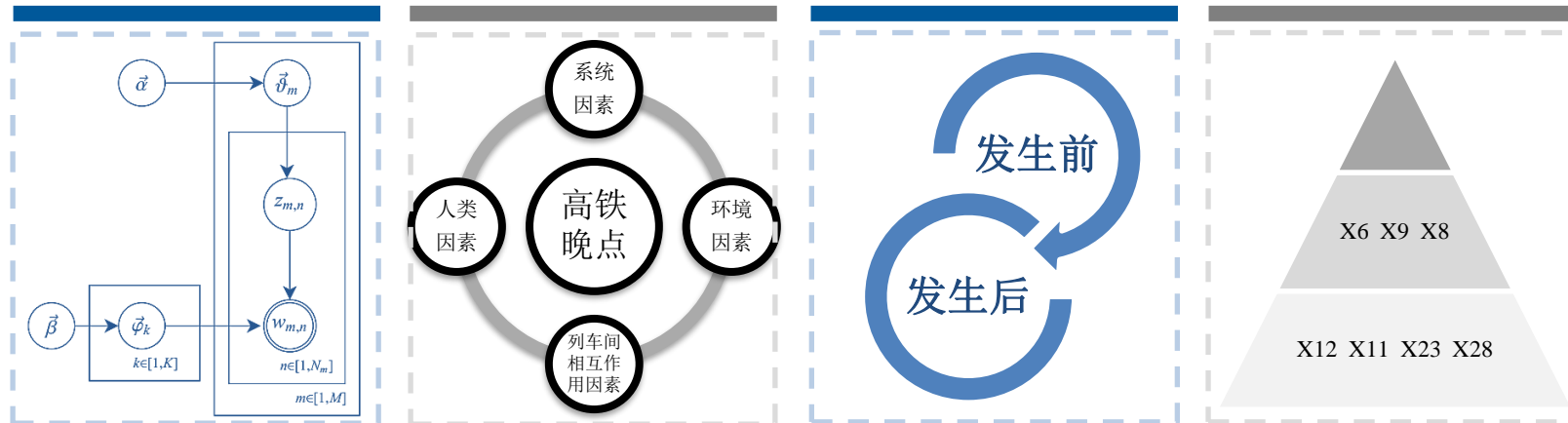
第四节 研究结论与展望

4.1 论文工作与结论

4.2 不足与展望



4.1 论文工作与结论



先验LDA模型

- 数据预处理
- 计算相关性矩阵
- 主题更新公式
- 得到主题-词分布

晚点影响因素

- 4大类别
- 28个基本事件
- 6个指标
- 3个重要度

针对性措施

- 2个时间节点
- 降低概率、及时发觉
- 避免传播、尽力恢复
- 分为5个角度

治理优先级

- 关注：
 - ① 管理因素
 - ② 列车故障
 - ③ 牵引供电故障
 - ④ 异物侵限



4.2 不足与展望

不足之处

- ❑ 列车间相互作用因素包含多种类型，本文数据集无法满足细分类型的需要；
- ❑ 三角模糊时，不可统计基本事件的模糊概率仅来自三位专家，具有主观性。

编号	专家1	专家2	专家3	m	σ
X_6	0.2764	0.2455	0.2639	0.2619	1.269E-02
X_7	0.0134	0.0119	0.0142	0.0132	9.534E-04
X_8	0.1433	0.1342	0.1313	0.1363	5.112E-03
X_{20}	0.0132	0.0127	0.0119	0.0126	5.354E-04
X_{22}	0.0079	0.0088	0.0086	0.0084	3.859E-04

未来展望

- ❑ 更新数据集，完善论文，以得到更精确的结论；
- ❑ 增加专家位数来降低主观影响，或直接采用更客观且适合的评分方式。





参考文献

1. 陆娅楠. 我国高铁运营里程超4万公里[N]. 人民日报,2021-12-31 (001).
2. 吴漫云.航空—高铁竞合关系分析[J]. 民航管理,2019 (10) :11-15.
3. UIC450-2. Assessment of the performance of the network related to rail traffic operation for the purpose of quality analyses - delay coding and delay cause attribution process[S]. Paris, France : International Union of Railways, 2009.
4. John Preston et al. Impact of Delays on Passenger Train Services : Evidence from Great Britain[J]. Transportation Research Record, 2009, 2117(1) : 14-23.
5. Nadjla Ghaemi et al. Impact of railway disruption predictions and rescheduling on passenger delays[J]. Journal of Rail Transport Planning & Management, 2018, 8(2) : 103-122.
6. J Wang, Granlf M , J Yu. Effects of winter climate on high speed passenger trains in Botnia-Atlantica region[J]. Journal of Rail Transport Planning & Management, 2020.
7. 翟恭娟. 高速铁路列车运行调整优化研究[D]. 西南交通大学,2013.
8. 汪静, 彭一川, 陆键. 基于ISM的高铁列车晚点影响因素分析[J]. 中国铁路,2020, (01) :48-52.
9. 纪媛媛. 基于特征选择与机器学习的列车晚点预测方法研究[D]. 北京交通大学,2020.
10. 石睿. 基于数据驱动的高速铁路列车晚点分析及预测方法研究[D]. 北京交通大学,2021.
11. Allahyari M , Pouriyeh S , Assefi M , et al. A Brief Survey of Text Mining: Classification, Clustering and Extraction Techniques[J]. 2017.
12. Hofmann T. Probabilistic latent semantic indexing: Proceedings of the 22nd Annual International SIGIR Conference. New York: ACM Press, 1999:50-57.



参考文献

13. Blei D ,Ng A, Jordan M. Latent Dirichlet Allocation .[J]. Journal of Machine Learning Research, 2003, 3 : 993-1022.
14. Chemudugunta C , Holloway A , Smyth P , et al. Modeling Documents by Combining Semantic Concepts with Unsupervised Statistical Learning[J]. Springer-Verlag, 2008.
15. Allahyari M , Kochut K . Semantic Context-Aware Recommendation via Topic Models Leveraging Linked Open Data[J]. 2016.
16. Blei D M , Griffiths T L , Jordan M I , et al. Hierarchical Topic Models and the Nested Chinese Restaurant Process[J]. Advances in neural information processing systems, 2004, 16.
17. Blei D M , Mcauliffe J D . Supervised topic models. In NIPS, 2007, Vol 7.121-128.
18. Petinot Y , Mckeown K R , Thadani K . A hierarchical model of web summaries[C]// The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference, 19-24 June, 2011, Portland, Oregon, USA - Short Papers. DBLP, 2011.
19. X Mao, Z Ming, T Chua, et al.; SSHLDA: A Semi-Supervised Hierarchical Topic Model[C].EMNLP, 2012.
20. 王峰. 基于文本挖掘的高铁车载设备故障诊断方法研究[D]. 北京交通大学,2016.
21. Griffiths T L, Steyvers M. Finding Scientific Topics[J]. Proceedings of the National Academy of Sciences of the United States of America, 2004, 101(S1): 5228-5235.
22. Teh Y, Jordan M, Beal M, et al. Hierarchical Dirichlet Processes [J]. Journal of the American Statistical Association, 2007, 101(476): 1566-1581.
23. 曹娟, 张勇东, 李锦涛, 等. 一种基于密度的自适应最优LDA模型选择方法[J]. 计算机学报, 2008, 31(10): 1780-1787.



参考文献

24. Arun R , Suresh V , Madhavan C E V , et al. On Finding the Natural Number of Topics with Latent Dirichlet Allocation: Some Observations[J]. Springer-Verlag, 2010.
25. D. Newman, J. H. Lau, K. Grieser, and T. Baldwin. Automatic evaluation of topic coherence. In Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics, pages 100–108. 2010
26. D. Mimno, H. M. Wallach, E. Talley, M. Leenders, and A. McCallum. Optimizing semantic coherence in topic models. In Proc. of the Conf. on Empirical Methods in Natural Language Processing, pages 262–272. 2011.
27. Michael Roder, Andreas Both, and Alexander Hinneburg. 2015. Exploring the space of topic coherence measures. In the Eighth ACM International Conference on Web Search and Data Mining, pages 39– 408.
28. 税昌锡.语义特征分析的作用和语义特征的提取[J]. 北方论丛,2005 (03) :66-70.
29. Calinski, T, Harabasz, et al. A Dendrite Method for Cluster Analysis[J]. Comm in Stats Simulation & Comp, 1974.
30. 胡思继. 列车运行图编制理论与方法[M]. 2013.
31. 袁强. 高速铁路列车晚点分布及传播模型研究[D]. 北京交通大学,2020.
32. 人民网 . 正点率达 98.9% 累计 50 亿人次坐高铁出行 [EB/OL]. (2016-7-22) [2022-5-7].
http://m.cnr.cn/news/20160722/t20160722_522759873.html.
33. 罗强. 面向“一日一图”的列车运行图与动车组交路协同优化研究[D]. 北京交通大学,2021.
34. 李津, 文超, 杜雨琪, 徐传玲. 基于梯度提升回归树的武广高铁区间晚点恢复策略研究[J]. 中国铁路,2021, (10) :76-84.

请各位老师批评指正！

——基于文本挖掘的高铁晚点影响因素分析



答辩人：赵冠华 导师：邢莹莹

答辩时间：2022年6月8日



同济大学
TONGJI UNIVERSITY

Algorithm 1 : 相关性知识提取算法步骤

输入数据: 处置过程与处置写实表记录 D ; 高铁常见晚点致因词典 Ω ; 晚点致因模式集 C ; 潜在晚点致因主题特征集 Z 。

输出结果: 晚点描述词项特征 $w_i \in W$ 与晚点致因主题特征 $z_k \in Z$ 相关性集合 $\Gamma(w_i, z_k)$ 。

- 1: $W \leftarrow$ 词库: 由加入常见晚点致因词典 Ω 的中文分词工具对 D 进行中文分词得到;
- 2: $M \leftarrow$ 词项文档矩阵: 基于 W 和 D ,由向量空间模型(VSM)表达得到;
- 3: for $w_i \in W$ 且 $c_j \in C$ do;
- 4: $T(i, j) \leftarrow$ 词项特征 w_i 和晚点致因模式 c_j 的相关性, 由公式(3-1)求得。
- 5: end for
- 6: $\bar{T} \leftarrow$ 通过公式(3-2)对 T 进行归一化。
- 7: $\Xi \leftarrow k(k=12)$ 个聚类类簇: 由K-means对 \bar{T} 聚类得到, 并降序排列。
- 8: $\Theta \leftarrow$ 复杂相关程度集合: Ξ 中除最高三个和最低三个类簇外, $\Theta_i(i=1,2,3,4,5,6)$ 表示剩下 $k-6=6$ 个类簇中第 i 个的中心点。
- 9: for $w_i \in W$ 且 $c_j \in C$ do;
- 10: if $\bar{T}(w_i, c_j)$ 属于 Ξ 最高的三个或最低的三个类簇then
- 11: $\bar{T}(w_i, c_j)$ 被指定为词项特征 w_i 和晚点致因模式 c_j 为强关联关系或弱关联关系,并且将对应的 S 矩阵的值分别设置为正数(>1)或极小数(≈ 0 且 >0)。
- 12: else
- 13: $\bar{T}(w_i, c_j)$ 被指定为词项特征 w_i 和晚点致因模式 c_j 为复杂关联关系, S 的值由公式(3-3)求得。
- 14: end if
- 15: end for
- 16: 为每个晚点致因模式预分配 $m(m=10)$ 个相应的潜在主题特征 $z_{10*i}, z_{10*(i+1)}, \dots, z_{10*(i+m)}(1 \leq i \leq |C|)$ 。
- 17: $\Gamma(w_i, z_k) \leftarrow$ 初始化词项特征 w_i 与主题特征 $z_k \in Z$ 相关性为0。
- 18: for $w_i \in W$ 且 $z_k \in Z$ do;
- 19: if $z_k \in c_j$, then
- 20: 将 $S(w_i, c_j)$ 的值分配给 $\Gamma(w_i, z_k)$ 。
- 21: end if
- 22: end for

输入数据: 晚点描述词项特征 $w_i \in W$ 与晚点成因主题特征 $z_k \in Z$ 相关性集合 $\Gamma(w_i, z_k)$; 处置过程与处置写实记录(分词后) D ; 超参数 α, β ; 最优主题个数 K ; 迭代次数 $iter\ times$;

输出结果: 文档主题分布 θ_d ; LDA主题与词的分佈 ϕ_k 。

1: $K \leftarrow$ 模型最优主题个数: 根据主题连贯性指标求得 $K = 74$ 。

2: 选择合适的超参数向量 α, β 。

3: 变量申请:

概率向量 p ; ; 词在类上的分布 nw ; 每个类上的词的总数 $nwsum$; 每篇文章中, 各个类的词个数分布 nd ; 每篇文章中的词的总个数 $ndsum$; 每个词分派一个类 Z ; 文章 \rightarrow 类的概率分布 θ ; 类 \rightarrow 词的
概率分布 φ

4: 初始化阶段: 对应语料库中每一篇文档的每一个词, 随机的赋予一个主题编号 z 。

for x in 文章数do:

统计 $ndsum[文章id][词$ 的个数]

for y in 每篇文章的词个数do:

给所有词随机分派一个类

end for

end for

5: 重新扫描语料库, 对于不同相关关系的词项, 选用修正后的Gibbs采样公式更新其主题编号, 并更新语料中该词的编号。

for i in 迭代次数:

for m in 文章数do:

for v in 文章中词:

取 $topic = Z[m][v]$

判断词项与主题相关关系, 根据公式(3.6)(3.7)(3.8)计算概率 p

for k in (1,类的个数-1) do:

$p[k] += p[k-1]$

end for

再随机分派一次, 记录被分派的新的topic

end for

end for

end for

6: 重复上一步基于坐标轴轮换的Gibbs采样, 直到Gibbs采样收敛。

7: 统计语料库中各个文档各个词的主题, 得到文档主题分布 θ_d , 统计语料库中各个主题词的分布, 得到LDA主题与词的分佈 ϕ_k 。

