

主题模型数据分析

Yuhang Ge

June 9, 2021

1 数据预处理

通过读取所有 Json 文件，获取每个 Json 数据的人物简介字段，对人物简介字段进行中文文本分词、去停用词、去数字等预处理。下面是一个样例：

预处理前：

周晓文，1954 年出生于北京市，中国内地导演、编剧、摄影师、制作人，毕业于北京电影学院摄影系。1986 年，执导战争片《他们正年轻》，从而开启了他的导演生涯 [1]。1988 年，自编自导犯罪片《疯狂的代价》[2]，他凭借该片获得夏威夷国际电影节荣誉奖。1991 年，凭借执导的剧情片《青春无悔》获得上海大学生电影节最佳导演奖。1994 年，执导剧情片《二嫫》，该片获得洛迦诺国际电影节评审团大奖 [3]。1998 年，执导古装宫廷剧《吕后传奇》[4]。2000 年，执导古装武侠剧《天龙八部》。2002 年，担任古装剧《大脚马皇后》的导演 [5]。2006 年，执导都市冒险剧《逃亡香格里拉》[6]。2009 年，执导现代情感剧《晚婚》[7]。2011 年，担任剧情片《百合》的导演和制作人，该片获得第 14 届电影华表奖优秀数字电影奖 [8]。

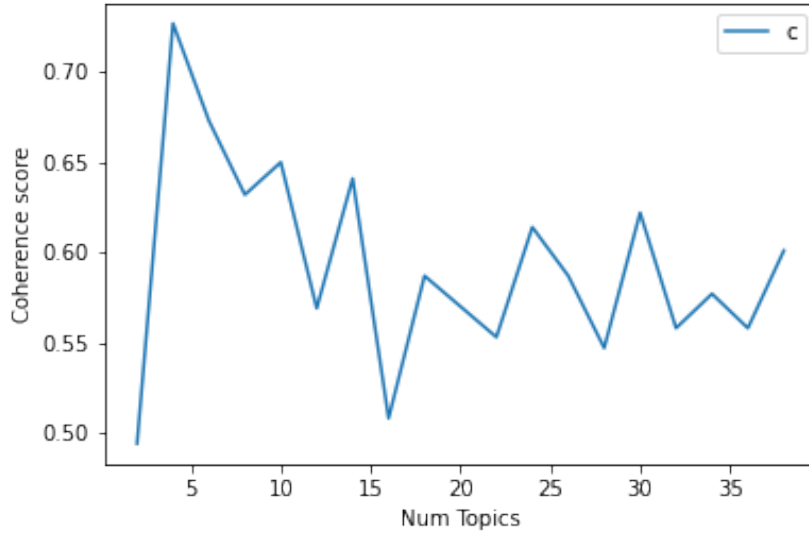
预处理后：

周晓文 年 出 生 于 北 京 市 中 国 内 地 导 演 编 剧 摄 影 师 制 作 毕 业 北 京 电 影 学 院 摄 影 系 执 导 战 争 片 年 轻 开 启 导 演 生 涯 自 编 自 导 犯 罪 疯 狂 代 价 该 片 获 得 夏 威 夷 国 际 电 影 节 荣 誉 奖 执 导 剧 情 片 青 春 无 悔 获 得 上 海 大 学 生 电 影 节 最 佳 导 演 奖 执 导 剧 情 片 该 片 获 得 洛 迦 诺 国 际 电 影 节 评 审 团 大 奖 执 导 古 装 宫 廷 吕 后 传 奇 执 导 古 装 武 侠 剧 天 龙 八 部 担 任 古 装 剧 大 脚 马 皇 后 导 演 执 导 都 市 冒 险 逃 亡 香 格 里 拉 执 导 现 代 情 感 晚 婚 担 任 剧 情 片 百 合 导 演 制 作 该 片 获 得 电 影 华 表 奖 优 秀 数 字 电 影

2 LDA 主题模型分析

2.1 最优主题数寻找

通过分别计算区间 $K = \{2 - 40\}$ 内主题的主题一致性 (coherence value)，绘制出下图，可以找到最优主题数为 $k = 6$ 。



将处理后的数据通过 LDA 主题模型进行数据分析，设置最优主题数为 6，并且每个主题的主题词数为 20，得到主题分布如下：

(0, '0.025*" 主演" + 0.018*" 出演" + 0.016*" 电视剧" + 0.015*" 饰演" + 0.012*" 电影" + '0.010*" 获得" + 0.010*" 中国" + 0.009*" 生于" + 0.008*" 同年" + 0.007*" 参演" + 0.006*" 古装" + 0.006*" 毕业" + 0.006*" 日出" + 0.006*" 都市" + 0.005*" 最佳" + 0.005*" 参加" + '0.005*" 个人" + 0.005*" 内地" + 0.005*" 情感" + 0.005*" 爱情"'),

(1, '0.027*" 主演" + 0.026*" 电影" + 0.016*" 获得" + 0.015*" 中国" + 0.011*" 最佳" + '0.009*" 电视剧" + 0.008*" 生于" + 0.007*" 执导" + 0.007*" 出演" + 0.007*" 饰演" + '0.006*" 同年" + 0.006*" 电影节" + 0.005*" 参演" + 0.005*" 导演" + 0.005*" 爱情" + '0.005*" 担任" + 0.005*" 上映" + 0.005*" 日出" + 0.005*" 毕业" + 0.004*" 公主"'),

(2, '0.015*" 中国" + 0.012*" 获得" + 0.011*" 专辑" + 0.008*" 发行" + 0.007*" 电视剧" + '0.007*" 生于" + 0.006*" 最佳" + 0.006*" 主演" + 0.006*" 电影" + 0.006*" 音乐" + 0.005*" 同年" + '0.005*" 演员" + 0.005*" 公元前" + 0.005*" 参加" + 0.005*" 日出" + 0.004*" 出演" + '0.004*" 担任" + 0.004*" 成为" + 0.004*" 推出" + 0.003*" 个人"'),

(3, '0.011*" 中国" + 0.010*" 获得" + 0.008*" 主演" + 0.006*" 生于" + 0.005*" 皇后" + 0.005*" 电影" + '0.005*" 最佳" + 0.005*" 电视剧" + 0.004*" 韩国" + 0.004*" 出演" + 0.004*" 担任" + '0.004*" 皇帝" + 0.004*" 元年" + 0.003*" 导演" + 0.003*" 时期"')

+ 0.003*” 日本” + 0.003*” 爱情” ’ ’+ 0.003*” 大赏” + 0.003*” 日出” + 0.003*” 女演员”’),

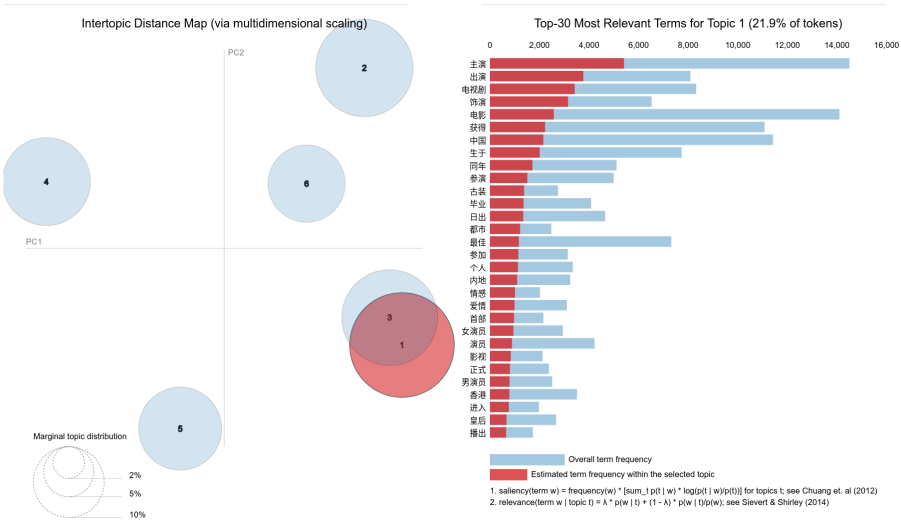
(4, ’0.022*” 电影” + 0.010*” 出演” + 0.009*” 参演” + 0.009*” 饰演” + 0.008*” 生于” + 0.007*” 中国” ’ ’+ 0.007*” 主演” + 0.007*” 获得” + 0.007*” 最佳” + 0.007*” 香港” + 0.006*” 电视剧” + ’0.005*” 演员” + 0.005*” 执导” + 0.004*” 日出” + 0.004*” 时期” + 0.004*” 成为” + 0.003*” 导演” ’ ’+ 0.003*” 谥号” + 0.003*” 毕业” + 0.003*” 大臣”’),

(5, ’0.017*” 电影” + 0.012*” 获得” + 0.011*” 最佳” + 0.010*” 中国” + 0.009*” 主演” + 0.009*” 生于” ’ ’+ 0.006*” 香港” + 0.006*” 演员” + 0.005*” 同年” + 0.005*” 参演” + 0.005*” 个人” + ’0.004*” 台湾” + 0.004*” 日出” + 0.004*” 电视剧” + 0.004*” 毕业” + 0.004*” 成为” + ’0.004*” 专辑” + 0.004*” 音乐” + 0.003*” 出演” + 0.003*” 男演员”’)

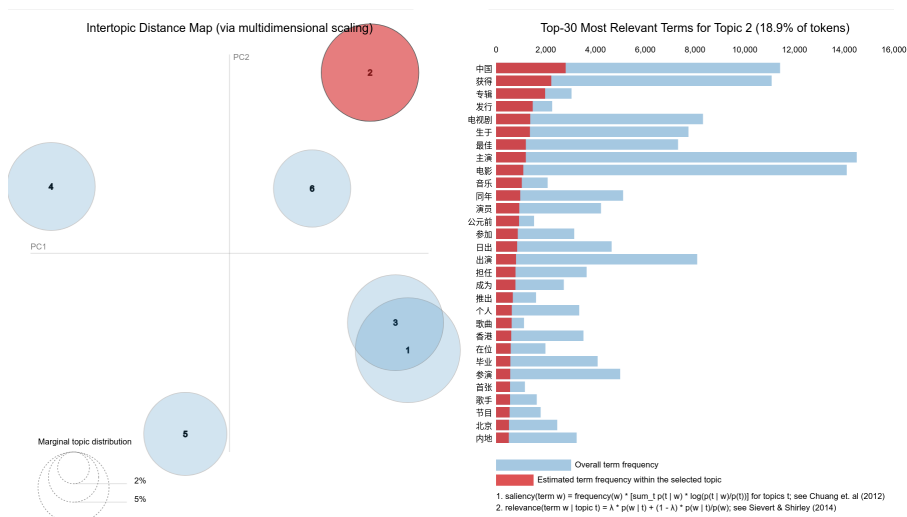
3 主题可视化

下面是分别 6 个主题分布的可视化图：

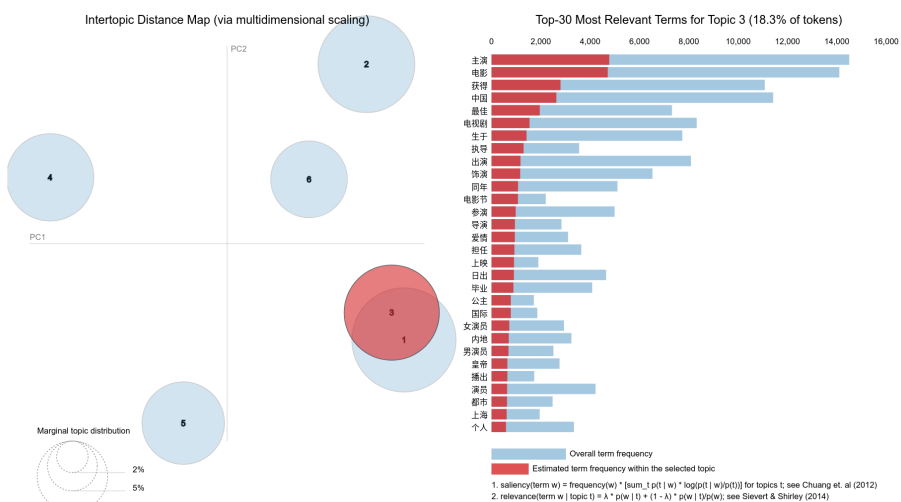
主题 1：



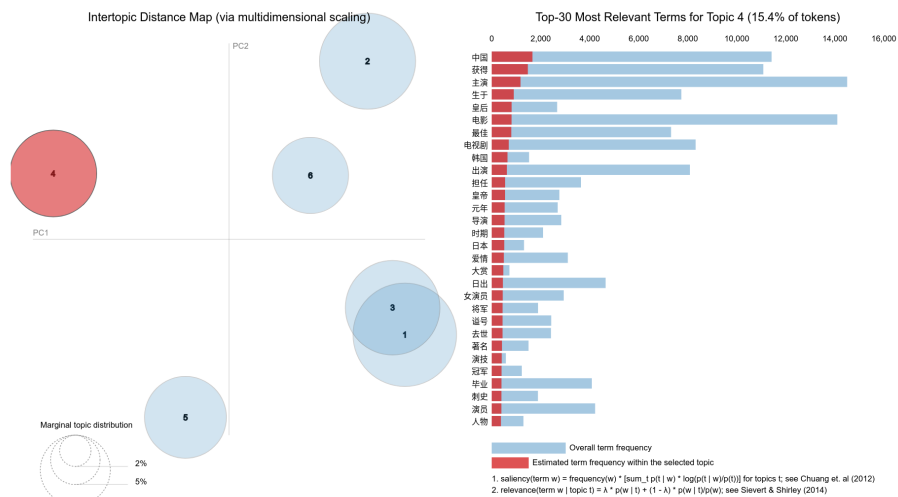
主题 2：



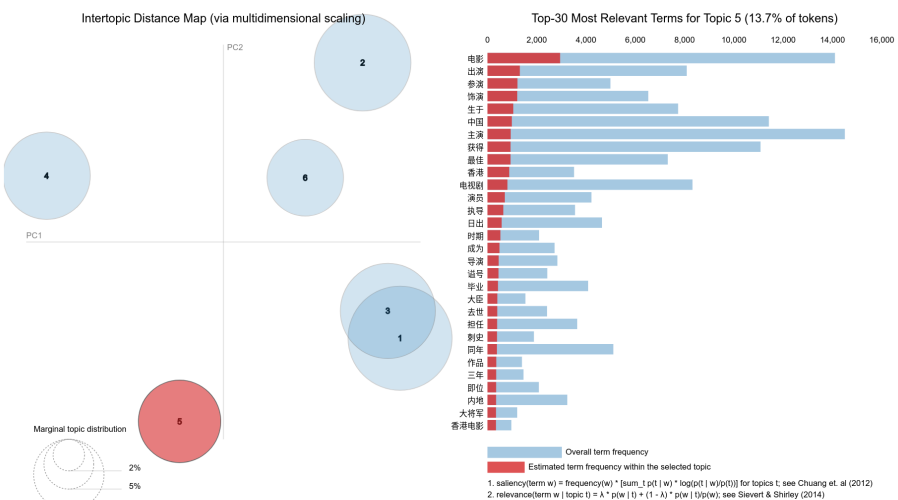
主题 3:



主题 4:



主题 5:



主题 6:

