

基于 word2vec 与 K-means 算法食品安全事件自动聚类研究

沈 思¹, 梁晓静²

(1.南京理工大学经济管理学院,江苏 南京 210094 2.金陵科技学院软件工程学院,江苏 南京 210000)

摘要:类别知识在食品安全事件的应急处理中不仅保证了相应应急措施推送的精准性而且可以提高食品安全事件应急处理的效率。为了更加精准的获取食品安全事件中的类别知识,结合统计、对比和分析的系统方法,通过食品安全事件的微博和信息数据,基于 word2vec 和 K-means 算法,文章探究了类别知识挖掘的具体流程、方法和相应的类别知识挖掘结果的评价等方面的相应问题。通过具体的实验,在所使用的数据上,确定了最优的聚类模型。

关键词:食品安全事件;word2vec;K-means;轮廓系数;TP393

中图分类号: 文献标识码:A 文章编号:1673-1131(2018)11-0008-03

Research of Automatic Clustering of Food Safety Events Based on Word2vec and K-means Algorithm

Si Shen¹, XiaoJing Liang²

(1.School of economics and management, Nanjing University of Science and Technology, Nanjing, 210094;

2. School of software engineering, Jinling institute of technology, Nanjing, 210000)

Abstract:The categories of knowledge not only ensures the precision of pushing the emergency measures but also can improve the efficiency of handling the food safety event emergency in handling the food safety event emergency. By using the method of statistics, comparison and analysis of system in the food safety incidents Weibo and information data, this paper explores the specific categories of knowledge mining process, method and the corresponding problems of evaluation of category knowledge mining results in order to more accurately get the category knowledge of food safety event. Through specific experiments, the best cluster model with category knowledge is determined based on the data used.

Keywords: Food safety event; Word2vec; K - means; Outline of the coefficient

0 引言

在食品安全事件的应急处理过程中,有效地知识是应对食品安全事件的关键因素,而再把有效的知识推送给用户的过程中,为了高效和精准的完成对应急知识的推送,需要获取有关食品安全事件的类别知识。目前获取类别知识最有效和主要的途径是通过聚类的方法对食品安全事件文本进行类别知识的获取。在上述这一背景下,针对相应的食品安全事件文本,通过把 word2vec 与 K-means 算法进行相融合的策略,针对领域食品安全事件的类别知识进行了相应的获取。

有关基于食品安全事件数据的探究目前主要停留在数据库的构建阶段,比较有代表性的研究如下。目前对食品安全事件的信息处理还停留在构建数据库的阶段,比较有代表性的研究如下。为了更好的完成对食品安全事件的预警,提高食品安全事件监管的科学性和高效性,张星联和吴云红分别提出^[1-2]提出构建食品安全事件动态数据库,并确保该数据库可以及时的更新,从而一方面可以有效地消除食品安全事件应对过程中的信息不对称问题,另一方面确保食品安全事件应对的透明度。在涉及到食品安全事件发生的整个过程中,食

情况提供了更有效的解决办法。通过多次实验证明,井下瞬变电磁仪在含水异常体探测中效果明显,结果准确,值得大量推广使用。该探测方法主要有以下几个优势:

(1)施工简单、仪器便于携带,工作时间短,不影响矿井正常生产作业,即能满足防治水工作需要,又不影响生产需要。

(2)井下由于特殊的支护环境,金属质地的网片、锚杆等难免会对仪器造成影响,故一定要将影响因素降到最低,孔中瞬变恰巧在一定程度上解决了这个问题。

(3)在不需要特殊布钻的情况下,将仪器接收探头送入现有钻孔,不仅不受周围环境的影响,而且使得仪器离异常更近,解释更清晰,更容易确定异常位置。

(4)降低了矿方盲目打钻的工作量,为矿方节省了成本。

参考文献:

- [1] 肖杰,焦当云.瞬变电磁法在煤矿防治水工作中的应用[J].江西煤炭科技,2018(01):113-114+119.
- [2] 王卫光.瞬变电磁仪在煤矿水文地质探测方面的应用研究[J].机电工程技术,2015,44(07):225
- [3] 薛国强,于于.瞬变电磁法在煤炭领域的研究与应用新

进展[J].地球物理学进展,2017,32(01):319-326.

- [4] 储韬玉.矿井孔中瞬变电磁测量方法及其应用研究[D].中国矿业大学,2015.
- [5] 陈丁.矿井全空间巷道孔中瞬变电磁波场特征数值模拟研究[D].中国矿业大学(北京),2016.
- [6] 张青.瞬变电磁法在瓦斯隧道隧底采空区探测中的应用[J].科技风,2018(16):91.
- [7] 王玺凯,王福生.瞬变电磁法在某水源井勘测中的应用[J].科技创新与应用,2018(15):171-172.
- [8] 李健军.瞬变电磁仪法与超前水平钻孔法在隧道施工中的综合运用[J].价值工程,2017,36(32):118-120.
- [9] 谢成梁.瞬变电磁仪在工作面底板水害探测中的应用[J].自动化与仪器仪表,2017(10):210-211+214.
- [10] 孟宪杨.巷道探测中瞬变电磁仪的应用[J].矿业装备,2017(05):52-53.

作者简介: 邓立博(1990-),男,陕西渭南人,助理研究员,硕士,研究方向:地球物理探测仪器开发。

品的加工是极其重要的一环,如何确保食品加工的安全性、可跟踪性和回溯性是控制食品安全事件发生的关键。基于此,余清等^[3]提出了构建食品加工风险数据库,为食品安全事件的检测奠定基础。在食品安全事件当中,生鲜农产品不仅涉及范围广泛而且在时间上对生鲜农产品也提出了更高的监管要求,鉴于此,以彭州市下面的一个镇为数据源,贾凯等^[4]提出了构建生鲜农产品数据库,从而实现对涉及到该镇的所有生鲜农产品可以实现有效的回溯。基于大数据的技术和方法,刘翠玲^[5]提出了一种应对食品安全事件的大数据监管平台,并基于该平台对农药残留这一食品安全事件进行了具体的对比探究和分析。

在对食品安全事件进行聚类的过程中,所选择的聚类方法应该充分考虑到中文文本的特殊性,以便于从食品安全事件的数据中挖掘出有效和针对性比较强的类别知识。通过逐一地计算样本数据的距离,夏长辉^[6]基于最大的距离样本为初始的聚类焦点,利用 k-means 聚类算法的特征逐步地完成对所选择数据的聚类,从而从文本中挖掘出相应的聚类知识,这种方法的缺陷是计算量比较大。针对微生物群落的数据分布,王侠林、贺建峰^[7]基于 OTUs 数据集的 27 个样本进行了聚类的探究,并通过与 PCA 这一方法的对比,发现 K-Means 的整体聚类效果更加突出。充分利用 HowNet 当中领域语义词汇知识,结合文档集合所呈现出来的特征空间,黄建宇和周爱武^[8]等给出了优化后的 k 值的 k-means 聚类策略,这在一定程度上有效提高了聚类的整体效果,但由于所依赖领域词汇库的局限性,在一定程度上限制了聚类效果的进一步提升。在云计算平台 Hadoop 下的并行 k-means 聚类算法基础上,赵卫中^[9]结合不同的数据源,充分而具体的验证了在大数据平台上聚类算法的加速比、扩展率和伸缩率等不同指标的整体性能。结合大数据的技术、方法和理念,利用最小相似度文本对的理念,通过构建 MapReduce 基础上的文本聚类模型,武森和冯小东^[10]提出了融合大数据的 k-means 聚类方法,所获取的聚类结果整体性能非常突出,并且具有较强的可迁移性。这一聚类的理念可以有效地拓展到以汉语为呈现形式上的非结构化文本上。

1 模型和食品安全事件文本语料介绍

为了从食品安全事件文本中获取更加有效的特征或者信息,本研究基于 word2vec^[11-13]获取食品安全事件文本中词汇的向量值,该算法是由 2013 年的谷歌公司提出的,同时对该算法进行工具包式的封装。该算法主要由 CBOW(Continuous Bag-Of-Words)和 Skip-gram 这两个具体的模型构成,这两种模型均通过某一当前的词汇去预测其他词汇,从而提高词汇单位关联的精准性。word2vec 算法最为突出的一点是可以进行向量间的有效转换,从而极大地提升了对文本进行语义相似度计算的精准性。基于这一计算出来的精准度较高的相似度值,对于改进相应的文本聚类、同义词查询、检索系统的构建、相似度情感词计算等探究具有较为重要的意义和价值。该算法的具体原理如公式 1 所示。

$$p(w_t|w_{t-2}^c, \dots, w_{t-1}, w_{t+1}, \dots, w_{t+2}^c) = \prod_{i=1}^c p(\text{context}(w_i)) \quad (1)$$

其中 CBOW 由输入层、投影层和输出层这三个层次构成,并且相互之间形成了有效的关联和嵌套。其中输入层是词向量的表示,投影层主要的功能是对输入的向量进行叠加,所使用的公式为 $\sum_{i=1}^c (\text{context}(w_i))$, $\text{context}(w)$ 表示词的概率值

而输出层是通过一个树形结构进行表示的,具体为二叉树。在具体实现的时候是通过构建 Huffman 树来完成的。随着 word2vec 被广泛应用于文本聚类、信息检索和主题挖掘等研究当中,各种封装的 word2vec 版本均被大家在不同的具体探究上进行了应用。结合本文与食品安全事件信息聚类相关的研究,gensim^[14]基础上的 word2vec 被选取应用于本文具体的探究当中。

在具体聚类算法选择上,基于本文食品安全事件文本数据的特征,k-means 聚类算法被选为针对食品安全事件数据的聚类知识挖掘算法,该聚类算法不仅设计加单,而且时间复杂度相对较低,并且可以有效地兼容不同的类型的数据。为了更好地利用该算法,针对食品安全事件数据,轮廓系数^[15-16]被用来对初始聚类值进行设定,具体轮廓系数的计算公式如下:

$$s_i = \frac{b_i - a_i}{\max(a_i, b_i)} \quad (2)$$

在这一公式中 a_i 代表某一样本到其他潜在统一类别样本的平均距离, b_i 代表某一样本表到与其最相近的类内部平均距离,因此 a_i 值越小则说明某一文本分类越合理,而 b_i 越大,则说明某一样本不应当归到其他类比,从而 s_i 值越靠近 1,则说明某一样本聚类合理,而 s_i 若靠近 0,说明某一样本聚类不合理。通过对所有样本的轮廓系数就平均值,进而会得到全局轮廓系数,具体计算公式如下。

$$s_k = \frac{1}{n} \sum_{i=1}^n s_i \quad (3)$$

通过关键词“甲醛白菜”获取了微博上相关的信息共 7901 条,具体年份如表 1 所示。首先对相应的数据进行了清洗,去除了编码不一致和有乱码的文本。其次,调用 pynlpir 包分词,去除“【”和“】”符号,建立二维列表,满足 word2vec 进行计算的需要。最后,对所获取的数据进行了统一的存储,具体经过分词后的数据样例如图 1 所示。

表 1 新浪微博中“甲醛白菜”主题数据集合

相关微博发表年份	去重前文本记录数	去重后文本记录数
2012	6547	5090
2013	2919	1948
2014	2057	53
总和	11524	7091

// @熊吉堂 DD: 这个猪有, 质量不好, 橡胶和甲醛味重。// @陈金武: @白菜党 @成都 568 休闲骑行 @中国人就是一化学周周表——2006 年, 苏丹红; 2008 年, 三聚氰胺; 2009 年, 瘦肉精; 2010 年, 地沟油; 2011 年, 塑化剂; 2012 年, 甲醛白菜; 2013 年, 毒胶囊; 2014 年, 毒奶粉; 2015 年, 毒大米; 2016 年, 毒鸡蛋; 2017 年, 毒月饼; 2018 年, 毒月饼; 2019 年, 毒月饼; 2020 年, 毒月饼; 2021 年, 毒月饼; 2022 年, 毒月饼; 2023 年, 毒月饼; 2024 年, 毒月饼; 2025 年, 毒月饼; 2026 年, 毒月饼; 2027 年, 毒月饼; 2028 年, 毒月饼; 2029 年, 毒月饼; 2030 年, 毒月饼; 2031 年, 毒月饼; 2032 年, 毒月饼; 2033 年, 毒月饼; 2034 年, 毒月饼; 2035 年, 毒月饼; 2036 年, 毒月饼; 2037 年, 毒月饼; 2038 年, 毒月饼; 2039 年, 毒月饼; 2040 年, 毒月饼; 2041 年, 毒月饼; 2042 年, 毒月饼; 2043 年, 毒月饼; 2044 年, 毒月饼; 2045 年, 毒月饼; 2046 年, 毒月饼; 2047 年, 毒月饼; 2048 年, 毒月饼; 2049 年, 毒月饼; 2050 年, 毒月饼; 2051 年, 毒月饼; 2052 年, 毒月饼; 2053 年, 毒月饼; 2054 年, 毒月饼; 2055 年, 毒月饼; 2056 年, 毒月饼; 2057 年, 毒月饼; 2058 年, 毒月饼; 2059 年, 毒月饼; 2060 年, 毒月饼; 2061 年, 毒月饼; 2062 年, 毒月饼; 2063 年, 毒月饼; 2064 年, 毒月饼; 2065 年, 毒月饼; 2066 年, 毒月饼; 2067 年, 毒月饼; 2068 年, 毒月饼; 2069 年, 毒月饼; 2070 年, 毒月饼; 2071 年, 毒月饼; 2072 年, 毒月饼; 2073 年, 毒月饼; 2074 年, 毒月饼; 2075 年, 毒月饼; 2076 年, 毒月饼; 2077 年, 毒月饼; 2078 年, 毒月饼; 2079 年, 毒月饼; 2080 年, 毒月饼; 2081 年, 毒月饼; 2082 年, 毒月饼; 2083 年, 毒月饼; 2084 年, 毒月饼; 2085 年, 毒月饼; 2086 年, 毒月饼; 2087 年, 毒月饼; 2088 年, 毒月饼; 2089 年, 毒月饼; 2090 年, 毒月饼; 2091 年, 毒月饼; 2092 年, 毒月饼; 2093 年, 毒月饼; 2094 年, 毒月饼; 2095 年, 毒月饼; 2096 年, 毒月饼; 2097 年, 毒月饼; 2098 年, 毒月饼; 2099 年, 毒月饼; 2100 年, 毒月饼; 2101 年, 毒月饼; 2102 年, 毒月饼; 2103 年, 毒月饼; 2104 年, 毒月饼; 2105 年, 毒月饼; 2106 年, 毒月饼; 2107 年, 毒月饼; 2108 年, 毒月饼; 2109 年, 毒月饼; 2110 年, 毒月饼; 2111 年, 毒月饼; 2112 年, 毒月饼; 2113 年, 毒月饼; 2114 年, 毒月饼; 2115 年, 毒月饼; 2116 年, 毒月饼; 2117 年, 毒月饼; 2118 年, 毒月饼; 2119 年, 毒月饼; 2120 年, 毒月饼; 2121 年, 毒月饼; 2122 年, 毒月饼; 2123 年, 毒月饼; 2124 年, 毒月饼; 2125 年, 毒月饼; 2126 年, 毒月饼; 2127 年, 毒月饼; 2128 年, 毒月饼; 2129 年, 毒月饼; 2130 年, 毒月饼; 2131 年, 毒月饼; 2132 年, 毒月饼; 2133 年, 毒月饼; 2134 年, 毒月饼; 2135 年, 毒月饼; 2136 年, 毒月饼; 2137 年, 毒月饼; 2138 年, 毒月饼; 2139 年, 毒月饼; 2140 年, 毒月饼; 2141 年, 毒月饼; 2142 年, 毒月饼; 2143 年, 毒月饼; 2144 年, 毒月饼; 2145 年, 毒月饼; 2146 年, 毒月饼; 2147 年, 毒月饼; 2148 年, 毒月饼; 2149 年, 毒月饼; 2150 年, 毒月饼; 2151 年, 毒月饼; 2152 年, 毒月饼; 2153 年, 毒月饼; 2154 年, 毒月饼; 2155 年, 毒月饼; 2156 年, 毒月饼; 2157 年, 毒月饼; 2158 年, 毒月饼; 2159 年, 毒月饼; 2160 年, 毒月饼; 2161 年, 毒月饼; 2162 年, 毒月饼; 2163 年, 毒月饼; 2164 年, 毒月饼; 2165 年, 毒月饼; 2166 年, 毒月饼; 2167 年, 毒月饼; 2168 年, 毒月饼; 2169 年, 毒月饼; 2170 年, 毒月饼; 2171 年, 毒月饼; 2172 年, 毒月饼; 2173 年, 毒月饼; 2174 年, 毒月饼; 2175 年, 毒月饼; 2176 年, 毒月饼; 2177 年, 毒月饼; 2178 年, 毒月饼; 2179 年, 毒月饼; 2180 年, 毒月饼; 2181 年, 毒月饼; 2182 年, 毒月饼; 2183 年, 毒月饼; 2184 年, 毒月饼; 2185 年, 毒月饼; 2186 年, 毒月饼; 2187 年, 毒月饼; 2188 年, 毒月饼; 2189 年, 毒月饼; 2190 年, 毒月饼; 2191 年, 毒月饼; 2192 年, 毒月饼; 2193 年, 毒月饼; 2194 年, 毒月饼; 2195 年, 毒月饼; 2196 年, 毒月饼; 2197 年, 毒月饼; 2198 年, 毒月饼; 2199 年, 毒月饼; 2200 年, 毒月饼; 2201 年, 毒月饼; 2202 年, 毒月饼; 2203 年, 毒月饼; 2204 年, 毒月饼; 2205 年, 毒月饼; 2206 年, 毒月饼; 2207 年, 毒月饼; 2208 年, 毒月饼; 2209 年, 毒月饼; 2210 年, 毒月饼; 2211 年, 毒月饼; 2212 年, 毒月饼; 2213 年, 毒月饼; 2214 年, 毒月饼; 2215 年, 毒月饼; 2216 年, 毒月饼; 2217 年, 毒月饼; 2218 年, 毒月饼; 2219 年, 毒月饼; 2220 年, 毒月饼; 2221 年, 毒月饼; 2222 年, 毒月饼; 2223 年, 毒月饼; 2224 年, 毒月饼; 2225 年, 毒月饼; 2226 年, 毒月饼; 2227 年, 毒月饼; 2228 年, 毒月饼; 2229 年, 毒月饼; 2230 年, 毒月饼; 2231 年, 毒月饼; 2232 年, 毒月饼; 2233 年, 毒月饼; 2234 年, 毒月饼; 2235 年, 毒月饼; 2236 年, 毒月饼; 2237 年, 毒月饼; 2238 年, 毒月饼; 2239 年, 毒月饼; 2240 年, 毒月饼; 2241 年, 毒月饼; 2242 年, 毒月饼; 2243 年, 毒月饼; 2244 年, 毒月饼; 2245 年, 毒月饼; 2246 年, 毒月饼; 2247 年, 毒月饼; 2248 年, 毒月饼; 2249 年, 毒月饼; 2250 年, 毒月饼; 2251 年, 毒月饼; 2252 年, 毒月饼; 2253 年, 毒月饼; 2254 年, 毒月饼; 2255 年, 毒月饼; 2256 年, 毒月饼; 2257 年, 毒月饼; 2258 年, 毒月饼; 2259 年, 毒月饼; 2260 年, 毒月饼; 2261 年, 毒月饼; 2262 年, 毒月饼; 2263 年, 毒月饼; 2264 年, 毒月饼; 2265 年, 毒月饼; 2266 年, 毒月饼; 2267 年, 毒月饼; 2268 年, 毒月饼; 2269 年, 毒月饼; 2270 年, 毒月饼; 2271 年, 毒月饼; 2272 年, 毒月饼; 2273 年, 毒月饼; 2274 年, 毒月饼; 2275 年, 毒月饼; 2276 年, 毒月饼; 2277 年, 毒月饼; 2278 年, 毒月饼; 2279 年, 毒月饼; 2280 年, 毒月饼; 2281 年, 毒月饼; 2282 年, 毒月饼; 2283 年, 毒月饼; 2284 年, 毒月饼; 2285 年, 毒月饼; 2286 年, 毒月饼; 2287 年, 毒月饼; 2288 年, 毒月饼; 2289 年, 毒月饼; 2290 年, 毒月饼; 2291 年, 毒月饼; 2292 年, 毒月饼; 2293 年, 毒月饼; 2294 年, 毒月饼; 2295 年, 毒月饼; 2296 年, 毒月饼; 2297 年, 毒月饼; 2298 年, 毒月饼; 2299 年, 毒月饼; 2300 年, 毒月饼; 2301 年, 毒月饼; 2302 年, 毒月饼; 2303 年, 毒月饼; 2304 年, 毒月饼; 2305 年, 毒月饼; 2306 年, 毒月饼; 2307 年, 毒月饼; 2308 年, 毒月饼; 2309 年, 毒月饼; 2310 年, 毒月饼; 2311 年, 毒月饼; 2312 年, 毒月饼; 2313 年, 毒月饼; 2314 年, 毒月饼; 2315 年, 毒月饼; 2316 年, 毒月饼; 2317 年, 毒月饼; 2318 年, 毒月饼; 2319 年, 毒月饼; 2320 年, 毒月饼; 2321 年, 毒月饼; 2322 年, 毒月饼; 2323 年, 毒月饼; 2324 年, 毒月饼; 2325 年, 毒月饼; 2326 年, 毒月饼; 2327 年, 毒月饼; 2328 年, 毒月饼; 2329 年, 毒月饼; 2330 年, 毒月饼; 2331 年, 毒月饼; 2332 年, 毒月饼; 2333 年, 毒月饼; 2334 年, 毒月饼; 2335 年, 毒月饼; 2336 年, 毒月饼; 2337 年, 毒月饼; 2338 年, 毒月饼; 2339 年, 毒月饼; 2340 年, 毒月饼; 2341 年, 毒月饼; 2342 年, 毒月饼; 2343 年, 毒月饼; 2344 年, 毒月饼; 2345 年, 毒月饼; 2346 年, 毒月饼; 2347 年, 毒月饼; 2348 年, 毒月饼; 2349 年, 毒月饼; 2350 年, 毒月饼; 2351 年, 毒月饼; 2352 年, 毒月饼; 2353 年, 毒月饼; 2354 年, 毒月饼; 2355 年, 毒月饼; 2356 年, 毒月饼; 2357 年, 毒月饼; 2358 年, 毒月饼; 2359 年, 毒月饼; 2360 年, 毒月饼; 2361 年, 毒月饼; 2362 年, 毒月饼; 2363 年, 毒月饼; 2364 年, 毒月饼; 2365 年, 毒月饼; 2366 年, 毒月饼; 2367 年, 毒月饼; 2368 年, 毒月饼; 2369 年, 毒月饼; 2370 年, 毒月饼; 2371 年, 毒月饼; 2372 年, 毒月饼; 2373 年, 毒月饼; 2374 年, 毒月饼; 2375 年, 毒月饼; 2376 年, 毒月饼; 2377 年, 毒月饼; 2378 年, 毒月饼; 2379 年, 毒月饼; 2380 年, 毒月饼; 2381 年, 毒月饼; 2382 年, 毒月饼; 2383 年, 毒月饼; 2384 年, 毒月饼; 2385 年, 毒月饼; 2386 年, 毒月饼; 2387 年, 毒月饼; 2388 年, 毒月饼; 2389 年, 毒月饼; 2390 年, 毒月饼; 2391 年, 毒月饼; 2392 年, 毒月饼; 2393 年, 毒月饼; 2394 年, 毒月饼; 2395 年, 毒月饼; 2396 年, 毒月饼; 2397 年, 毒月饼; 2398 年, 毒月饼; 2399 年, 毒月饼; 2400 年, 毒月饼; 2401 年, 毒月饼; 2402 年, 毒月饼; 2403 年, 毒月饼; 2404 年, 毒月饼; 2405 年, 毒月饼; 2406 年, 毒月饼; 2407 年, 毒月饼; 2408 年, 毒月饼; 2409 年, 毒月饼; 2410 年, 毒月饼; 2411 年, 毒月饼; 2412 年, 毒月饼; 2413 年, 毒月饼; 2414 年, 毒月饼; 2415 年, 毒月饼; 2416 年, 毒月饼; 2417 年, 毒月饼; 2418 年, 毒月饼; 2419 年, 毒月饼; 2420 年, 毒月饼; 2421 年, 毒月饼; 2422 年, 毒月饼; 2423 年, 毒月饼; 2424 年, 毒月饼; 2425 年, 毒月饼; 2426 年, 毒月饼; 2427 年, 毒月饼; 2428 年, 毒月饼; 2429 年, 毒月饼; 2430 年, 毒月饼; 2431 年, 毒月饼; 2432 年, 毒月饼; 2433 年, 毒月饼; 2434 年, 毒月饼; 2435 年, 毒月饼; 2436 年, 毒月饼; 2437 年, 毒月饼; 2438 年, 毒月饼; 2439 年, 毒月饼; 2440 年, 毒月饼; 2441 年, 毒月饼; 2442 年, 毒月饼; 2443 年, 毒月饼; 2444 年, 毒月饼; 2445 年, 毒月饼; 2446 年, 毒月饼; 2447 年, 毒月饼; 2448 年, 毒月饼; 2449 年, 毒月饼; 2450 年, 毒月饼; 2451 年, 毒月饼; 2452 年, 毒月饼; 2453 年, 毒月饼; 2454 年, 毒月饼; 2455 年, 毒月饼; 2456 年, 毒月饼; 2457 年, 毒月饼; 2458 年, 毒月饼; 2459 年, 毒月饼; 2460 年, 毒月饼; 2461 年, 毒月饼; 2462 年, 毒月饼; 2463 年, 毒月饼; 2464 年, 毒月饼; 2465 年, 毒月饼; 2466 年, 毒月饼; 2467 年, 毒月饼; 2468 年, 毒月饼; 2469 年, 毒月饼; 2470 年, 毒月饼; 2471 年, 毒月饼; 2472 年, 毒月饼; 2473 年, 毒月饼; 2474 年, 毒月饼; 2475 年, 毒月饼; 2476 年, 毒月饼; 2477 年, 毒月饼; 2478 年, 毒月饼; 2479 年, 毒月饼; 2480 年, 毒月饼; 2481 年, 毒月饼; 2482 年, 毒月饼; 2483 年, 毒月饼; 2484 年, 毒月饼; 2485 年, 毒月饼; 2486 年, 毒月饼; 2487 年, 毒月饼; 2488 年, 毒月饼; 2489 年, 毒月饼; 2490 年, 毒月饼; 2491 年, 毒月饼; 2492 年, 毒月饼; 2493 年, 毒月饼; 2494 年, 毒月饼; 2495 年, 毒月饼; 2496 年, 毒月饼; 2497 年, 毒月饼; 2498 年, 毒月饼; 2499 年, 毒月饼; 2500 年, 毒月饼; 2501 年, 毒月饼; 2502 年, 毒月饼; 2503 年, 毒月饼; 2504 年, 毒月饼; 2505 年, 毒月饼; 2506 年, 毒月饼; 2507 年, 毒月饼; 2508 年, 毒月饼; 2509 年, 毒月饼; 2510 年, 毒月饼; 2511 年, 毒月饼; 2512 年, 毒月饼; 2513 年, 毒月饼; 2514 年, 毒月饼; 2515 年, 毒月饼; 2516 年, 毒月饼; 2517 年, 毒月饼; 2518 年, 毒月饼; 2519 年, 毒月饼; 2520 年, 毒月饼; 2521 年, 毒月饼; 2522 年, 毒月饼; 2523 年, 毒月饼; 2524 年, 毒月饼; 2525 年, 毒月饼; 2526 年, 毒月饼; 2527 年, 毒月饼; 2528 年, 毒月饼; 2529 年, 毒月饼; 2530 年, 毒月饼; 2531 年, 毒月饼; 2532 年, 毒月饼; 2533 年, 毒月饼; 2534 年, 毒月饼; 2535 年, 毒月饼; 2536 年, 毒月饼; 2537 年, 毒月饼; 2538 年, 毒月饼; 2539 年, 毒月饼; 2540 年, 毒月饼; 2541 年, 毒月饼; 2542 年, 毒月饼; 2543 年, 毒月饼; 2544 年, 毒月饼; 2545 年, 毒月饼; 2546 年, 毒月饼; 2547 年, 毒月饼; 2548 年, 毒月饼; 2549 年, 毒月饼; 2550 年, 毒月饼; 2551 年, 毒月饼; 2552 年, 毒月饼; 2553 年, 毒月饼; 2554 年, 毒月饼; 2555 年, 毒月饼; 2556 年, 毒月饼; 2557 年, 毒月饼; 2558 年, 毒月饼; 2559 年, 毒月饼; 2560 年, 毒月饼; 2561 年, 毒月饼; 2562 年, 毒月饼; 2563 年, 毒月饼; 2564 年, 毒月饼; 2565 年, 毒月饼; 2566 年, 毒月饼; 2567 年, 毒月饼; 2568 年, 毒月饼; 2569 年, 毒月饼; 2570 年, 毒月饼; 2571 年, 毒月饼; 2572 年, 毒月饼; 2573 年, 毒月饼; 2574 年, 毒月饼; 2575 年, 毒月饼; 2576 年, 毒月饼; 2577 年, 毒月饼; 2578 年, 毒月饼; 2579 年, 毒月饼; 2580 年, 毒月饼; 2581 年, 毒月饼; 2582 年, 毒月饼; 2583 年, 毒月饼; 2584 年, 毒月饼; 2585 年, 毒月饼; 2586 年, 毒月饼; 2587 年, 毒月饼; 2588 年, 毒月饼; 2589 年, 毒月饼; 2590 年, 毒月饼; 2591 年, 毒月饼; 2592 年, 毒月饼; 2593 年, 毒月饼; 2594 年, 毒月饼; 2595 年, 毒月饼; 2596 年, 毒月饼; 2597 年, 毒月饼; 2598 年, 毒月饼; 2599 年, 毒月饼; 2600 年, 毒月饼; 2601 年, 毒月饼; 2602 年, 毒月饼; 2603 年, 毒月饼; 2604 年, 毒月饼; 2605 年, 毒月饼; 2606 年, 毒月饼; 2607 年, 毒月饼; 2608 年, 毒月饼; 2609 年, 毒月饼; 2610 年, 毒月饼; 2611 年, 毒月饼; 2612 年, 毒月饼; 2613 年, 毒月饼; 2614 年, 毒月饼; 2615 年, 毒月饼; 2616 年, 毒月饼; 2617 年, 毒月饼; 2618 年, 毒月饼; 2619 年, 毒月饼; 2620 年, 毒月饼; 2621 年, 毒月饼; 2622 年, 毒月饼; 2623 年, 毒月饼; 2624 年, 毒月饼; 2625 年, 毒月饼; 2626 年, 毒月饼; 2627 年, 毒月饼; 2628 年, 毒月饼; 2629 年, 毒月饼; 2630 年, 毒月饼; 2631 年, 毒月饼; 2632 年, 毒月饼; 2633 年, 毒月饼; 2634 年, 毒月饼; 2635 年, 毒月饼; 2636 年, 毒月饼; 2637 年, 毒月饼; 2638 年, 毒月饼; 2639 年, 毒月饼; 2640 年, 毒月饼; 2641 年, 毒月饼; 2642 年, 毒月饼; 2643 年, 毒月饼; 2644 年, 毒月饼; 2645 年, 毒月饼; 2646 年, 毒月饼; 2647 年, 毒月饼; 2648 年, 毒月饼; 2649 年, 毒月饼; 2650 年, 毒月饼; 2651 年, 毒月饼; 2652 年, 毒月饼; 2653 年, 毒月饼; 2654 年, 毒月饼; 2655 年, 毒月饼; 2656 年, 毒月饼; 2657 年, 毒月饼; 2658 年, 毒月饼; 2659 年, 毒月饼; 2660 年, 毒月饼; 2661 年, 毒月饼; 2662 年, 毒月饼; 2663 年, 毒月饼; 2664 年, 毒月饼; 2665 年, 毒月饼; 2666 年, 毒月饼; 2667 年, 毒月饼; 2668 年, 毒月饼; 2669 年, 毒月饼; 2670 年, 毒月饼; 2671 年, 毒月饼; 2672 年, 毒月饼; 2673 年, 毒月饼; 2674 年, 毒月饼; 2675 年, 毒月饼; 2676 年, 毒月饼; 2677 年, 毒月饼; 2678 年, 毒月饼; 2679 年, 毒月饼; 2680 年, 毒月饼; 2681 年, 毒月饼; 2682 年, 毒月饼; 2683 年, 毒月饼; 2684 年, 毒月饼; 2685 年, 毒月饼; 2686 年, 毒月饼; 2687 年, 毒月饼; 2688 年, 毒月饼; 2689 年, 毒月饼; 2690 年, 毒月饼; 2691 年, 毒月饼; 2692 年, 毒月饼; 2693 年, 毒月饼; 2694 年, 毒月饼; 2695 年, 毒月饼; 2696 年, 毒月饼; 2697 年, 毒月饼; 2698 年, 毒月饼; 2699 年, 毒月饼; 2700 年, 毒月饼; 2701 年, 毒月饼; 2702 年, 毒月饼; 2703 年, 毒月饼; 2704 年, 毒月饼; 2705 年, 毒月饼; 2706 年, 毒月饼; 2707 年, 毒月饼; 2708 年, 毒月饼; 2709 年, 毒月饼; 2710 年, 毒月饼; 2711 年, 毒月饼; 2712 年, 毒月饼; 2713 年, 毒月饼; 2714 年, 毒月饼; 2715 年, 毒月饼; 2716 年, 毒月饼; 2717 年, 毒月饼; 2718 年, 毒月饼; 2719 年, 毒月饼; 2720 年, 毒月饼; 2721 年, 毒月饼; 2722 年, 毒月饼; 2723 年, 毒月饼; 2724 年, 毒月饼; 2725 年, 毒月饼; 2726 年, 毒月饼; 2727 年, 毒月饼; 2728 年, 毒月饼; 2729 年, 毒月饼; 2730 年, 毒月饼; 2731 年, 毒月饼; 2732 年, 毒月饼; 2733 年, 毒月饼; 2734 年, 毒月饼; 2735 年, 毒月饼; 2736 年, 毒月饼; 2737 年, 毒月饼; 2738 年, 毒月饼; 2739 年, 毒月饼; 2740 年, 毒月饼; 2741 年, 毒月饼; 2742 年, 毒月饼; 2743 年, 毒月饼; 2744 年, 毒月饼; 2745 年, 毒月饼; 2746 年, 毒月饼; 2747 年, 毒月饼; 2748 年, 毒月饼; 2749 年, 毒月饼; 2750 年, 毒月饼; 2751 年, 毒月饼; 2752 年, 毒月饼; 2753 年, 毒月饼; 2754 年, 毒月饼; 2755 年, 毒月饼; 2756 年, 毒月饼; 2757 年, 毒月饼; 2758 年, 毒月饼; 2759 年, 毒月饼; 2760 年, 毒月饼; 2761 年, 毒月饼; 2762 年, 毒月饼; 2763 年, 毒月饼; 2764 年, 毒月饼; 2765 年, 毒月饼; 2766 年, 毒月饼; 2767 年, 毒月饼; 2768 年, 毒月饼; 2769 年, 毒月饼; 2770 年, 毒月饼; 2771 年, 毒月饼; 2772 年, 毒月饼; 2773 年, 毒月饼; 2774 年, 毒月饼; 2775 年, 毒月饼; 2776 年, 毒月饼; 2777 年, 毒月饼; 2778 年, 毒月饼; 2779 年, 毒月饼; 2780 年, 毒月饼; 2781 年, 毒月饼; 2782 年, 毒月饼; 2783 年, 毒月饼; 2784 年, 毒月饼; 2785 年, 毒月饼; 2786 年, 毒月饼; 2787 年, 毒月饼; 2788 年, 毒月饼; 2789 年, 毒月饼; 2790 年, 毒月饼; 2791 年, 毒月饼; 2792 年, 毒月饼; 2793 年, 毒月饼; 2794 年, 毒月饼; 2795 年, 毒月饼; 2796 年, 毒月饼; 2797 年, 毒月饼; 2798 年, 毒月饼; 2799 年, 毒月饼; 2800 年, 毒月饼; 2801 年, 毒月饼; 2802 年, 毒月饼; 2803 年, 毒月饼; 2804 年, 毒月饼; 2805 年, 毒月饼; 2806 年, 毒月饼; 2807 年, 毒月饼; 2808 年, 毒月饼; 2809 年, 毒月饼; 2810 年, 毒月饼; 2811 年, 毒月饼; 2812 年, 毒月饼; 2813 年, 毒月饼; 2814 年, 毒月饼; 2815 年, 毒月饼; 2816 年, 毒月饼; 2817 年, 毒月饼; 2818 年, 毒月饼; 2819 年, 毒月饼; 2820 年, 毒月饼; 2821 年, 毒月饼; 2822 年, 毒月饼; 2823 年, 毒月饼; 2824 年, 毒月饼; 2825 年, 毒月饼; 2826 年, 毒月饼; 2827 年, 毒月饼; 2828 年, 毒月饼; 2829 年, 毒月饼; 2830 年, 毒月饼; 2831 年, 毒月饼; 2832 年, 毒月饼; 2833 年, 毒月饼; 2834 年, 毒月饼; 2835 年, 毒月饼; 2836 年, 毒月饼; 2837 年, 毒月饼; 2838 年, 毒月饼; 2839 年, 毒月饼; 2840 年, 毒月饼; 2841 年, 毒月饼; 2842 年, 毒月饼; 2843 年, 毒月饼; 2844 年, 毒月饼; 2845 年, 毒月饼; 2846 年, 毒月饼; 2847 年, 毒月饼; 2848 年, 毒月饼; 2849 年, 毒月饼; 2850 年, 毒月饼; 2851 年, 毒月饼; 2852 年, 毒月饼; 2853 年, 毒月饼; 2854 年

[-6.88628220e-05 -1.19989728e-03 -2.36077557e-04 2.49695818e-03
-8.79188088e-04 -1.48803692e-03 1.81370428e-03 -1.79939355e-03....
-5.38808817e-04 -4.01521463e-03 2.63561918e-03 -1.31672230e-03
1.45825999e-04 6.91009550e-04 2.14562120e-04 -3.06116202e-03]

其次,根据“肘部法则”来确定最优“簇值”K,即当图形变化很快,出现类似“肘部”时,及K较为适合,具体实验过程中“肘部法则”的最优化选择过程如图2所示。

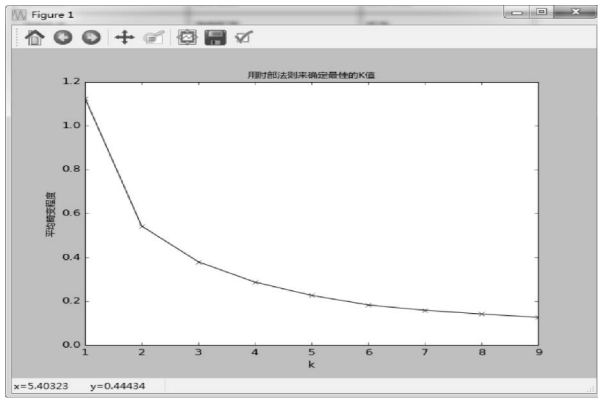


图2 “肘部法则”来确定最优“簇值”K过程图

最后,调用 sklearn 包来进行聚类,并得到“轮廓系数”(轮廓系数的值是介于 [-1,1],越趋近于 1 代表内聚度和分离度都相对较优)。

在上述整个聚类的流程下,对于“甲醛白菜”微博文本当中所蕴含的类别知识进行了预估计和逐步的实验,并观察在给定的类下轮廓系数的不同。具体实验结果如表2所示。其中超参结果的选择为 $\alpha=0.001$, $\text{min_count}=40$, $\text{size}=200$,最佳聚类类别 $k=2$ 。

表2 基于 word2vec 的聚类结果与分析

相关微博发表年份	对应轮廓系数值
2012	0.6607
2013	0.7153
2014	0.8808

根据三个不同的轮廓系数和图3所示的聚类结果示意图可以看出,利用 word2vec 与 K-means 算法进行相融合的策略,能有效对真实的与甲醛白菜相关的所有微博文本中蕴含的类别知识进行聚类。

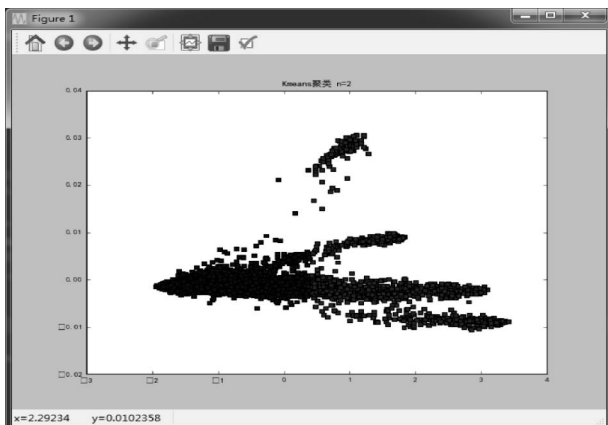


图3 $k=2$ 时 2012 年微博文本聚类结果图

3 结语

为了从食品安全事件微博文本中获取相应的类别知识,本文在新的文本向量表示下,结合已有的相应聚类算法,针对甲醛白菜中所蕴含的类别知识进行了系统的挖掘和分析。通过具体实验的验证,本文发现所聚出来的类为2的时候,可以看作最有效的聚类结果,并符合甲醛白菜微博文本中所包含的类别知识。

参考文献:

- [1] 张星联.我国食品安全预警数据库系统的建设与实现[J].食品科技,2008,33(12):250-254.
- [2] 吴云红.食品监管改革的关键--基于互联网的动态第三方数据库[J].食品工业科技,2009(09):272-274.
- [3] 余清,洪源.加工食品风险数据库的构建思路[J].价值工程,2013(30):174-175.
- [4] 贾凯,彭培好,阮伟玲.四川省彭州市三界镇农民专业合作社调查研究[J].北京农业,2014,1:247-248.
- [5] 刘翠玲,徐莹莹,孙晓荣,等.基于多源大数据食品安全监测预控系统的设计与实现[J].食品科学技术学报,2018,36(3):88-94.
- [6] 夏长辉.一种改进的 K-means 聚类算法[J].信息与电脑(理论版),2017,(14):40-42.
- [7] 王侠林,贺建峰.基于 K-Means 聚类的微生物群落结构研究[J].软件导刊,2018(1):146-148.
- [8] 黄建宇,周爱武,肖云,谭天诚.基于特征空间的文本聚类[J/OL].计算机技术与发展,2017,(08):1-3.
- [9] 赵卫中,马慧芳,傅翔翔,等.基于云计算平台 Hadoop 的并行 k-means 聚类算法设计研究[J].计算机科学,2011,38(10):166-168.
- [10] 武森,冯小东,杨杰,等.基于 MapReduce 的大规模文本聚类并行化[J].北京科技大学学报,2014(10):1411-1419.
- [11] Mikolov T, Chen K, Corrado G, et al. Efficient Estimation of Word Representations in Vector Space[J]. Computer Science, 2013.
- [12] Mikolov T, Sutskever I, Chen K, et al. Distributed Representations of Words and Phrases and their Compositionality[J]. Advances in Neural Information Processing Systems, 2013, 26:3111-3119.
- [13] 姜霖,王东波.采用连续词袋模型(CBOW)的领域术语自动抽取研究[J].现代图书情报技术,2016,32(2):9-15.
- [14] eh ek R. models.word2vec-Deep learning with word2vec [EB/OL]. [2018-07-26]. <https://radimrehurek.com/gensim/models/word2vec.html>.
- [15] 张冬梅.基于轮廓系数的层次聚类算法研究[D].秦皇岛:燕山大学,2009.
- [16] 朱连江,马炳先,赵学泉.基于轮廓系数的聚类有效性分析[J].计算机应用,2010(s2):139-141.

基金项目 本文系江苏省社会科学基金“时间感知大数据特征下的食品安全突发事件应对策略挖掘研究”(项目编号:15TQC003)的研究成果之一。

作者简介 沈思(1983-),女,讲师,博士,主要研究方向为信息检索,学术文本挖掘,时间语义分析;梁晓静(1996-),女,主要研究方向为学术文本检索。