

VSM-LSI 分析

1. 工作概要

第二次作业是对 VSM 和 LSI 的简单使用，并体会它们的用处和不同，感谢前面几位同学的作业，从而了解到可以使用 gensim 这个开源库，并在此基础上进行了 VSM-LSI 的相似度计算和比较

2. 语料预处理

首先，由于老师给的文档不是连续的，如下图，一篇文档分成了若干部分，并以空格隔开，所以首要任务就是把原来是一个文档的数据放到一起。

```
19980101-01-001-001/m 迈向/v 充满/v 希望/n 的/u 新/a 世纪/n ——/w 一九九八年/t 新年/t 讲话/n (/w 附/v 图片/n 1/r ^
19980101-01-001-002/m 中共中央/nt 总书记/n 、/w 国家/n 主席/n 江/nr 泽民/nr
19980101-01-001-003/m (/w 一九九七年/t 十二月/t 三十一日/t ) /w
19980101-01-001-004/m 1 2月/t 3 1日/t , /w 中共中央/nt 总书记/n 、/w 国家/n 主席/n 江/nr 泽民/nr 发表/v 1 9 9 8年
19980101-01-001-005/m 同胞/n 们/k 、/w 朋友/n 们/k 、/w 女士/n 们/k 、/w 先生/n 们/k : /w
19980101-01-001-006/m 在/p 1 9 9 8年/t 来临/v 之际/f , /w 我/r 十分/m 高兴/a 地/u 通过/p 中央/n 人民/n 广播/vn 电
19980101-01-001-007/m 1 9 9 7年/t , /w 是/v 中国/ns 发展/vn 历史/n 上/f 非常/d 重要/a 的/u 很/d 不/d 平凡/a 的/u -
19980101-01-001-008/m 在/p 这/r 一/m 年/q 中/f , /w 中国/ns 的/u 改革/vn 开放/vn 和/c 现代化/vn 建设/vn 继续/v 向
19980101-01-001-009/m 在/p 这/r 一/m 年/q 中/f , /w 中国/ns 的/u 外交/n 工作/vn 取得/v 了/u 重要/a 成果/n 。/w 通
19980101-01-001-010/m 1 9 9 8年/t , /w 中国/ns 人民/n 将/d 满怀信心/l 地/u 开创/v 新/a 的/u 业绩/n 。/w 尽管/c 我
19980101-01-001-011/m 实现/v 祖国/n 的/u 完全/a 统一/vn , /w 是/v 海内外/s 全体/n 中国/ns 人/n 的/u 共同/b 心愿/n 。
19980101-01-001-012/m 台湾/ns 是/v 中国/ns 领土/n 不可分割/l 的/u 一/m 部分/n 。/w 完成/v 祖国/n 统一/vn , /w 是/v
19980101-01-001-013/m 环顾/v 全球/n , /w 日益/d 密切/a 的/u 世界/n 经济/n 联系/vn , /w 日新月异/i 的/u 科技/n 进步,
19980101-01-001-014/m 中国/ns 政府/n 将/d 继续/v 坚持/v 奉行/v 独立自主/i 的/u 和平/n 外交/n 政策/n , /w 在/p 和平
19980101-01-001-015/m 在/p 这/r 辞旧迎新/l 的/u 美好/a 时刻/n , /w 我/r 祝/v 大家/r 新年/t 快乐/a , /w 家庭/n 幸福/a
19980101-01-001-016/m 谢谢/v ! /w (/w 新华社/nt 北京/ns 1 2月/t 3 1日/t 电/n ) /w

19980101-01-002-001/m 在/p 十五大/j 精神/n 指引/vn 下/f 胜利/vd 前进/v ——/w 元旦/t 献辞/n
19980101-01-002-002/m 我们/r 即将/d 以/p 丰收/vn 的/u 喜悦/an 送/v 走/v 牛年/t , /w 以/p 昂扬/a 的/u 斗志/n 迎来/v
19980101-01-002-003/m 刚刚/d 过去/v 的/u 一/m 年/q , /w 大气磅礴/i , /w 波澜壮阔/i 。/w 在/p 这/r 一/m 年/q , /w
和/c 国际/n 威望/n 进一步/d 提高/v 。/w 实践/v 使/v 亿万/m 人民/n 对/p 邓小平理论/n 更加/d 信仰/v , /w 对/p 以/p 江
19980101-01-002-004/m 1 9 9 8年/t , /w 是/v 全面/ad 贯彻/v 落实/v 党/n 的/u 十五大/j 提出/v 的/u 任务/n 的/u 第一/i
```

即如第一个同学分享的 txt:

```
NO: 0 Title:迈向充满希望的新世纪——一九九八年新年讲话 (附图片 1 张) /Author:中共中央总书记、国家主席江泽民/main Body: (一
共同心愿。通过中葡双方的合作和努力,按照“一国两制”方针和澳门《基本法》, 1 9 9 9 年 1 2 月澳门的回归一定能够顺利实现。台湾
NO: 1 Title:在十五大精神指引下胜利前进——元旦献辞/Author:我们即将以丰收的喜悦送走牛年,以昂扬的斗志迎来虎年。我们伟大祖国
坚持这个中心,是完全正确的。今后,我们能否把建设有中国特色社会主义伟大事业全面推向 2 1 世纪,关键仍然要看能否把经济工作搞上
NO: 2 Title:北京举行新年音乐会/Author:江泽民李鹏乔石朱鎔基李瑞环刘华清尉健行李岚清与万名首都各界群众和劳动模范代表一起辞旧
留念。李铁映、贾庆林、曾庆红等领导同志也出席了今晚音乐会。
NO: 3 Title:李鹏在北京考察企业/Author:向广大职工祝贺新年,对节日坚守岗位的同志们表示慰问/main Body:新华社北京十二月三十一
,当李鹏了解到老闫退休前一直都是厂里的先进工作者、曾经被评为北京市“五好职工”,退休后仍然为改善职工的住房而奔波时,十分高
NO: 4 Title:挂起红灯迎新年 (图片) /Author:元旦来临,安徽省合肥市长江路悬挂起 3 3 0 0 盏大红灯笼,为节日营造出“千盏灯笼凌空
NO: 5 Title:全总致全国各族职工慰问信/Author:勉励广大职工发挥工人阶级主力军作用,为企业改革发展建功立业/main Body:本报北京
NO: 6 Title:忠诚的共产主义战士,久经考验的无产阶级革命家刘澜涛同志逝世/Author: (附图片 1 张) /main Body:新华社北京 1 2 月 3
NO: 7 Title:党中央国务院关心西藏雪灾救灾工作/Author:灾区各级政府全力组织抗灾力争降低灾害损失/main Body:据新华社北京 1 2 月
NO: 8 Title:明天天气预报 (1 9 9 8 年元月 1 日 2 0 时—元月 2 日 2 0 时) /Author:天气趋势分析/main Body:受暖气流影响,今天晚
```

不过这里没必要弄成如上的样子,不然还得划分,上面的可以用来对结果比对时分析,下面的代码将相同文章的词进行了合并,得到词袋:

```
"""
```

函数说明：语料预处理

输入：要打开的文件名

输出：

documents - 处理过后的词空间列表

```
"""
```

```
def load_data(filename = '199801_clear.txt'):  
    # 文档分词列表  
    documents = []  
    #文档分词  
    doc = []  
    # 打开语料库文件，逐行读取语料，将同一篇文章的不同章节合并，  
    并过滤掉无意义的词和符号  
    fr_open = open(filename, 'r', encoding='GBK')  
  
    for line in fr_open:  
        if line.strip(): # 跳过空行, 这里 strip () 无参数是去掉  
            首尾空格  
            words = line.split()  
            # 过滤掉无意义的词和符号  
            for i in range(len(words)):  
                word = words[i]  
                if 'w' not in word and 'y' not in word and 'u'  
not in word \   
                    and 'c' not in word and 'm' not in word:  
doc.append(word)  
            else: #出现空行，即将开始下一篇文章读入，所以把上一篇文章  
            档添加进 documents，并将 doc 置为初始  
                documents.append(doc)  
                doc = []  
    return documents
```

3. TF-IDF 权重

再对文本进行重组后就可以统计它的词频了，并根据词袋计算它们的 tfidf 值，这里使

用 gensim 的 corpora 进行词频统计和计算权重，代码注释已经很详细了

```
"""
函数说明：统计并计算 tf_idf 权重值
输入：documents - 处理过后的词空间列表
输出：corpus_tfidf - tfidf 的词向量空间
"""

def tf_idf(documents):
    # 将 documents 放入一个字典, 这个字典默认按照一定的顺序排列
    documents

    dictionary = corpora.Dictionary(documents)
    #转换成词频, 会去掉每个文本的重复词, 并自动生成词向量空间
    corpus = [dictionary.doc2bow(text) for text in documents]
    # 将 corpus 作为初始化的参数, 得到 tfidf 关于 corpus 的模型
    tfidf = models.TfidfModel(corpus)
    # 通过 tfidf 【】 将词频转换成 tfidf 值
    corpus_tfidf = tfidf[corpus]
    return corpus_tfidf
```

4. VSM 及结果分析

将上一步得到的词向量空间计算它们的相似度，这里不能直接用 similarities.MatrixSimilarity () 函数，会出现内存不够，应该先将其持久化到磁盘，再计算相似度，相当于扩大了虚拟内存。然后再将相似度输出到 csv 文件里面，因为相似度矩阵是一个对称阵，所以我输出了一个下三角型矩阵（节约存储空间），最后可以查找某一个文本和其他文本的相似度，我这里返回前 3 个最大相似度的文本的序号

```
"""
函数说明：进行 LSI 分析，并计算各个文档之间的相似度
输入：corpus_tfidf - tfidf 的词向量空间
输出：CSV 格式的 VSM 下三角相似度矩阵
"""
```

```
def VSM(corpus_tfidf):
    # 必须先把 corpus_tfidf 写入磁盘，因为矩阵太大，直接在内存计算相似度会导致内存不够
    corpora.MmCorpus.serialize('./vsm_test.mm', corpus_tfidf)
    corpus_tfidf = corpora.MmCorpus('./vsm_test.mm')
    # 用文档向量初始化一个相似度计算的对象
    # index = similarities.MatrixSimilarity(corpus_tfidf)
    index = similarities.SparseMatrixSimilarity(corpus_tfidf)
    # 输出下三角矩阵
    with open('./vsm.csv', 'w') as fr:
        i = 1
        for sims in index[corpus_tfidf]:
            a = sims[:i]
            fr.write(','.join(map(str, a)) + '\n')
            i += 1
        num = 1 # 查找第 num 篇文档的相似度，可自行修改
        query = sorted(enumerate(index[corpus_tfidf[num]]),
            key=lambda x: x[1], reverse=True)
        # 返回最相似的前三个文档的序号
        print(query[:4])
```

输出的下三角形矩阵：

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	1																						
2	0.277739	0.999999																					
3	0.00499	0.071296	1																				
4	0.080266	0.06591	0.10779	1																			
5	0.018304	0.005512	0.028968	0.016821	1																		
6	0.116771	0.244071	0.037285	0.077703	0.008671	1																	
7	0.048026	0.025518	0.036603	0.022286	0.003329	0.008943	1																
8	0.039736	0.048319	0.014968	0.022193	0.001602	0.018292	0.020783	1															
9	0.014384	0.003374	0.008322	0.009667	0.003259	0	0.06495	0															
10	0.185868	0.118165	0.04025	0.060419	0.014509	0.067145	0	0.003361	1														
11	0.091891	0.123377	0.05202	0.074131	0.004515	0.079895	0.043223	0.064256	0.023401	0.039557	1												
12	0.053671	0.108765	0.048245	0.038034	0.003325	0.07446	0.037272	0.016788	0	0.06315	0.058373	1											
13	0.029435	0.029472	0.012391	0.015865	0.008921	0.023889	0.002204	0.008871	0.008022	0.010102	0.022351	0.004905	1										
14	0	0	0	0	0	0	0	0	0	0	0	0	0	1									
15	0.064122	0.029058	0.039773	0.016065	0.011615	0.016541	0.030723	0.00836	0.00653	0.031876	0.029775	0.020499	0.022184	0	1								
16	0.048339	0.008599	0.02311	0.019504	0.037749	0.005564	0.006336	0.010371	0	0.01245	0.011634	0.005211	0.00315	0	0.004545	1							
17	0.301804	0.130087	0.018956	0.024654	0.002125	0.056275	0.022621	0.035984	0.018448	0.060392	0.053885	0.04835	0.012635	0	0.021649	0.012885	1						
18	0.026565	0.008473	0.002489	0.020707	0	0.011709	0.000552	0.002807	0	0.001834	0.007359	0.00151	0.000858	0	0.022216	0	0.001892	1					
19	0.018215	0.014346	0.001077	0.011362	0.001245	0.006104	0.000406	0.002475	0.000775	0.002893	0.026385	0.004918	0.008057	0	0.006844	0.000958	0.00276	0.006454	1				
20	0.014814	0.00605	0.01115	0.012352	0.000497	0.003215	0.000493	0.002036	0	0.009929	0.001478	0.00049	0.002178	0	0.015627	0.000383	0.00049	0.118694	0.085791	1			
21	0.048822	0.020448	0.011186	0.015519	0.003562	0.017899	0.004699	0.006854	0.001744	0.012796	0.017837	0.010715	0.00806	0	0.008117	0.006066	0.027957	0.005073	0.004202	0.014411	1		
22	0.019537	0.003607	0.011824	0.016546	0.000372	0.004009	0.040417	0.007896	0.022065	0.007057	0.016229	0.010174	0.002346	0	0.013706	0.000745	0.019701	0.002096	0.000512	0.002224	0.012966	1	
23	0.187144	0.070361	0.029429	0.035114	0.003722	0.013593	0.025666	0.012912	0.051395	0.047132	0.024738	0.015976	0	0.038821	0.007486	0.124033	0.004244	0.006341	0.002358	0.002668	0.007988	0	1
24	0.168836	0.047325	0.020213	0.030382	0.002063	0.023193	0.018690	0.025111	0.010356	0.040487	0.048715	0.023212	0.013087	0	0.024614	0.014408	0.127324	0.001604	0.007773	0.001927	0.02096	0.021593	0.779778
25	0.152706	0.085385	0.019639	0.039826	0.002529	0.021246	0.023163	0.04131	0.018098	0.043003	0.053004	0.037375	0.013833	0	0.029056	0.005004	0.182528	0.005197	0.005131	0.019977	0.037671	0.033218	0.129007
26	0.060409	0.053279	0.023315	0.035791	0.001704	0.022256	0.023882	0.006978	0.008441	0.030464	0.020699	0.008793	0.009845	0	0.027248	0.022926	0.025228	0.002585	0.003564	0.002918	0.01211	0.006967	0.021986
27	0.153394	0.050961	0.022062	0.027419	0.005552	0.027711	0.022135	0.026024	0.011649	0.033327	0.068418	0.000148	0.021113	0	0.044892	0.034685	0.065751	0.016773	0.042582	0.003477	0.034849	0.013577	0.31831
28	0.154317	0.073479	0.021759	0.026395	0.000709	0.027121	0.041228	0.026645	0.002633	0.044852	0.036719	0.034123	0.006026	0	0.030813	0.007211	0.0818	0.000482	0.003941	0.000329	0.013821	0.016247	0.108442
29	0.025593	0.009349	0.028048	0.007874	0.030881	0.009904	0.012965	0.000888	0	0.019188	0.002602	0.001677	0.005435	0	0.011611	0.039194	0.001899	0.000101	0.000428	0.000409	0.020209	0.001505	0.002194
30	0.106086	0.091794	0.018717	0.020767	0.001793	0.071419	0.020989	0.027104	0.012432	0.023387	0.080315	0.024477	0.018551	0	0.018217	0.008235	0.065467	0.000852	0.027116	0.002954	0.021835	0.014007	0.002815

可以再 excel 里面查找相似度较高的文章，这里我搜寻了大于 0.8 的值，如 101 行 BH 列的值为 0.823，即第 101 篇和第 60 篇的相似度为 0.823，它们是：

NO: 59 Title:居民生活如常事业蒸蒸日上香港特区政府发表回归半年大事录/Author:据新华社香港 1 月 1 日电香港特区政府在 1 9 9 7 年的最后一天，发表了《香港特别行政区——首六个月大事录》，认为香港回归后的情况“粉碎了一些对 7 月 1 日后香港变得灰暗的预言”。/main Body:大事录详细记录了香港回归以来政治、经济、法治、人权、教育、基本建设等方面的重大事项。包括：香港成为中国

的一个特别行政区后，在“一国两制”方针政策下运作；基本法保障了香港资本主义制度及原有生活方式五十年不变；发行自己的货币；保持原有税制、独立关税区地位和公务员制度等。这份特区政府的报告长达20页。文中说，香港顺利回归祖国，香港居民生活如常，正如基本法所保证的一样。报告说，回归两个月后，香港成功举行了世界银行和国际货币基金组织理事会年会，数千名来自各国和地区的财长、央行行长及官员和2000多传媒人士，在香港得到了第一手的特区运作观感，香港回归后的情况赢得了人们的信心。报告认为，虽然亚洲出现了金融风波，但香港经济仍然跨步向前。报告重申了特区政府维持联系汇率制度的决心，指出香港拥有多种有利因素，包括：比邻近地区通胀率低、香港银行业采取审慎贷款及维持高资本充足比率的政策、租金和楼价下调、联系汇率可稳定汇率等。因此，“假若情况好转，本港的复原速度亦会较快”。报告还记录了行政长官董建华去年9至11月的外访活动和10月发表首份施政报告。

NO: 100 Title:谱写脱贫致富新篇章——吕梁扶贫开发工作与计划生育相结合工作纪实/Author:赖林嵩周美林范双民张振明/main Body:山西省吕梁地区是个贫困大区。全区13个县市中有9个是贫困县，其中国家重点扶持的贫困县有7个。1971年建区时，总人口226.6万，出生率高达33.10‰，死亡率8.4‰，人口自然增长率24.70‰。1976年开始推行计划生育，无计划生育势头有所遏制。近几年来，全区上下结合实施“八七扶贫攻坚计划”，按“经济要上去，人口要下来”的方针，计划生育工作连年取得好成绩。1996年，全区总人口控制在324万人，出生率下降到14.55‰，自然增长率9.43‰。贫困人口由1992年的105万人减少到1996年的80.4万人。吕梁地区的计划生育工作坚持服从、服务于经济建设这个中心，服从、服务于扶贫攻坚这个历史性任务，努力实现计划生育工作的“两个转变”，走出了一条扶贫开发工作与计划生育相结合的希望之路、成功之路。党政重视，步入人口发展的良性循环轨道吕梁地区坚持不懈、持之以恒地严格控制人口增长。他们采取的主要措施有：——明确一把手负总责。11年来，一直把计划生育工作作为“一把手工程”来抓，地委、行署每年召开计划生育工作专题会议都在三次以上，每次会议都要切实解决几个实质性问题。每年还要定期召开一把手“交账会议”，既交经济账，又交计划生育账。——加强基层基础工作不松动。全区各县市都建立了计划生育服务站；236个乡镇全部建立了计划生育服务所；4626个村配齐了计划生育服务员，人员、经费、待遇得到落实。——确保经费投入逐年增加。地区计划生育经费已由1990年前每年20万元逐年增加到1997年的300万元。县市财政用于计划生育的经费由1990年的117.6万元（人均0.60元），增长到1996年的590万元（人均1.85元）。穷则思变，走扶贫开发与计生相结合之路吕梁的特点是贫困人口占全省的1/3左右。吕梁的干部群众认识到，扶贫开发工作与计划生育相结合，是解决贫困地区人口与发展问题的必由之路。1990年，地委、行署在总结经验的基础上正式出台《关于扶贫开发工作与计划生育相结合的意见》，在全区实施。结合学习外地经验，吕梁地区扶贫开发工作与计划生育相结合工作开始形成党政牵头、部门配合、计生协调、群众参与的格局。柳林县下三交镇坪上村原是个有名的老大难村，党支部、村委会班子瘫痪，多年来人均收入徘徊在一二百元之间。1986年全村242名已婚妇女，有半数多超计划生育。这一年新党支部书记刘步飞上任，他带头只生了一个女孩，拿出自己放电影补助费捐赠给计划生育贫困户，垫支7000元兴修水利、引水浇地，还带领群众创办了“少生快富甜蜜公司”

看标题 NO: 59Title:居民生活如常事业蒸蒸日上香港特区政府发表回归半年大事录和 NO: 100 Title:谱写脱贫致富新篇章，知两个都是与人民生活越来越好相关。还是比较相关。

5. LSI 及结果分析

前面的数据处理过程都一样，就是得到 tfidf 权重值。LSI 主要是对矩阵用 SVD 进行了一个降维，根据 SVD，任何 $m \times n$ 的矩阵 M ，都存在以下 SVD 分解：

$$M_{m \times n} = U_{m \times m} \Sigma_{m \times n} (V_{n \times n})^T$$

通过几何意义对其证明如下：

假设存在 $M \times N$ 矩阵 A ，事实上， A 矩阵将 n 维空间中的向量映射到 k ($k \leq m$) 维空间中， $k = \text{Rank}(A)$ 。现在的目标就是：在 n 维空间中找一组正交基，使得经过 A 变换后还是正交的。假设已经找到这样一组正交基：

$$\{v_1, v_2, \dots, v_n\}$$

则 A 矩阵将这组基映射为：

$$\{Av_1, Av_2, \dots, Av_n\}$$

如果要使他们两两正交，即

$$Av_i \cdot Av_j = (Av_i)^T Av_j = v_i^T A^T A v_j = 0$$

根据假设，存在

$$v_i^T v_j = v_i \cdot v_j = 0$$

所以如果正交基 v 选择为 $A^T A$ 的特征向量的话，由于 $A^T A$ 是对称阵， v 之间两两正交，那么

$$\begin{aligned} v_i^T A^T A v_j &= v_i^T \lambda_j v_j \\ &= \lambda_j v_i^T v_j \\ &= \lambda_j v_i \cdot v_j = 0 \end{aligned}$$

这样就找到了正交基使其映射后还是正交基了，现在，将映射后的正交基单位化：

因为

$$Av_i \cdot Av_i = \lambda_i v_i \cdot v_i = \lambda_i$$

所以有

$$|Av_i|^2 = \lambda_i \geq 0$$

所以取单位向量

$$u_i = \frac{Av_i}{|Av_i|} = \frac{1}{\sqrt{\lambda_i}} Av_i$$

由此可得

$$Av_i = \sigma_i u_i, \sigma_i(\text{奇异值}) = \sqrt{\lambda_i}, 0 \leq i \leq k, k = \text{Rank}(A)$$

当 $k < i \leq m$ 时, 对 u_1, u_2, \dots, u_k 进行扩展 $u_{(k+1)}, \dots, u_m$, 使得 u_1, u_2, \dots, u_m 为 m 维空间中的一组正交基, 即

将 $\{u_1, u_2, \dots, u_k\}$ 正交基扩展成 $\{u_1, u_2, \dots, u_m\}$ R^m 空间的单位正交基

同样的, 对 v_1, v_2, \dots, v_k 进行扩展 $v_{(k+1)}, \dots, v_n$ (这 $n-k$ 个向量存在于 A 的零空间中, 即 $Ax=0$ 的解空间的基), 使得 v_1, v_2, \dots, v_n 为 n 维空间中的一组正交基, 即

在 A 的零空间中选取 $\{v_{k+1}, v_2, \dots, v_n\}$ 使得 $Av_i=0, i>k$, 并取 $\sigma_i=0$

则可得到

$$A[v_1 \ v_2 \ \dots \ v_k \ | \ v_{k+1} \ \dots \ v_n] = [u_1 \ u_2 \ \dots \ u_k \ | \ u_{k+1} \ \dots \ u_m] \begin{bmatrix} \sigma_1 & & & & 0 \\ & \ddots & & & \\ & & \sigma_k & & 0 \\ \hline & & & 0 & 0 \end{bmatrix}$$

继而可以得到 A 矩阵的奇异值分解:

$$A = U \Sigma V^T$$

Σ 是一个对角矩阵, 其中的元素值就是奇异值, 并且按照从大到小的顺序排列。在很多情况下, 前 10% 甚至更少的奇异值的平方和就占了全部奇异值平方和的 90% 以上了, 因此可以用前 k 个奇异值来近似描述矩阵:

$$M_{m \times n} \approx U_{m \times k} \Sigma_{k \times k} (V_{n \times k})^T$$

使用 `models.LsiModel()` 函数构造 LSI 模型并计算相似度:

"""

函数说明: 进行 LSI 分析, 并计算各个文档之间的相似度

输入: `corpus_tfidf` - tfidf 的词向量空间

输出: CSV 格式的 LSI 下三角相似度矩阵

"""

```
def LSI(corpus_tfidf):
    # 构造 LSI 模型并将待检索的 query 和文本转化为 LSI 主题向量
    # 转换之前的 corpus 是 BOW 向量
    lsi_model = models.LsiModel(corpus_tfidf, num_topics=80)
    corpus_lsi = lsi_model[corpus_tfidf]
    # 转换成潜语义文档向量列表并持久化
    corpora.MmCorpus.serialize('./lsi_test.mm', corpus_lsi)
    corpus_lsi = corpora.MmCorpus('./lsi_test.mm')
    # 构建文档相似度矩阵索引用于查询，再使用文档列表本身进行相似度查询（默认使用 Cosine）
    index = similarities.MatrixSimilarity(corpus_lsi)
    # 输出下三角矩阵
    with open('./lsi.csv', 'w') as fr:
        i = 1
        for sims in index[corpus_lsi]:
            a = sims[:i]
            fr.write(','.join(map(str, a)) + '\n')
            i += 1
```

在降维后，再来计算它们的相似度，可以得到：

#	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W
1	1																						
2	0.59993	1																					
3	0.244101	0.200858	1																				
4	0.25338	0.286082	0.48346	1																			
5	0.10598	0.05562	0.292494	0.216225	1																		
6	0.487305	0.789222	0.167468	0.371148	0.118048	1																	
7	0.127699	0.276457	0.31954	0.219979	-0.00435	0.185731	1																
8	0.108629	0.11913	0.074304	0.108496	0.038474	0.118144	0.080063	1															
9	0.038622	0.006729	0.018288	0.055862	0.016541	#####	0.023828	0.189162	1														
10	0.818132	0.607392	0.366317	0.353969	0.150242	0.532899	0.249853	0.075568	0.040852	1													
11	0.252267	0.398427	0.248441	0.273492	0.014779	0.346781	0.254633	0.281582	0.036092	0.263701	1												
12	0.23499	0.476246	0.144837	0.219444	0.037895	0.449748	0.287494	0.136221	-0.00368	0.310543	0.275881	1											
13	0.258643	0.164654	0.110537	0.290746	0.205284	0.322374	0.048364	0.147869	0.107663	0.225775	0.222418	0.067254	1										
14	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
15	0.316017	0.195846	0.247539	0.043868	0.115573	0.126214	0.309341	0.063058	0.055635	0.344455	0.093726	0.156034	0.297339	0	1								
16	0.250414	0.151144	0.202152	0.08781	0.575775	0.158556	0.125995	0.100736	0.032341	0.203025	0.088104	0.112191	0.167322	0	0.149822	1							
17	0.747127	0.510313	0.081805	0.160444	0.03887	0.143072	0.073534	0.135474	0.04501	0.533738	0.216556	0.227672	0.248074	0	0.174942	0.161321	1						
18	0.029818	0.068289	-0.04563	0.11153	0.031912	0.057511	0.034815	0.04901	0.046655	0.082167	0.033254	0.086306	0.166114	0	0.174469	-0.03723	0.053728	1					
19	0.098	0.163644	-0.00603	0.139368	0.026269	0.066358	0.013645	0.120161	0.02931	0.119964	0.247737	0.098957	0.349932	0	0.192743	0.042223	0.147238	0.501349	1				
20	0.044126	0.022554	0.020096	0.193638	0.102558	-0.01229	0.067859	0.142477	0.064761	0.116427	-0.03485	0.050241	0.212289	0	0.30455	0.043838	0.092013	0.649557	0.527186	1			
21	0.315443	0.156202	-0.0183	0.040356	-0.00563	0.097677	0.062639	0.066559	0.014951	0.135166	0.112976	0.028525	0.188909	0	0.186371	0.114471	0.284791	0.048959	0.14008	0.097512	1		
22	0.077932	0.00963	0.028291	0.061393	0.003712	0.015923	0.153581	0.035544	0.05112	0.04062	0.068328	0.083454	0.122169	0	0.197599	0.116574	0.172059	0.037547	0.096985	0.044351	0.026407	1	
23	0.453283	0.159162	0.071168	0.10723	0.032096	0.148075	0.100482	0.073804	0.043108	0.340351	0.122108	0.095963	0.100208	0	0.159754	0.15496	0.44246	0.015725	0.064367	0.020693	0.288831	0.048807	1
24	0.423061	0.096019	0.037695	0.08482	0.020856	0.09183	0.091766	0.062463	0.035003	0.308859	0.090737	0.080267	0.063622	0	0.097613	0.10676	0.461354	0.021302	0.064546	0.007065	0.241252	0.061106	0.953239
25	0.40937	0.236612	0.044899	0.127417	0.052715	0.192212	0.15484	0.072803	0.058908	0.238417	0.159294	0.110563	0.138042	0	0.138379	0.098502	0.631936	0.063102	0.104338	0.201112	0.313911	0.164765	0.252472
26	0.403748	0.297286	0.202417	0.27873	-0.03044	0.13814	0.319922	-0.00096	0.052901	0.344505	0.162254	0.049363	0.241446	0	0.257947	0.116643	0.306101	0.012826	0.087169	0.077247	0.283412	0.037566	0.206603
27	0.457602	0.178414	0.101627	0.147808	0.06942	0.145085	0.042946	0.070983	0.065997	0.321669	0.226172	0.064403	0.200955	0	0.268021	0.179636	0.440336	0.009868	0.104403	0.027022	0.270182	0.064997	0.84256
28	0.644483	0.324242	0.225097	0.189792	0.017846	0.285749	0.280663	0.111165	0.02575	0.456459	0.185995	0.275129	0.059325	0	0.246545	0.072309	0.672603	-0.00555	0.011885	-0.0449	0.286284	0.145223	0.475537
29	0.084692	0.077896	0.278317	0.017643	0.626855	0.057067	0.185513	0.113038	0.014891	0.13724	0.040373	0.047848	0.162744	0	0.174302	0.56164	0.048482	0.051481	0.145845	0.181255	0.314374	0.032689	0.052629
30	0.479914	0.498792	0.127292	0.231914	0.052849	0.417657	0.104504	0.109107	0.008015	0.298264	0.508265	0.161839	0.207282	0	0.210249	0.065519	0.466916	0.003721	0.163268	0.031662	0.29623	0.183418	0.206833

下面是从中选取的两篇相似度较高的文档的分析：

NO: 0 Title:迈向充满希望的新世纪——一九九八年新年讲话（附图片 1 张） /Author:中共中央总书记、国家主席江泽民/main Body:

（一九九七年十二月三十一日） 1 2 月 3 1 日，中共中央总书记、国家主席江泽民发表 1 9 9 8 年新年讲话《迈向充满希望的新世纪》。（新华社记者兰红光摄） 同胞们、朋友们、女士们、先生们：在 1 9 9 8 年来临之际，我十分高兴地通过中央人民广播电台、中国国际广播电台和中央电视台，向全国各族人民，向香港特别行政区同胞、澳门和台湾同胞、海外侨胞，向世界各国的朋友们，致以诚挚的问候和良好的祝愿！ 1 9 9 7 年，是中国发展历史上非常重要的很不平凡的一年。中国人民决心继承邓小平同志的遗志，继续把建设有中国特色社会主义事业推向前进。中国政府顺利恢复对香港行使主权，并按照“一国两制”、“港人治港”、高度自治的方针保持香港的繁荣稳定。中国共产党成功地召开了第十五次全国代表大会，高举邓小平理论伟大旗帜，总结百年历史，展望新的世纪，制定了中国跨世纪发展的行动纲领。在这一年中，中国的改革开放和现代化建设继续向前迈进。国民经济保持了“高增长、低通胀”的良好发展态势。农业生产再次获得好的收成，企业改革继续深化，人民生活进一步改善。对外经济技术合作与交流不断扩大。民主法制建设、精神文明

建设和其他各项事业都有新的进展。我们十分关注最近一个时期一些国家和地区发生的金融风波，我们相信通过这些国家和地区的努力以及有关的国际合作，情况会逐步得到缓解。总的来说，中国改革和发展的全局继续保持了稳定。在这一年中，中国的外交工作取得了重要成果。通过高层互访，中国与美国、俄罗斯、法国、日本等大国确定了双方关系未来发展的目标和指导方针。中国与周边国家和广大发展中国家的友好合作进一步加强。中国积极参与亚太经合组织的活动，

NO: 47 Title:北国有寒趣南疆闻花香各地群众喜庆元旦/Author:据新华社北京1月1日电1998年元旦之际，从寒峭的北国大地，到如春的海南宝岛；从美丽的东海之滨，到莽莽的西部边陲，全国各地各族人民纷纷喜庆佳节，决心在新的一年里，高举邓小平理论的伟大旗帜，在党的十五大精神指引下胜利前进。/main Body:元旦的北京，节日气氛格外浓厚。清晨，来自祖国各地及港澳台地区的4万多人，冒着零下7摄氏度的严寒，自发地聚集到了祖国的心脏——天安门广场，参加了庄严的元旦天安门广场升旗仪式。在长安街沿线及大街小巷的高大建筑物上，国旗、彩旗飘扬。京城各大公园、商场及各种娱乐场所也都充满了欢乐气氛。巨大的氢气红灯笼和各种颜色的彩球飘荡在冬日京城的天空中。今天，哈尔滨市的最低气温达到零下25摄氏度，但市民们欢庆新年的热情却没有丝毫减弱。在冰雪覆盖的松花江上，刚刚开放的江上游乐园成了人们游玩的好去处。一座座冰雕构成了另一个崭新的世界，人们穿梭其间，尽情地享受冰雪带来的无穷乐趣。与寒冷的哈尔滨形成鲜明对照，南国春来早，花城今更艳。'98广州风情贺年缤纷大巡游今天在广州隆重举行，为广州新年第一天增添了喜庆的节日气氛，同时也拉开了'98华夏城乡游（广州）活动的序幕。来自香港特区10多个团体的400多位香港同胞首次来广州参加此次贺年活动。海口市阳光明媚。

这两篇文档一个是喜庆新年，一个喜庆元旦，具有一定的相似度。

6. VSM 和 LSI 比较

首先看运算时间，这里因为前面的数据预处理和，tfidf 权重计算都一样，所以就没加入，只考虑了后面 VSM 和 LSI 处理的计算时间

```
"""
函数说明：进行测试
"""

def test():
    #加载数据并生成 tfidf 词向量空间
    documents = load_data('199801_clear.txt')
    corpus_tfidf = tf_idf(documents)
    #对词向量空间进行 VSM 的相应处理并计时
    VSMstart_time = time.time()
    VSM(corpus_tfidf)
    VSMend_time = time.time()
    print('VSM 耗时：' + str(VSMend_time - VSMstart_time) + 's')
    # 对词向量空间进行 LSI 的相应处理并计时
    LSIstart_time = time.time()
```

```
LSI(corpus_tfidf)
LSIend_time = time.time()
print('LSI 耗时: ' + str(LSIend_time - LSIstart_time) + 's')
```

测试结果:

VSM耗时: 11.93734097480774s

LSI耗时: 22.029853582382202s

可以看出 LSI 大概比 VSM 多用 10s，基本上是矩阵进行 SVD 分解，和奇异值计算造成的。但原本的向量空间太过于稀疏，有很多的 0，进行降维后压缩了矩阵，也去除了噪音，使得结果更加精确。

并且 VSM 无法处理文本中隐含的潜语言和一次多意，而 LSI 在压缩矩阵时，减轻了此问题，但 SVD 主动舍去了一部分不重要的特征，所以在降维的同时也要承担数据损失的风险。

7. 参考

<https://www.cnblogs.com/lzlllovesyl/p/5243370.html> SVD 推导

<https://blog.csdn.net/u011630575/article/details/80159100> SVD 推导

<https://blog.csdn.net/zhongkejingwang/article/details/43053513> gensim