

Problems Found in Submissions for Assignment 1

Here are some notes about problems or suggestions on the basis of submissions for Data Mining course Assignment 1. Before reading this note, please read the course and assignment page carefully.

Course web: <http://cs.nju.edu.cn/zlj/Courses.html>

Assignment web: <http://lamda.nju.edu.cn/qianh/dm15.html>

Note 1: About opening raw datasets files

Please use Editors like *Sublime* or *Notepad++* to open raw data files. It may be hard to distinguish between lines posts if we open them via *notepad* (*NOT suggested*). For example, if we use Sublime or Notepad++ to open them, the raw datasets are shown as below:

```
1 南京行货和水货手机的主要购买渠道
大陆行货：（厂家客服保修一年，七天包退，十五天包换，有正规发票，港行不是行货）网购价格比实体店便宜，而且购买过程简单可靠易讯：http://
www.icson.com/（货到付款，价格最低）京东：http://www.360buy.com/（货到付款）新蛋：http://www.newegg.com.cn/（货到付款）苏宁易购：
http://www.suning.com/国美：http://www.gome.com.cn/ http://www.coo8.com/（两个貌似是一回事）欧酷：http://www.ouku.com/
（不能货到付款，必须在线支付）水货：水货有风险，机器无保修，容易纠纷的人，慎重买水货
2 各版本手机客户端下载帖(10S, 安卓, WP)
由于手机客户端版本越做越多，导致置顶不够，故开此贴，将各个版本的手机客户端链接集合于此：i05:SteveJobs版：http://bbs.nju.edu.
cn/bbs0an?path=/groups/GROUP_3/Mobile/D8A298568/D7A3F506BSteveJobs个人网站下载：http://www.pgyer.com/lilyBadPanda版：http://bbs.nju.edu.
cn/bbs0an?path=/groups/GROUP_3/Mobile/D8A298568/DA1698434fetter版
微百合5再次感谢各个版本的开发者对百合的贡献！此后将会把最新更新的版本（不限系统）占满置顶，有新客户端帖发出后，自动取代之旧旧的帖子。--
※ 修改：. Adonis 於 Sep 12 10:38:07 2014 修改本文。[FROM: 210.28.138.10
3 关于小米手机等饥饿营销类产品的进一步处理意见 此前，本版曾有针对包括但不限于小米的购买资格相关帖的规定，链接：http://bbs.nju.edu.
cn/vd46061/bbstcon?board=Mobile&file=M.1358003470.A
```

Note 2: The TF-IDF output

Instead of computing TF-IDF for each word in each class, we need to compute TF-IDF for each post and combine them as a sparse vector. It must be sparse because most words in the word dictionary may not appear in a post and their TF-IDF values will be zero. Then, we will have a sparse feature vector for each post, which can easily be used in following tasks such as classification.

A Wrong result (computing values for each word separately) is shown below:

In D_Computer.txt

```
1 maven 1.0
2 认识 7.0
3 书写 1.0
4 大处 4.0
5 明确 2.0969100130080567
6 学会 2.0969100130080567
7 寻觅 1.0
8 好了 1.3979400086720377
9 较多 3.0
10 招收 10.0
11 最酷 2.0
12 因病 2.0
13 百思买 1.0
```

A Right result:

	Word 1	Word 2	Word 3	Word 4	Word 5	Word 6	Word 7	Word 8
Post1	0	0	0.3356	0	0	0	0.3999		
Post2	0.1246	0	0	0	0	0	0.2378		
Post3	0	0.2467	0	0.6835	0.1376	0	0	0	
...									

1. The TF-IDF feature result only need to contain their values, and the *word, post and line number should not be contained.*
2. Please only use *three or four decimals* (like 3.120, 1.1132) to show your results.
3. Directly write sparse vectors into a file may occupy a lot of spaces. We **RECOMMEND** you to use a **sparse representation** of the sparse feature vectors:
`<index1>:<value1> <index2>:<value2> ...`
(whether use sparse representation or not will not affect your scores)

Note 3: Dimension for each vector

Dimensionality for each sparse vector must be same. Suppose there are N words in the word dictionary, then the length of TF-IDF feature vector for each post should be N.

Note 4: Please check your TF-IDF results carefully

Please check your results *carefully* after preprocessing the data. There are some values like `0.002090063018401983YYYYYYYYYYYYYYYYYY` and `-INFINITY` in the existing uploaded files. All the TF-IDF value *must be* a real number (type of double).

Note 5: About stop words

There may be some messy code in your dictionary or in some posts, you can regard them as stop words and delete them. **We have provided a stop word list, but you can use your own one for better results.**

Note 6: Please check your file before uploading it

Please check your file carefully before uploading. Some uploaded zip file miss the code or data results. The final score will depend on your code, data results and report.

Note 7: The file format

Please use zip file format to compress your file. The same format will help us evaluate your results in a batch style.

Note 8: File and folder name

MG1533001.zip

- code folder(contains all your code or project files)
 - code file (such as tf_idf.py)
- ReadMe.txt (about how to execute your code)
- report.doc
- result folder (10 txt file for TF-IDF feature vectors)

-- Basketball.txt
-- Mobile.txt
-- ...