

R 语言环境下的文本挖掘

刘思喆

© 上海财经大学

2012 年 11 月 3 日



CAPITAL OF STATISTICS
PROFESSION, HUMANITY & INTEGRITY

目录

- ① 文本挖掘的概述
- ② 网页数据抓取的利器 -XML
- ③ tm 包及相关应用
- ④ 应用的实例

目录

① 文本挖掘的概述

② 网页数据抓取的利器 -XML

③ tm 包及相关应用

④ 应用的实例

文本挖掘的处理流程
相关的 R 包

文本挖掘的一般流程

对于文本处理过程首先要拥有分析的语料 (text corpus), 比如报告、出版物、网页文章等。而后根据这些语料建立半结构化的文本库 (text database), 生成包含词频的结构化的词条 - 文档矩阵 (term-document matrix)。

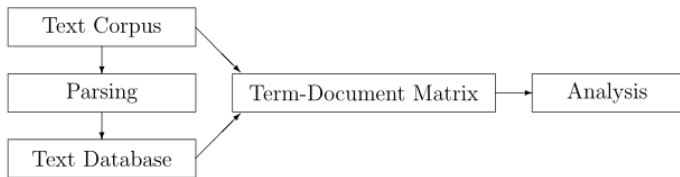


Figure: 文本挖掘的处理流程

文本挖掘的适用范围

解析后的结构化数据会被用于后续的分析，比如：

- 语法分析；
- 信息提取和修复；
- 文档信息汇总，比如提取相关有代表性的关键词、句子等。
- 文本分类，比如根据现有的文本分类情况，对未知文本进行归类；
- 其他

文本挖掘相关的 R 包

- XML
- tm
- topicmodels
- RWeka, lsa, RTextTools, zipfR, TextRegression, wordcloud

目录

① 文本挖掘的概述

② 网页数据抓取的利器 -XML

网页数据的获取函数

实际的例子

③ tm 包及相关应用

④ 应用的实例

XML 文件的解析

在 R 中对网页解析 (XML、HTML 文件, 或包含 XML、HTML 的字符串) 有多种方法, 比较成熟的方法是使用 **XML** 包。这个包能够将 XML、HTML 网页树 (tree) 解析成 R 结构数据。

解析函数

对标准 XML 文件的解析函数 `xmlParse`, 以及适应性更强的 `htmlTreeParse` 函数, 这些函数都拥有大量的参数来适应解析需要。

解析的 XML 文件

```
<doc>
  <a xmlns:omegahat="http://www.omegahat.org">
    <b>
      <c>
        <b/>
      </c>
    </b>
    <b omeegahat:status="foo">
      <r:d>
        <a status="xyz"/>
        <a/>
        <a status="1"/>
      </r:d>
    </b>
  </a>
</doc>
```

解析过程

```
1 doc <- xmlParse(system.file("exampleData", "tagNames.xml",  
2                          package = "XML"))  
3 els <- getNodeSet(doc, "/doc//a[@status]")  
4 sapply(els, function(el) xmlGetAttr(el, "status"))  
5 [1] "xyz" "1"
```

XML 支持的功能远不止如此，可参考 XML 包的帮助文档

解析过程

```
1 doc <- xmlParse(system.file("exampleData", "tagNames.xml",
2                             package = "XML"))
3 els <- getNodeSet(doc, "/doc//a[@status]")
4 sapply(els, function(el) xmlGetAttr(el, "status"))
5 [1] "xyz" "1"
```

XML 支持的功能远不止如此，可参考 XML 包的帮助文档

获得国家地震科学数据共享中心的地震数据

- 我们先观察一下它的网页
- 再看看怎么抓取
- 其实这就是一个“爬虫”

获得国家地震科学数据共享中心的地震数据

- 我们先观察一下它的网页
- 再看看怎么抓取
- 其实这就是一个“爬虫”

获得国家地震科学数据共享中心的地震数据

- 我们先观察一下它的网页
- 再看看怎么抓取
- 其实这就是一个“爬虫”

目录

① 文本挖掘的概述

② 网页数据抓取的利器 -XML

③ tm 包及相关应用

④ 应用的实例

中文分词

词条 -文档关系矩阵

中文分词工具的 R 包 rmmseg4j

rmmseg4j 调用了 Java 分词程序对语句进行分词。下面是一个例子：

```
library(rmmseg4j)
> mmseg4j('花儿为什么这样红')
[1] "花儿 为什么 这样 红"
```

这段话被分为了“花儿”，“为什么”，“这样”，“红”。这个空格分隔的结果集，后期就可以按照英文的方式进行后续处理。

由于 **rmmseg4j** 这个包已经发布 CRAN 仓库上，正常安装：

```
install.packages("rmmseg4j")
```


语料库

在 tm 中主要的管理文件的结构被称为语料库 (Corpus) , 代表了一系列的文档集合。

动态语料库 :

```
1 Corpus(x,  
2   readerControl = list(reader = x$DefaultReader, language = "en"),  
3   ...)
```

静态语料库 (需要 filehash 包的支持) :

```
1 PCorpus(x,  
2   readerControl = list(reader = x$DefaultReader, language = "en"),  
3   dbControl = list(dbName = "", dbType = "DB1"),  
4   ...)
```

语料库的处理

作用	函数
转化纯文本	<code>tm_map(reuters, as.PlainTextDocument)</code>
去除空白	<code>tm_map(reuters, stripWhitespace)</code>
小写变化	<code>tm_map(reuters, tolower)</code>
停止词去除	<code>tm_map(reuters, removeWords, stopwords("english"))</code>
剔除数字	<code>tm_map(ovid, removeNumbers)</code>
...	...

词条 - 文档关系矩阵

词条 - 文档关系矩阵是后续构建模型的基础。假设我们有两个文档分别是 text mining is funny 和 a text is a sequence of words，那么其映射到空间上对应的矩阵为：

	a	funny	is	mining	of	sequence	text	words
Doc 1	0	1	1	1	0	0	1	0
Doc 2	2	0	1	0	1	1	1	1

在 tm 包里，根据词条、文档分别作为行、列或反之，对应两种矩阵：

- TermDocumentMatrix
- DocumentTermMatrix

文档 - 词条关系矩阵的相关操作

- `findFreqTerms(dtm, 5)`
- `findAssocs(dtm, "opec", 0.8)`
- `removeSparseTerms(dtm, 0.4)`

延伸阅读

- ① tm 的扩展，包括并行计算
- ② 根据条件进行文档过滤 (`tm_filter`)
- ③ 元数据的管理 (`meta` , `DublinCore`)
- ④ 字典 (`Dictionary`) 的应用

目录

① 文本挖掘的概述

② 网页数据抓取的利器 -XML

③ tm 包及相关应用

④ 应用的实例

wordcloud

文档识别和聚类

主题模型 (topic model)

关键词网络

文档识别的原理

假如我们有下表这样的矩阵，那能做分类么？

Table: 矩阵示例

Y	data	big	market	last	google	new	...
1	0	1	1	0	0	0	
1	0	0	0	0	1	1	
0	1	0	9	1	2	2	
0	0	0	2	0	1	1	
1	0	0	1	0	0	0	
0	0	1	1	0	1	0	
...							...

一般的 Classification 类的方法都可以使用在这个矩阵上，比如从最简单的 knn，稍稍高级点的 regression、decision tree，甚至 SVMs 等方法都可以使用。

主题模型简介

主题模型是专门抽象一组文档所表达“主题”的统计技术。最早的模型是 probabilistic latent semantic indexing (PLSI)，后来 Latent Dirichlet allocation (LDA，潜在狄利克雷分配模型) 模型成为了最常见的主题模型，它可以认为是 PLSI 的泛化形式。

LDA 主题模型涉及到贝叶斯理论、Dirichlet 分布、多项分布、图模型、变分推断、EM 算法、Gibbs 抽样等知识。

详细参考统计之都主站文章：

http://cos.name/2010/10/lda_topic_model/

应用实例 (一)

在某个门户网站上收集了 720 篇文章，对这些文档构建 LDA 模型，参数设置为 6 类，模型提取出 6 个类别的关键词分别是美食、旅游、拍摄、购物、汽车和招聘，每个关键词对应的文档数量如下图所示：

Table: 预测种类及主题关键词

Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6
美食	旅游	拍摄	购物	汽车	招聘
102	140	118	110	129	121

应用实例 (二)

实际上收集的这些文章分属不同频道: " 购物类", " 旅游类", " 美食类", " 汽车类", " 数码类", " 招聘类", 共计六类。对比网站编辑的人工标记同模型标记的差异情况:

Table: 预测种类同人工判别种类的混淆矩阵

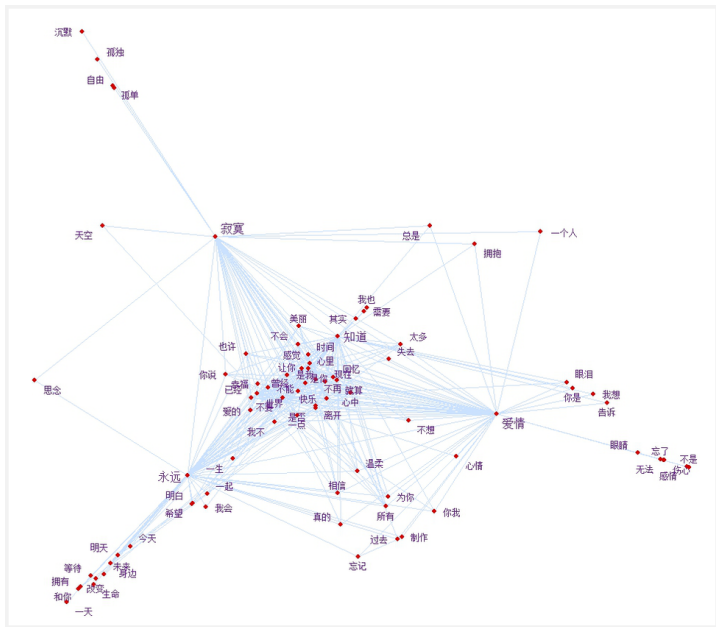
	Topic 1	Topic 2	Topic 3	Topic 4	Topic 5	Topic 6
购物类	9	1	3	91	15	1
旅游类	0	113	0	4	0	3
美食类	90	24	0	5	1	0
汽车类	3	0	0	7	110	0
数码类	0	2	115	0	3	0
招聘类	0	0	0	3	0	117

注: 在使用 topic model 过程中设置了 6 类的参数是为了和实际情况进行比对

关键词网络

- 22996 首中文流行歌曲的歌词 (844 位歌手)
- 这些歌词是怎么组织在一起的？

爱情
永远
寂寞



- 邮件: `sunbjt<at>gmail.com`
- 博客: `http://www.bjt.name`
- 微博: @刘思喆

Jump to first slide