

主流视频网站评论分析系统

OLAP 系统搭建

沈祥壮

中山大学

2019 年 10 月 29 日

1 项目概览

2 OLAP 概览

3 数据集描述

4 OLAP 系统搭建流程

5 登录 OLAP 系统

6 可视化结果展示

项目概览

- 项目的核心是对主流的一些视频网站的评论系统进行深入的分析, 力图增强对用户反馈的理解, 以便更好的完善社区的服务。
- 本次项目以国内最大的视频网站 bilibili 为研究对象, 手动抓取网站的评论数据, 并将数据导入数据仓库, 建立 OLAP 系统。
- 此外, 整个项目前期主要是以探索性数据分析为主, 后面会陆续将新的数据融合到数据仓库, 并对数据进行深入分析。

OLAP 概览

本次 OLAP 系统的搭建过程中, 使用的数据源来自之前自己写的 B 站视频评论的爬虫, 即 Bilibili-Comments-Spider 。OLAP 系统的搭建采用了 Python 的 Cubes 框架, 可视化部分采用 Cubesviewer 。在使用过程中, 对框架本身的执行效率进行了优化, 主要是数据库插入数据的优化, 极大提升了程序的运行效率。

数据集描述

本次数据集 (*BilibiliGaomu.csv*) 来自 B 站其中一部番剧的评论数据, 共包含 *mid*, *username*, *rpId*, *gender*, *content*, *ctime*, *likes*, *rcount* 八个字段, 共 77474 条评论数据。因为框架的具有十分良好的可扩充性, 经过简单的预处理 (*Preprocessing.py*), 可以十分方便地将之前爬虫采集的数据集成进去。处理后的数据集 (*data.csv*) 包含如下字段。

| 字段 | 含义 |
|-----------------|-------|
| <i>sex</i> | 评论人性别 |
| <i>username</i> | 评论人昵称 |
| <i>likes</i> | 评论获赞数 |
| <i>wordnum</i> | 评论字数 |
| <i>month</i> | 评论月份 |

OLAP 系统搭建流程

- ❶ 将抓取的数据导入数据仓库
- ❷ 构建数据模型
- ❸ 配置本地 OLAP 服务器
- ❹ 配置可视化服务
- ❺ 将服务部署到远程服务器
- ❻ 登录到 OLAP 系统
- ❼ 开始进行联机分析

登录 OLAP 系统

目前已经将项目部署到服务器之上，登录搭建好的 OLAP 系统流程只需要如下三个步骤：

- 1 使用浏览器打开文件
`bilibili/cubesviewer/html/studio.html`
- 2 在弹出的对话框中输入服务器地址
`http://119.23.209.92:5009/`
- 3 开始进行联机分析

可视化结果展示

