

基于多源异构数据的“新高考”平行志愿智能填报系统的研发

张彦俊¹ 简卓然² 陈丽煌³ 张继涛⁴

中山职业技术学院 广东 中山 528404

[摘要]本文是在实行平行志愿投档录取模式的背景下,根据平行志愿投档录取模式的特点、规律等,结合“新高考”的政策、要求等,运用大数据技术,基于多源异构的高考大数据,开发出一套既适合平行志愿投档录取模式,又适应“新高考”政策下平行志愿智能填报综合服务系统,面向所有实行平行志愿及“新高考”省份参加高考的考生或家长,提供“新旧高考”政策解读、“新旧高考”平行志愿智能填报、大学学业规划、职业生涯规划等咨询服务活动。

[关键词]多源数据; 高考志愿; 系统开发与应用

[DOI] 10.12252/j.issn.2096-6261.2021.08.1471

一、研究背景

新一轮高考制度的改革正带来高考志愿智能填报思维方式的转变。作为中央部署全面深化改革的重大举措之一,关于考试招生制度改革的实施意见。目前高考改革方案的亮点在于两方面,一是“2015年起推行自主招生安排在全国统一高考后进行”;二是“创造条件逐步取消高校招生录取批次。改进投档录取模式,推进并完善平行志愿投档方式,增加高校和学生的双向选择机会。2015年起在有条件的省份开展录取批次改革试点。”原因很简单,如果录取制度不进行深层次改革,在当前分批次按计划集中录取制度框架下,文理不分科、英语一年两次考,并不能改变基础教育的应试格局。因此,在当前分批次按计划集中录取制度框架下,研究各批次录取投档的规律,对于考生和家长来说,具有十分重要的现实意义。

目前,共有三批的省份实行了新高考。第一批新高考(2014-2017年),上海市:实行“3+3”和“6选3”,即语文、数学、外语三科必考;浙江省:实行“3+3”和“7选3”,即语文、数学、外语(英语、俄语、德语、法语、日语等)三科必考。第二批新高考(2017-2020年),北京、天津、海南、山东等4个省份,实行“3+3”和“6选3模式”,即语文、数学、外语三科必考,学生要从物理、化学、生物、政治、历史、地理等6个学科中,选出3个学科作为高考选考科目。第三批新高考(2018-2021年),河北、辽宁、湖北、湖南、重庆、江苏、福建、广东等8个省份作为第三批“新高考”试点省份,推行了“3+1+2”的新高考模式。第四批新高考(2021-2024年),黑龙江、吉林、安徽、江西、广西、贵州、甘肃等7个省份,将成为第四批“新高考”试点省份,推行“3+1+2”的新高考模式。

二、研究内容

在实行平行志愿投档录取模式的背景下,根据平行志愿投档录取模式的特点、规律等,结合“新高考”的政策、要求等,运用大数据技术,基于多源异构的高考大数据,开发出一套既适合平行志愿投档录取模式,又适应“新高考”政策下平行志愿智能填报综合服务平台——“大学导航”,面向所有实行平行志愿及“新高考”省份参加高考的考生或家长,提供“新旧高考”政策解读、“新旧高考”平行志愿智能填报、大学学业规划、职业生涯规划等咨询服务活动。

1. 面对多源异构的高考大数据采集、处理

项目采用数据爬取的方法对各省往年和今年高考大数据进行采集、处理与建库,在数据爬取中,多种数据获取渠道与多种数据渲染形式为高考大数据自动采集带来挑战,项目拟通过构建“中间件”,将这些多源异构数据进行自动采集并统一为同一类型的数据,具体包括以下几种格式:

(1) PDF 表格提取

采用技术为 PDFplumber: PDFplumber 找到可见或不可见的表格线;PDFplumber 根据表格边界线确定表格之间的交点;根据交点,确认包围的最小单元格;把单元格融合在一起,孵化检测出表格对象。

(2) Excel 文档表格提取

采用 Python 中的依赖库 pandas,基于 read_excel(url) 方法,读取页面表格,再采用建立 csv 文件存入,即可实现页面表格提取归一化。

(3) 图片表格提取

选择特别的图像元素采用膨胀腐蚀方法提取表格结构,将表格结构继续细化,提取完全的表格图像;提取图像后,定位表格单元格,将其切分为小单元结构。

(4) HTML/XML 数据提取

运用正则化、JSON 串、DOM 树、Webmagic 等工具对 HTML、XML 等格式网页源码进行提取。

2. 建立核心基本库

(1) “批次控制线”

运用所学数据爬取方法,在老师的指导下,对各省(市)当年往年的高考“批次控制线”大数据进行采集、整理、分类、加工、处理、制表,从而建立起各省(市)高考“批次控制线”数据库。

(2) “一分一段表”

运用所学数据爬取方法,在老师的指导下,对各省(市)当年往年的高考“一分一段表”大数据进行采集与网络爬虫、清洗与存储挖掘、分析与应用处理,从而建立起各省(市)高考“一分一段表”数据库。

(3) “最低投档分”

运用所学专业知识和在老师的指导下,对各省(市)当年往年的高考“最低投档分”大数据进行采集与网络爬虫、清洗与存储挖掘、分析与应用处理,从而建立起各省(市)高考“最低投档分”数据库。

3. “大学导航”高考志愿智能填报系统研究开发与完善

(1) “大学导航”高考志愿智能填报系统数据结构建立

运用所学专业知识,在老师的指导下,以上述“批次控制线”、“一分一段表”、“最低投档分”为依据,利用高考数据运算数学模型,在前期开发的“大学导航”高考志愿智能填报系统 1.0 版的基础上,进一步完善系统数据结构,为 2.0 版系统开发升级搭建框架。

(2) “大学导航”高考志愿智能填报系统研究与开发

高考志愿智能填报系统运用所学的《WEB 前端应用程序开发》、《数据库设计与应用》、《程序设计与应用》、《Spring Boot 企业级开发》、《JavaScript 高级程序设计》等专业课程知识与技能,在老师的指导下,围绕高考志愿智能填报准备、高考志愿关键信息、高考志愿“海选”方

