

# 不同领域的词库构建

目前存在的问题：

- 词袋的难点在于目前人工构建的词袋较小，扩展起来较为麻烦。

解决方法：

- 尝试通过提取目前常用的输入法的词库进行热词、不同领域的词的扩展，
  - 优势：快速、高效的进行词库扩充

## 领域词汇表研究

---

词汇表定义：

- 词汇表是专业或技术词汇及其含义的列表。

参考：

- [维基百科--目录](#)

分类

### 1. 文化艺术

- 文化
- 艺术
- 小说和文学
- 运动
- 收藏和爱好

### 2. 地理和地方

- 地方

### 3. 健康和健身

- 药物、临床研究、精神病学、解剖学

### 4. 历史和事件

- 历史、考古

### 5. 数学和逻辑

- 数学领域、几何、

### 6. 自然科学和物理科学

- 生物科学
  - 生物学、植物学、生态学、病毒学
- 物理科学
  - 天文(气象)、化学、气候变化、地理、地质学、物理学

### 7. 哲学与思考

- 哲学
- 8. 社会和社会科学
  - 商业、政治、法律
  - 语言和俚语
- 9. 技术与应用技术
  - 技术与应用科学
  - 计算和信息技术
  - 运输

## 领域含义说明

---

数据源：

- wikipedia

### 维基百科：内容/文化与艺术

文化，意思是“居住、培养或荣誉”，一般来说文化指人类活动，文化的不同定义反映了不同的理解理论或评价人类活动的标准。现在的人类学家使用该术语来指代人类对经验进行分类并以符号方式对其进行编码和交流的普遍能力。

艺术是文化的一个巨大分支，由许多创造性的努力和学科组成。从更广泛的意义上讲，艺术是创造力或想象力的表达，艺术也可以被理解为与创造力、美学和情感的产生有关。

### 文化

- 艺术
  - 文化
  - 视觉艺术
  - 表演艺术
- 美食
  - 食物
  - 饮食
- 娱乐
  - 节假日
  - 游戏
  - 体育
    - 球类
    - 格斗运动
    - 赛车运动
- 人文学科
  - 地区研究
  - 古典研究

## 维基百科：内容/地理和地点

地理学，是对地球及其特征、居民和现象的研究。地理研究的四大历史传统是空间分析、自然和人类现象的研究(地理作为分布研究)，区域研究(地方和地区)，人地关系研究和地球科学研究。尽管如此，现代地理学是一门包罗万象的学科，它首先寻求了解世界及其所有人类和自然的复杂性--不仅是物体在哪里，而且需要解释它们是如何变化和形成的。作为“人文科学与自然科学之间的桥梁”，地理学分为两个主要分支--人文地理学和自然地理学。

### 地理

- 地理学分支

- 自然地理学
- 人文地理学
- 区域地理学

- 地理和地方

地球、世界、水域、城市、社区、大陆、国家、沙漠、湖泊、地貌、山脉、海洋、人口稠密的地方、保护区、地下、城市、村庄

- 自然地理特征

- 地貌
- 水体

- 人造地理特征

- 世界

- 按半球
- 按文化区域
- 按经济发展
- 按大陆

## 维基百科：目录/社会和社会科学

社会是一群人组成的一个半封闭的系统，最简单的社会一词指一大群分享自己的文化和制度的人，改词起源于拉丁语，意指“与他人的友好关系”。社会的意义与被认为是社会的东西密切相关。隐含在社会的含义是：它的成员可能分享一些共同的关注和利益，共同的目标或共同的特征。社会科学是一组研究世界人类方面的学科。它们与艺术和人文科学的不同之处在于，社会科学倾向于强调在人文研究中使用科学方法，包括定量和定性方法。

### 社会科学

- 社会科学

- 人类学、考古学、认知科学、传播研究、批判理论、文化研究、经济学
- 教育、地理、历史、语言学、法律、政治学

- 社会

- 社区
- 社会发展
- 社会机构

- 家庭
- 基础设施
- 经济与商业
- 教育
- 文明社会
- 政府与政治
- 社会网络

## 维基百科：目录/自然与物理科学

在科学中，自然科学一词是指研究宇宙的理性的方法，被理解为遵守自然起源的规律或规则。“自然科学”一词也区别于使用科学方法研究自然的领域与使用科学方法研究人类行为和社会的社会科学领域，并且来自使用不同方法论的形式科学，例如数学和逻辑。与生物科学相比，物理科学是自然科学和研究非生命系统的一个涵盖性术语。

## 生物学

- 解剖学、生物化学、植物学、细胞生物学、计算生物学、生态学、进化生物学、
- 基因学、组织学、免疫学、微生物学、生命起源、生理学、动物学

## 物理科学

- 天文学
- 化学
- 物理学
  - 原子、基本粒子、力、引力、质量、物质、光学、时间、声音、能量、经典力学

## 维基百科：内容/历史和事件

历史是对过去事件、社会和文明的解释。历史一词来自希腊语，“一个人的调查记录”。1911年的《大英百科全书》指出：“广义上的历史就是发生的一切，不仅是人类生活的所有现象，还有自然界的现象，它是发生变化的一切。”

## 历史

- 按时期划分
  - 史前
  - 古代历史
    - 古代东方
    - 古代西方
  - 中世纪
  - 文艺复兴
  - 早期现代史

- 现代历史
- 全球化
- 按地区划分
  - 古埃及
  - 古罗马
  - 拜占庭帝国
  - 奥斯曼帝国
- 按主题划分
  - 按领域划分
    - 艺术史
    - 商业历史
    - 地理学史
    - 科学史
    - 技术史
  - 历史科学
    - 考古学
    - 地质学
    - 生物学
    - 天文学

## 参考资料:

---

- **THUOCL: 清华大学开放中文词库**
  - IT、财经、成语、地名 历史名人、诗词、医学、饮食、法律、汽车、动物
- **搜狗输入法-词库**
  - 城市信息、自然科学、人文科学、社会科学、工程与应用科学、农林渔畜、医学、艺术、生活.....

## 词库分类

<b>城市信息大全</b> 全国 广东 北京 上海 云南 安徽 四川 江苏 辽宁 国外地名  <b>你所在的城市:广东省广州市</b> 广州市城市信息精选 广州市公交站名	<b>电子游戏</b> 单机游戏 网络游戏 网页游戏  王者荣耀 英雄联盟 梦幻西游 魔兽世界 阴阳师 剑网3 龙之谷	<b>自然科学</b> 生物 化学 数学 物理 其它 天文学 气象学 海洋学 地理地质  动物词汇大全 生物词汇大全 地理地质词汇大全 物理词汇大全	<b>人文科学</b> 语言 文学 宗教 历史 哲学 神学 考古 其它 伦理学  网络流行新词 宋词精选 成语俗语 唐诗300首 古诗词名句
<b>社会科学</b> 法律 军事 其它 心理学 政治学 房地产 社会学 伦理学 档案学  法律词汇大全 财经金融词汇大全 心理学词汇大全	<b>工程与应用科学</b> 建筑 化工 其它 造纸 包装 计算机 机械工程 电子工程 钢铁冶金  计算机词汇大全 建筑词汇大全 电力词汇大全 船舶港口词汇大全	<b>农林渔畜</b> 农业 林业 渔业 畜牧业  农业词汇大全 林业树种名词库 农学 土壤学名词	<b>医学</b> 中医 疾病 中药 医疗 其它 外科 兽医 针灸 西药学  医学词汇大全 中医中药大全 中外药品名称大全 医疗器械大全
<b>艺术</b> 绘画 曲艺 音乐 摄影 舞蹈 陶瓷 雕塑 书法篆刻 刺绣织染  绘画美术词汇大全 书法词汇大全 戏剧戏曲词汇大全 二人转词汇大全	<b>运动休闲</b> 球类 武术 其它 垂钓 奥运 气功 轮滑 棋牌类 太极拳  象棋 篮球 太极拳 羽毛球 足球 F1赛车	<b>生活</b> 理财 饮食 旅游 习俗 日常 服饰 礼品 美发 办公文教  股票基金 家居装修词汇大全 饮食大全 旅游词汇大全	<b>娱乐</b> 动漫 明星 汽车 收藏 烟草 宠物 其它 魔术 模型  日剧、动漫大全 汽车词汇大全 歌手人名大全 热门电影大全

- 百度输入法
  - 人文社科、生活百科、文化艺术、人名专区.....

## 词库分类

<b>城市区划</b> >  全国 广东 江苏 陕西 四川 山东 安徽 湖南 北京 浙江 辽宁 广西	<b>理工行业</b> >  医学 计算机 土木建筑 生物学 化学化工 军事学 机械工程 物理学	<b>人文社会</b> >  文学 语言 历史学 经济金融 宗教学 法律 哲学 公共管理 心理学
<b>电子游戏</b> >  网络游戏 单机游戏 手游 网页游戏 其他	<b>生活百科</b> >  旅游 理财投资 品牌 餐饮 医药 服饰 日常 家居 宠物 烟酒茶	<b>娱乐休闲</b> >  颜文字 卡通动漫 网络聊天 影视 音乐 汽车 棋牌 表演
<b>人名专区</b> >  常用 明星 名人 其他	<b>文化艺术</b> >  流行小说 戏曲 绘画 摄影 报刊广告 书法 工艺 相声评书 写作	<b>体育运动</b> >  球类 射击 水上 赛车 自行车 奥运 垂钓 登山 极限 滑冰 武术

- QQ输入法
  - 自然科学、社会科学、医学医药、生活百科.....

搜索

词库分类

城市信息

自然科学

社会科学

工程应用

农林渔畜

医学医药

电子游戏

艺术设计

生活百科

运动休闲

文学

默认

全部

言情

武侠

水浒传

诗词歌赋

三国演义

奇幻玄幻

民间文学

流行作品

科幻

经典名著

红楼梦

最炫文言风

说明：把一些常用语转化成文言文，使别人听不懂，也显得自己有才。  
词汇：16条      大小：10432KB      下载次数：1682  
更新时间：2015-03-09 19:52:05

大道争锋

说明：起点连载小说《大道争锋》词库  
词汇：740条      大小：31124KB      下载次数：1226  
更新时间：2016-12-11 21:42:16

天下3文韵墨香答题词库

说明：文韵墨香问题答案  
词汇：12条      大小：10520KB      下载次数：1613  
更新时间：2015-05-09 09:21:42

dsfgdf

搜索算法的词库解码方案：

- ☒ 搜狗输入法
  - scelParser
- ☒ 百度输入法
  - bdictParser
- ☐ QQ输入法