



R环境短文本分类器与电商品类数据研究

Bird 张翔

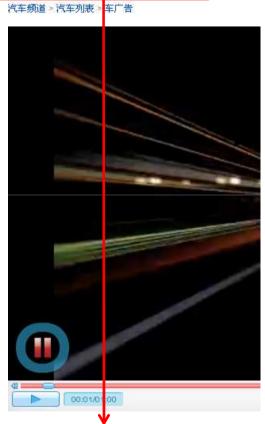
北京•上海•广州•深圳•东京•硅谷•香港

www.iresearch.com.cn

为什么是短文本



视频: 全新奥迪A6L 巅峰感受



- 1. 提取标题
- 2. 分词 与自动分类
- 3. 类目与词频统计





短文本分类的核心目的





2012年5月面部护理类商品销售额标签云



Source: 视频数据来自iVideoTracker, 2012.4. 基于对20万名家庭及办公(不含公共上网地点)样本网络视频行为的长期监测数据获得,不包括在线视频客户端数据。

电商数据来自艾瑞对20万名家庭及办公网民网络购物行为监测所得。

短文本分类难点





扩充词库和关键词人工列表



● 淘宝专用词库【官方推荐】 🖁

词条样例:阿半、阿迪板鞋、阿迪杰、阿尔岱雪、阿尔法、阿尔卡特、阿尔塔、阿福贝贝、阿格利司、阿九的店、阿卡、阿卡邦、阿卡住、阿卡随心手艺、阿珂姆、阿拉蕾、阿里通、阿里扎、阿玛迪斯...

📕 淘宝专用词库【官方推荐】

文件(P) 编辑(E) 格式(Q) 查看

阿半

阿迪板鞋

阿迪杰

阿尔岱雪

阿尔法

阿尔卡特

阿尔塔

阿福贝贝

阿格利司

阿九的店

阿卡

阿卡邦

阿卡佳

阿卡随心手艺

阿珂姆

阿拉蕾

阿里通

阿里扎 阿玛迪斯

网玛迪斯 阿玛尼

阿玛施

阿墨瓜丝

Term	Category
保温杯	家居用品/杯具
笔记本电脑	电脑/笔记本
餐桌	家具/桌子
车载充电器	汽车用品/电器/配件
低音炮	音响
电话卡	手机/号码卡-电话卡
儿童座椅	座椅

R文本分析包:

Tm

Rmmseg4j

Smartcn

wordcloud

遗留的问题

BB霜 美的电器

很美的电器贴图

处理流程(借鉴KDDCUP2005)



谷歌关闭了分类目录 使用人工特征词列表建立 最初训练集

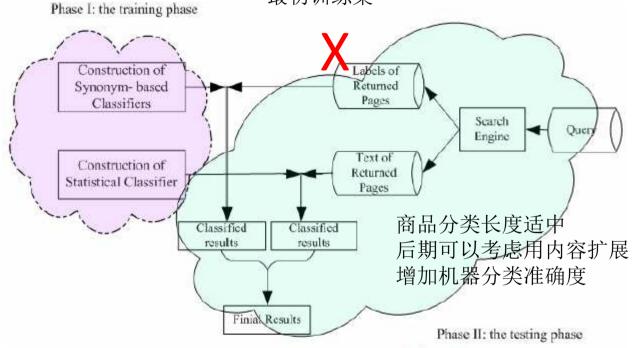
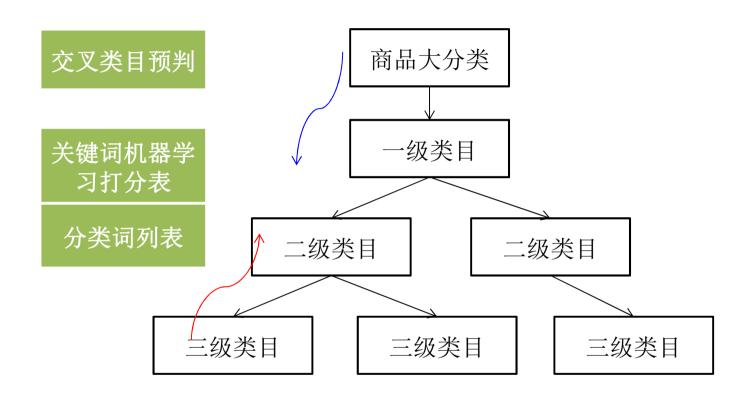


Figure 1. The architecture of our approach Q2C@UST

混合分类器流程





有人工分类表参与的分类准确性接近90%,但是词库建立耗时,目前积累的词库商品覆盖率达70% 只依赖机器自学习的分类准确性为60%-80%,但是在做大类 预判的时候覆盖率高

贝叶斯从朴素到"奢华"



- 朴素贝叶斯的假设条件
 - 词之间没有任何依赖关系,即条件独立;

$$P(d \mid c_i) = P(w_1, w_i, ..., w_n \mid c_i) = \prod_{1 \le k \le n} P(w_k \mid c_i)$$

- 贝叶斯网络
 - 人的思维是有依赖关系的
 - " 苹果 iphone "
 - "____苹果__红富士__"
 - 用一个有向无环图和一个条件概率表集合来表示一组变量的联合概率分布

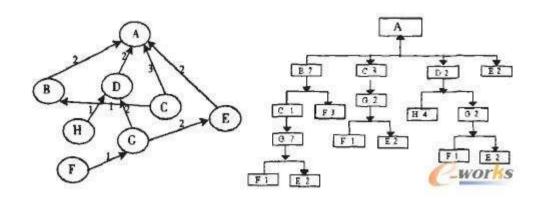
$$P(x_1, x_2, ..., x_n) = \prod_{i=1}^{n} P(x_i | Parents(x_i))$$

- 现实是残酷的
 - 对于文本的词量,条件网络过于复杂,可以直接使用"词对"预测(半朴素贝叶斯)

网状类目结构



- 普通目录结构是树状结构固定深度
 - 目前艾瑞使用的是3+1级目录结构,列表形式储存
- 真实类目前存在交叉,是网状关联结构
 - 计算机语言中可以做面向对象结构设计
 - R语言中可以使用sna或igraph的图矩阵格式



结果应用-店铺数据



▶ 剃须护理、汽车电子淘宝店铺垄断性强

- 马太效应显著的类目有剃须护理及汽车电子,两类销售额10%的店铺占据了约9成的业绩 份额。
- 相比之下,部分家具、饮料类目店铺竞争相对充分。

2012年5月淘宝部分类目销售额在店铺的分布情况

淘宝集市					天猫				
类目	TOP10%店铺 销售额占比		TOP30%店铺 销售额占比		类目	TOP10%店铺 销售额占比		TOP30%店铺 销售额占比	
剃须护理		92%		98%	剃须护理		96%		99%
音响播放		89%		98%	汽车电子		94%		99%
汽车电子		88%		98%	休闲食品		86%		97%
住宅家具		69%		92%	办公家具		63%		90%
碳酸饮料		67%		91%	含乳饮料		63%		91%
办公家具		66%		89%	矿泉水		39%		83%
	剃须护理 音响播放 汽车电子 住宅家具 碳酸饮料	类目TOP1 销售剃须护理 音响播放 汽车电子住宅家具 碳酸饮料	共同TOP10%店铺 销售额占比剃须护理92%音响播放89%汽车电子88%住宅家具69%碳酸饮料67%	共日TOP10%店舗 TOP30 销售额占比 销售额 3剃须护理92%音响播放89%汽车电子88%住宅家具69%碳酸饮料67%	共日TOP10%店舗 TOP30%店舗 销售额占比 销售额占比销售额占比 销售额占比剃须护理92%音响播放89%汽车电子88%住宅家具69%碳酸饮料67%	类目TOP10%店舗 TOP30%店舗 销售额占比 销售额占比 销售额占比类目剃须护理92%98%剃须护理音响播放89%98%汽车电子汽车电子88%98%休闲食品住宅家具69%92%办公家具碳酸饮料67%91%含乳饮料	类目TOP10%店舗 TOP30%店舗 YEDP1销售额占比销售额占比销售额占比剃须护理92%98%剃须护理音响播放89%98%汽车电子汽车电子88%98%休闲食品住宅家具69%92%办公家具碳酸饮料67%91%含乳饮料	类目TOP10%店舗 TOP30%店舗 销售额占比共享日本 TOP10%店舗 销售额占比剃须护理92%98%剃须护理96%音响播放89%98%汽车电子94%汽车电子88%98%休闲食品86%住宅家具69%92%办公家具63%碳酸饮料67%91%含乳饮料63%	类目TOP10%店舗 TOP30%店舗 销售额占比 销售额占比 销售额占比 销售额占比 销售额占比 销售额占比 销售额占比 销售额占比 销售额占比 销售额的

|汪祥:表甲列至了此次研究品奕中,10%店铺销售钡占比载多及载少的二个商品奕目。

来源:基于情报通抓取的公开交易数据。

细分商品-电视品类分布



▶商城买家电 C店淘配件

- 在主要B2C网站中,用户购买的电视类商品以电视机本身为主,电视机品牌(索尼、创 维、TCL)、配置说明(全高清、LED)等描述的商品名称销售额较高。
- 淘宝集市中,除了电视机外,电视周边商品也有较高的关注度:机顶盒、电视开关、遥控器、贴纸等品类也占据了较多的份额。

2012年5月电视商品销售额标签云





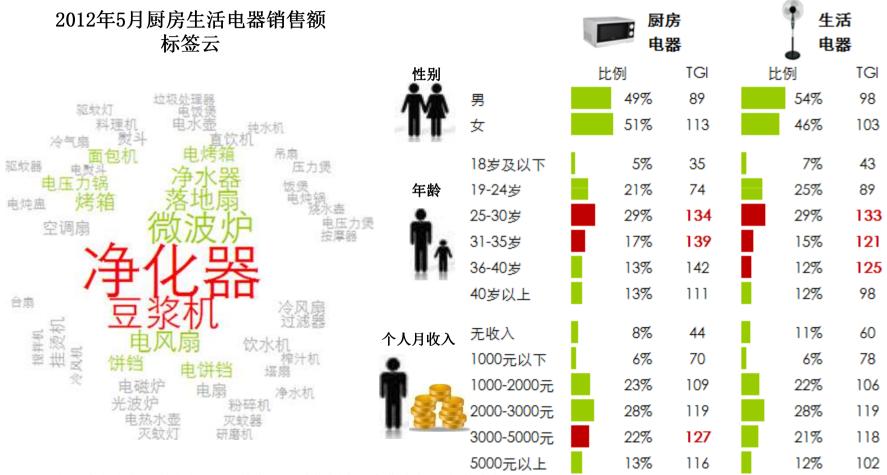
来源:基于艾瑞对20万名家庭及办公网民网络购物行为监测所得。

细分类目-厨房&生活电器



▶ 年轻、高收入者是厨房生活电器类主要访问者

2012年5月特定类目访问者基本属性



注释: 左图中仅包括了销售额TOP50的商品; 销售额为淘宝集市与网上商城合计值;

右表中TGI是目标群体指数,TGI指数=(目标群体中某一特征的群体所占比例/总体中相同特征的群体所占的比例)*标准数 契额:淘宝集市数据基于情报通抓取的公开交易数据,网上商城总体数据基于艾瑞对20万名家庭及办公网民网络购物行为监测所得。 12

热销品牌-厨房电器类



▶淘宝集市与网上商城热销品牌重合度高

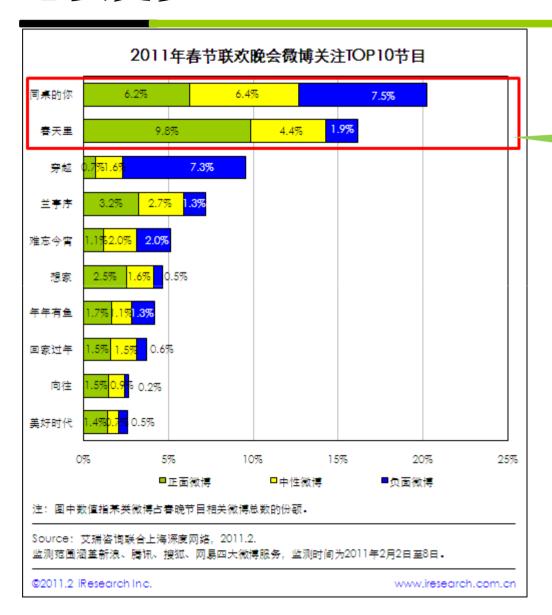
厨房电器销售额较高的品牌中,九阳、美的及苏泊尔在淘宝集市及网上商城排名均居前列,两类电商网站的热销品牌存在较高的重合。

2012年5月厨房电器类销售额第一阵营品牌墙



想要更多





草根节目在网民评论中好评高

感谢聆听,期待合作!



选择艾瑞,

选择值得信任的合作伙伴!

