

垂直知识图谱构造工具与行业应用

阮彤 自然语言处理与大数据挖掘实验室 主任 华东理工大学

主要内容

- ■为什么垂直行业需要知识图谱
- ■垂直知识图谱的特点
- ■垂直知识图谱工具——VKG Builder介绍
- ■垂直知识图谱应用

我们从通用知识图谱开始。。。。

1. SSCO

使用机器学习获得的知识网络,包括260,345个概念,5,602,180个实例,526,219个同义关系,下图显示了网络的部分节点



http://ssco.zhishimofang.com/

2.zhishi.me





对不同的中文数据源进行合并,组成统一的知识库。共有621万实例,73万类别,上亿的事实。成果在ISWC以及SCI期刊上发表。

华东理工大学

上海交大

为了让知识图谱有用,准备构造一个搜索引擎?

谷歌知识图谱

"知心" "知立方 搜狗

目页 图片 10,95 前局 地图 亚多十 按索工具

批例的 40,300,000 奈结果 (用計 0.13 秒)

劉德華官方網站

www.andylau.com/ - 转为简体问页

1930-2014醛属世界高海... 天王刘婧华轻松穿出大长腿. 祝賀《華仔天地》成立26直年... 一花一世界。對應草陣突迅李進杰型點十大。 按 规键笔记67 生氣 價換品種

刘德华 百度百科

baike baidu.com/view/1758.htm .

刘德华1961年9月27日生于香港、演员、数手、调作人、制片人、影视歌多幅出居的代表 艺人之一,也是位有使命题的电影人。1981年以全优组缮毕业于TVE艺训班签约... 刘向惠、晋程、无线五贵、世界十大杰出青年

刘德华的新闻搜索结果



刘德华周杰伦比排盘点演唱会"失足"的大腕(图)

蘇华阿 4 小时期

此前有消息称,天王黎富城曾因清堪会彩排排倒导致跟键撕裂。其实 和郭天王一样,很多大黄明星都会因为场地原因或自身原因在清噶会 班场当中

刘德华周珂发娶衙门千金揭像上白富美的男星 人民网-4小时后

则德华家连杰或龙罗志祥携十大舆星的"干锅"

台海門 - 19 小村前

更多关于"刘德华"的新国

劉德華-維基百科,自由的百科全书

zh.wikipedia.org/zh-cn/對機華 *

劉德華、MH。JP(Andy Lau Takwah。1961年9月27日-)。香港著名演員兼款手。 1990年代获到香港乐坛「四大天王」之一。也是吉尼斯世界纪录大全中摄频最多的

刘德华 - 一听音乐网

www.fting.com/singer/65/singer_114.html =

刘德华全部歌曲、刘德华所有希腊、刘德华最新单曲、刘德华跟片、刘德华资料、刘德华2014年 好新的歌曲,歌手主页,一听音乐网,每天听一听

刘德华_腾讯娱乐。腾讯网

datalib ent gg.com/star/1333/ *

封至今日,刘谦华仍然是歌迁的一级巨星。他对工作孜孜不得。仍是为娱乐覆当红偶像。 司濟魅力无边。 生尚: 牛鞋子尺寸: 42号衣服尺寸: 中荷学历: 中学预料小学: 萬 ...

刘德华 Andy Lau - Mtime时光网 - 电影



刘德华

刘德华。MH。JP,香港署名演员兼数手。1990年代获封香港乐坛"四大天 王"之一,也是吉尼斯世界纪录大全中就装载多的香港歌手;电影方面他已 三股香酒电影全像奖载佳舆主角奖。并获得两座金马奖影响,截止2013年 参湾超过140部电影。刘德华现在是缺艺集团的老板。作为投资人参与制作 720多部华语电影。 维基森科

生于: 1961年9月27日(52岁),大埔

身高: 1.75米

专备: Global DSD Audio King . The Best of Andy Lau, 等等

子女: 汉纳拉岛

所赞奖项: 香港电影全像双悬任男主角,香港电影全像双悬任电影,香 滑电影全像奖载任国献角



用户还搜索了

无瑕酒









还有45+项

铁炮 2011年

还有15+顷



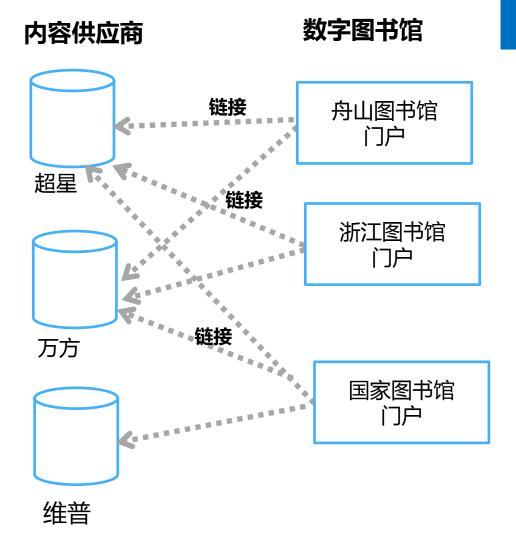


+ 为什么行业需要知识图谱

1.图书馆行业的故事

谁拥有更多的资源? 大图书馆!





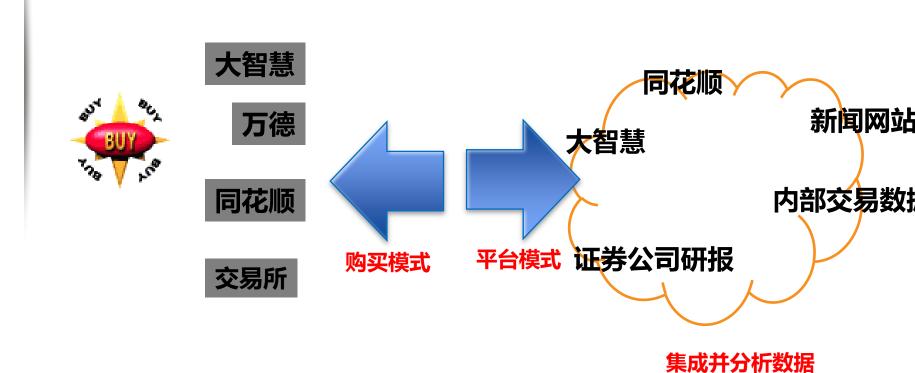
+ 地方图书馆利用知识图谱实现转型

严重的内容相似 与地方经济结合 缺乏内容控制 自有内容 专有技术 缺乏竞争力

寻觅新的机会!

+ 为什么行业需要知识图谱

2.证券行业的故事——购买数据VS自己处理数据





+ 为什么行业需要知识图谱

2.证券行业的故事——现有搜索引擎的困惑



为什么行业需要知识图谱

3. 医疗行业的例子——难点

- ■电子病历搜索
 - 某类患者,如心衰并患有高血 压患者的患者?
 - 与某病人相似的患者?



■ 电子病历文本中有大量的数据,如何进行结构化,以更好地进行电子病历数据的大数据挖掘?

- 不同来源知识库之间如何关联?
 - 疾病、药品、检查的关联
 - 中西医疾病名称关联
 - 中西药成分关联

* 为什么行业需要知识图谱

3. 医疗行业的例子——病历结构化的必要性

传统非结构化病历数据,只能通过文本匹配来进行查询

高血压 开始搜索◎所有◎现病史◎既往史◎家族史大肠癌早期便血开始搜索》所有◎现病史◎既往史◎家族史

病例查询结果		病例查询结果无
1.xxx,高血压患者	~	
2.最高血压160/90	×	
3.无高血压病史	×	

词汇二义性

无法理解相同文字但不同 含义的词汇

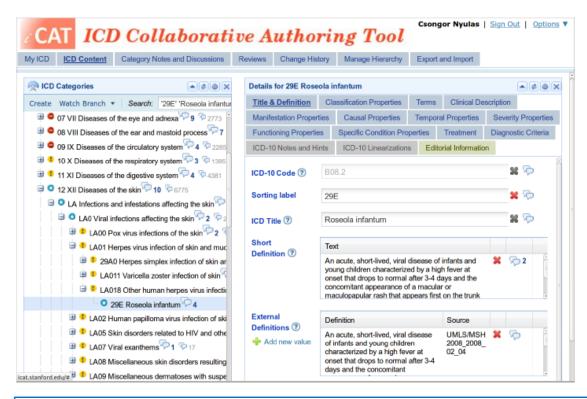
关联查询

无法精确切分查询词, 理解 查询意图

+

为什么行业需要知识图谱

3. 医疗行业的例子——ICD-11的构造



ICD 10以及以前版本,基本上没有结构。分类码通常是一长串的疾病与它们相关的代码,以及少量属性,如同义词等。

WHO在11版的ICD开发中使用语义Web 技术,支持协同编辑的语义Web 平台。 4年中,有270个来自世界各地的领域专家使用iCAT编辑了45,000个类,执行了260,000个变更。构造了17,000个链接,链到外部医学术语。

ICD 11使用了类/子类关系,子属性,定义域与值域,等价类。ICD 11表达能力是 SHOIN(D)。可以使用DL 推理程序去推理系统之间的非一致性。

+ 为什么行业需要知识图谱

总结

- 行业需要数据(语义)集成能力
 - 文本 与 结构化数据的集成
 - 不同来源、不同格式的大量数据
 - 自动/半自动的集成
- 行业需要(语义)数据查询能力
 - 更丰富的表达
 - 更精准的结果



十为什么行业需要知识图谱

知识图谱的技术优势

渐增式数据模式设计

初始设计的时候,很难 清楚所有的概念,而知 识图谱的动态可扩充性 以及"无模式"特性使 得用户很容易增加或修 改模式。

数据集成更轻松

本体的语义互操作特性 以及"链接数据"原则, 使得来自不同供应商的 数据集成更为方便。

现有标准支持

有RDF(S),OWL, SPARQL等标准,可 以逐渐要求内容供应 商支持。

语义搜索

用户可以查询具有某 类特征的某类实体 比起基于基于关键词 的搜索,更为精准。

主要内容

- ■为什么需要行业需要知识图谱
- ■垂直知识图谱的特点
- ■垂直知识图谱工具
- ■垂直知识图谱应用



- 1. 更为丰富与精确的领域数据
- 领域对数据质量要求更高,例如
 - 药品名称
 - 处方当中,药品的克数
 - 企业名称
 - 企业股票价格
- 领域数据字段与数据关联更丰富
 - 一个企业包含的信息字段可能有上百个字段
 - ICD11当中,每个疾病有56个属性,而其中52个是用填实例的。(不是随意填写一个文本,而是这个填的值本身也是RDF的一个实例。
- 只有满足了上述条件,才能用于商业分析与决策支持

VS 姚明的身高?

VS 电影的字段



2. 自顶向下

- 普通的KG,使用一种自底向上的方法,更强调数据的宽度。如 DBpedia在构造过程中,先有数据,后有本体。
- 对于行业本体来说,由于数据质量以及行业本身的规范要求,将使用自顶向下的方式。



- 3. 丰富的内部数据来源
- 企业/组织结构自有数据,如:
 - 对证券公司而言,自有用户交易数据
 - 医院,电子病历,付费、检查记录
 - 图书馆 地方政府数据
- ■行业数据
 - 如行业标准、规范, 如Medical Guideline
 - 如第三方收集的企业数据
- 诸多以RDB方式存储



- 4. 需要可扩充的第三方工具支持
- 比起互联网企业,垂直行业客户相对投入低,能力弱。
- ■需要可配置、图形化界面。
- 需要容易地面向不同行业做客户化。

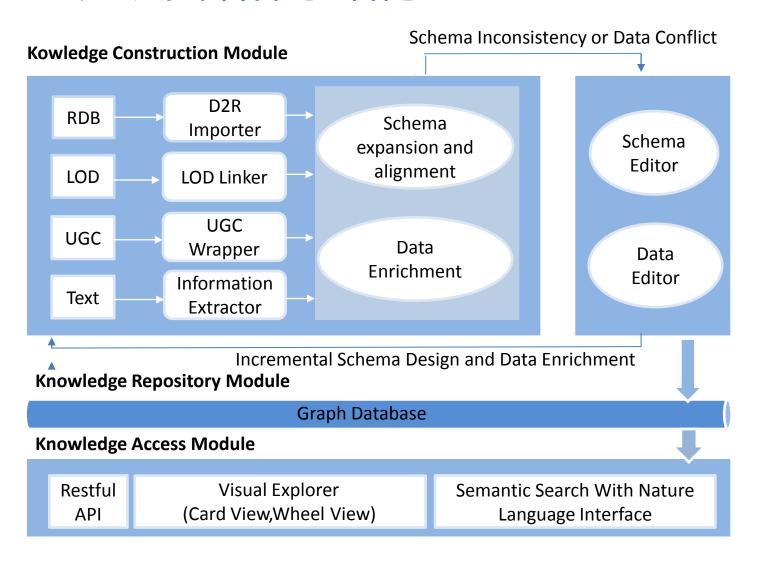


- 5. 与决策支持集成——语义搜索?Rule Engine?大数据挖掘?
- 医疗:如果我构造了医学知识图谱,下一步就是,如何基于这个图谱做诊疗?
- 证券:如何我构造了一个企业知识图谱,下一步就是,如何找到满足某一类条件的企业?

主要内容

- ■为什么需要行业需要知识图谱
- ■垂直知识图谱的特点
- ■垂直知识图谱工具
- ■垂直知识图谱应用

垂直知识图谱体系结构



+

构造海洋KG的例子——多种数据来源







+从Web抽取数据



+ 从百科网站抽取数据



十不同数据来源冲突解决



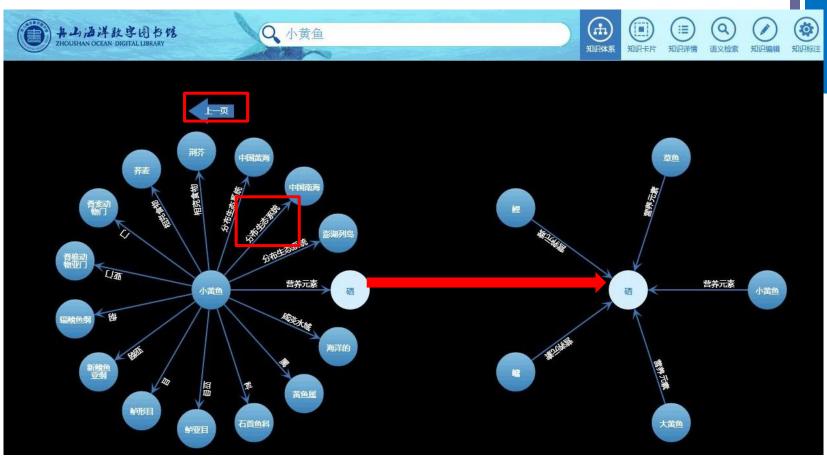
+ 语义检索



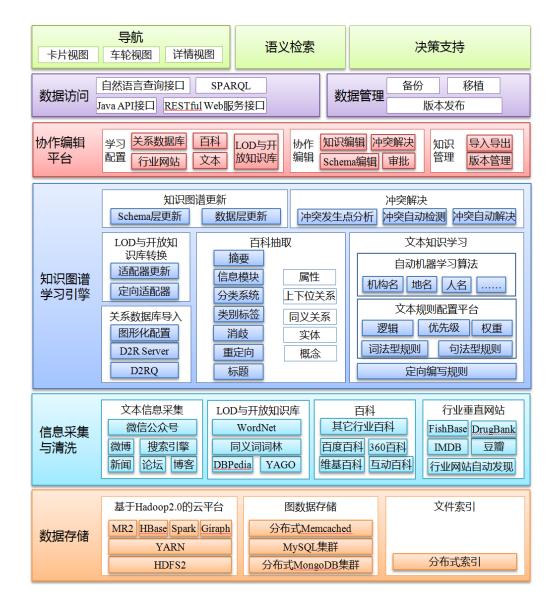
+ 浏览KG——卡片视图



+ 浏览KG——轮子视图



*未来产品框架

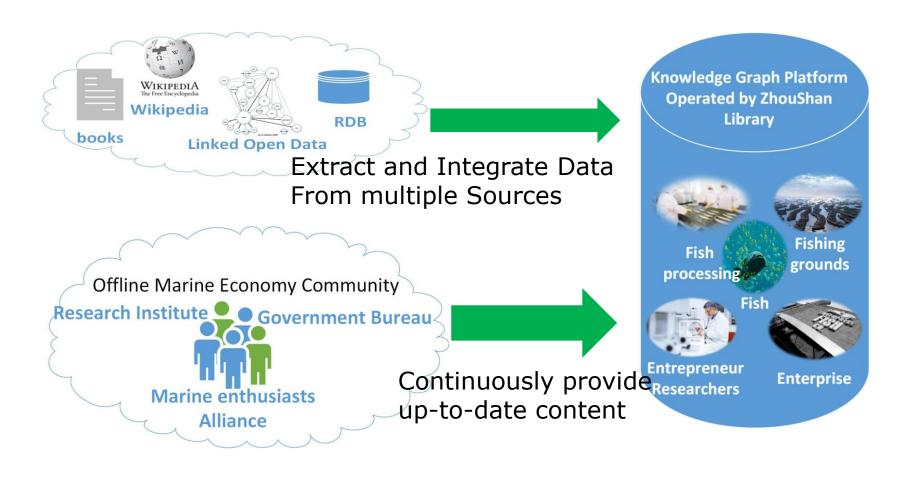


主要内容

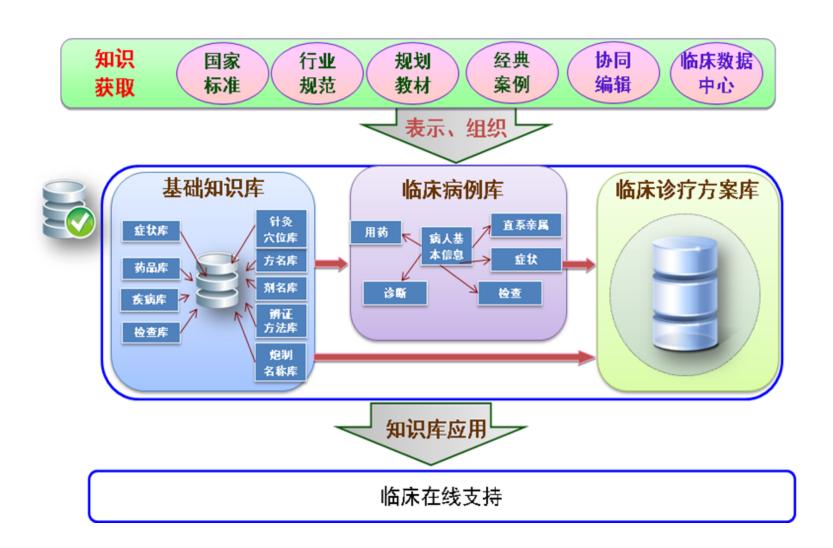
- ■为什么垂直行业需要知识图谱
- ■垂直知识图谱的特点
- ■垂直知识图谱工具
- ■垂直知识图谱应用

知识图谱帮助图书馆模式变迁

内容供应商+平台运营商



基于知识图谱技术的医疗知识库一一正在进行。。。



+ 总结

- 垂直知识图谱有巨大的前景
- 垂直知识图谱有诸多难点问题
- ■我们在图书馆、证券等行业做了部分探索





