

目录

目录

- 1.摘要
- 2.结果展示与知识图谱基本信息
- 3.实现思路与代码组织
- 4.代码细节
 - 4.1 GETDATA
 - 4.2 TI
- 5.创新点
- 6.实现规则
 - 6.1 rdf,rdfs规则
 - 6.2 部分owl规则
 - 6.3 基于频率推断subclassof和subpropertyof
 - 6.4 不一致性(得到Error时会输出错误信息)
 - 6.5 生成sameas用于同义消歧
- 6.小组每人工作量
- 7.引用

1.摘要

在本次课程设计中我们小组对金庸小说中人物进行知识图谱构建（本体构建，事实抽取，类别推断）。其中，类别推断部分我们构建了一个推理机并在此基础上加入了不一致性检测的功能。

在后面的章节中，[第2节](#)中描述了如何运行测试代码并提供了知识图谱的基本信息（例如规模）；[第3节](#)介绍了所提交代码的代码结构和各个模块的功能；[第4节](#)介绍了各个模块实现的具体细节[第5节](#)介绍了我们小组提出的创新想法；[第6节](#)描述了小组每人的工作量，[第7节](#)为引用部分

2.结果展示与知识图谱基本信息

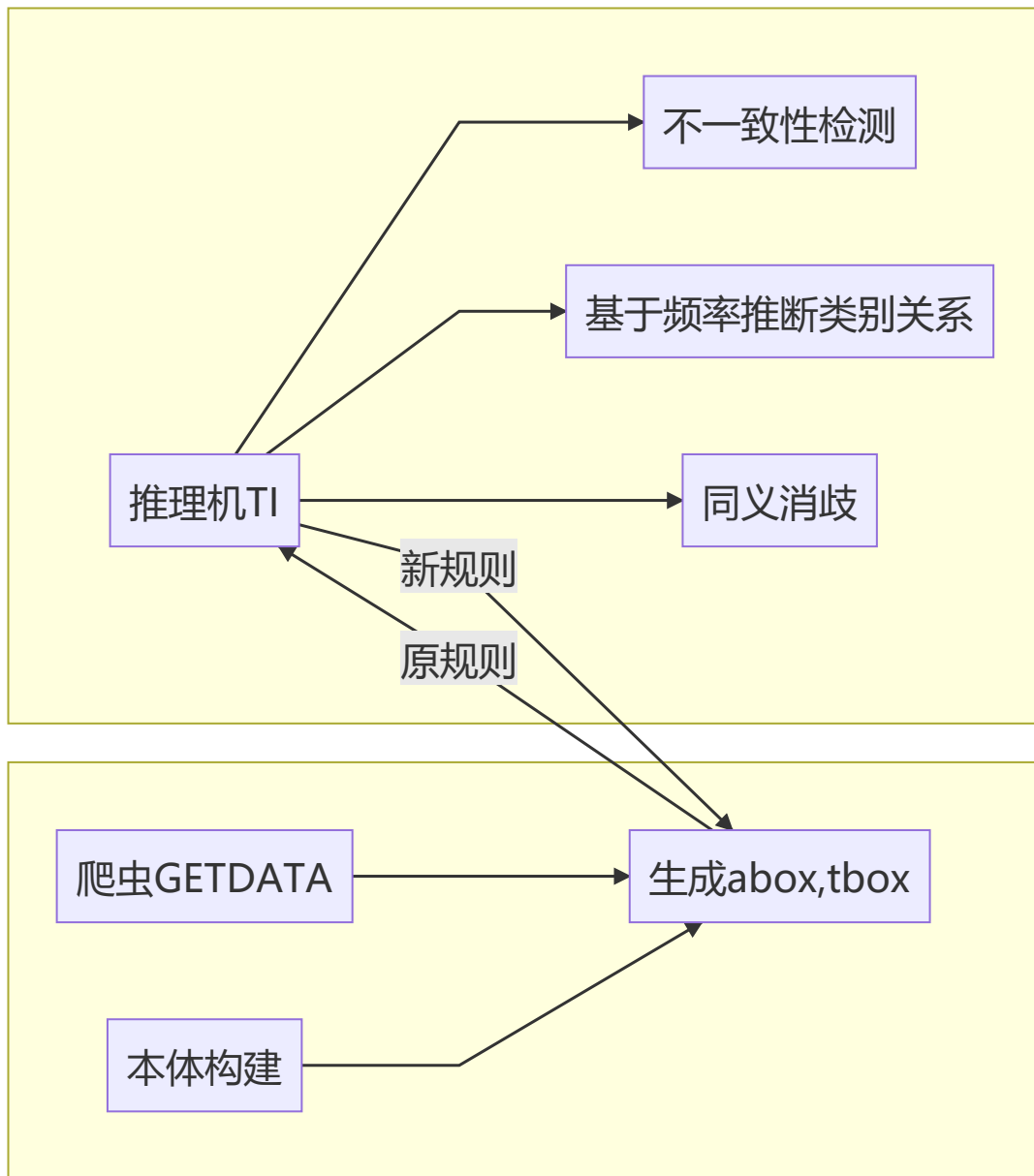
文件：result/金庸小说人物知识图谱.rdf

注明：result文件夹内其他文件为生成知识图谱的过程文件

规模：3508个实体。abox规则18366条；tbox规则6687条；节点数634个(金庸小说人物数量)

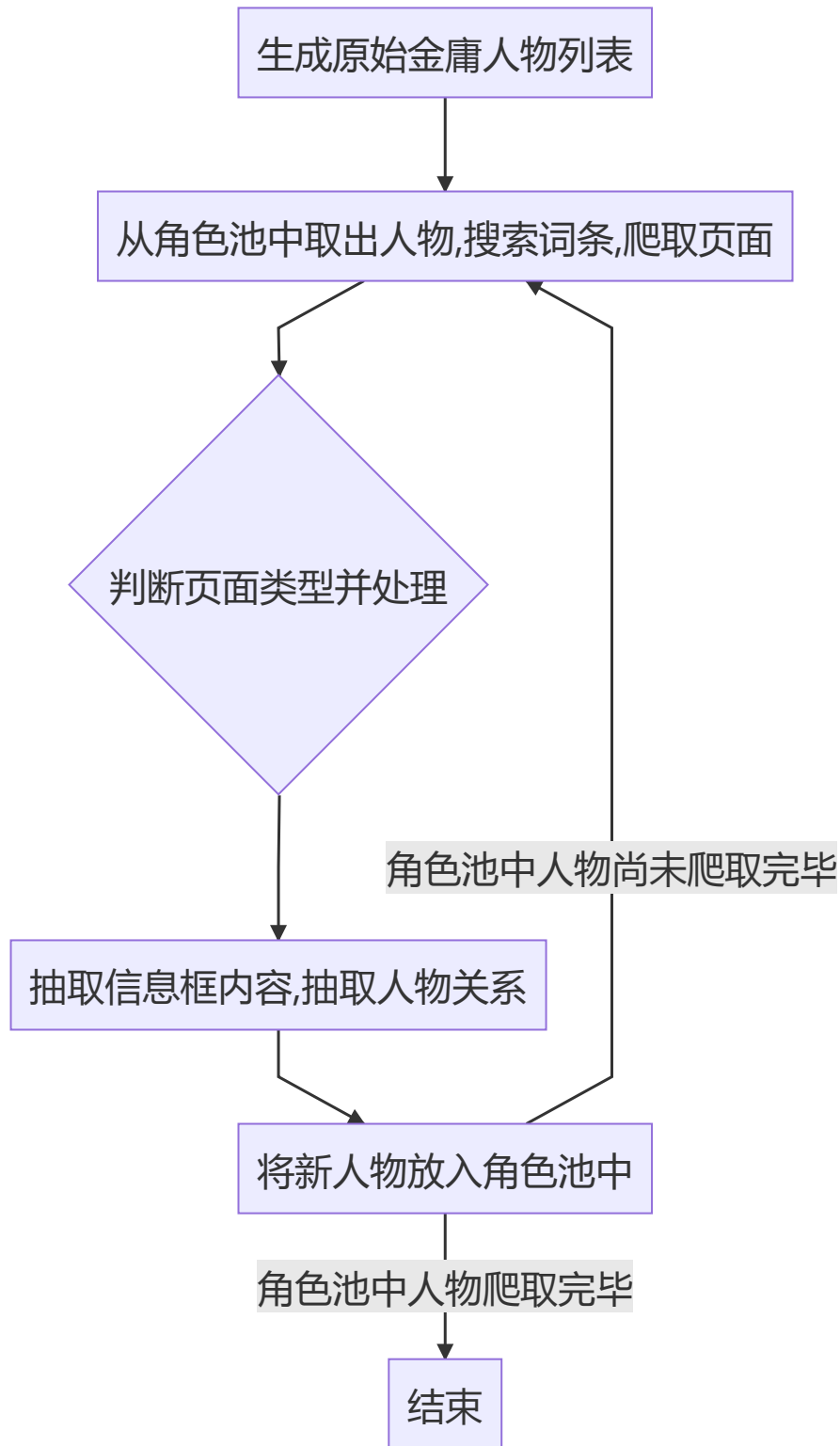
注明：知识图谱实际规模与答辩时所述不一致。规模大小的检查可以使用‘金庸小说人物知识图谱.rdf’，节点数检查可以使用‘abox_test_1.csv’

3.实现思路与代码组织

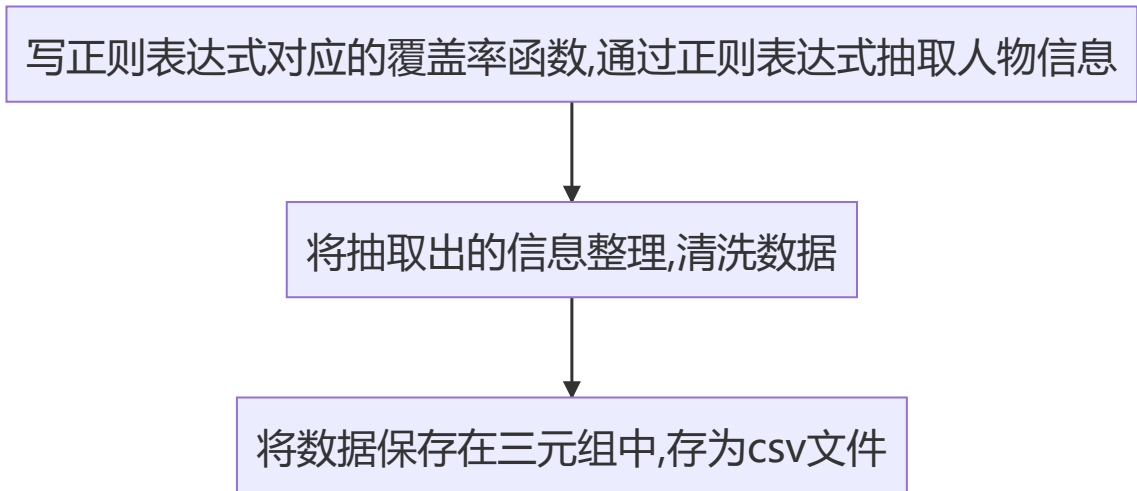


爬虫和事实抽取封装为库GETDATA;推理机封装为库TI;本节后面的内容首先介绍GETDATA的代码结构和执行流程，再介绍TI的代码结构和执行流程

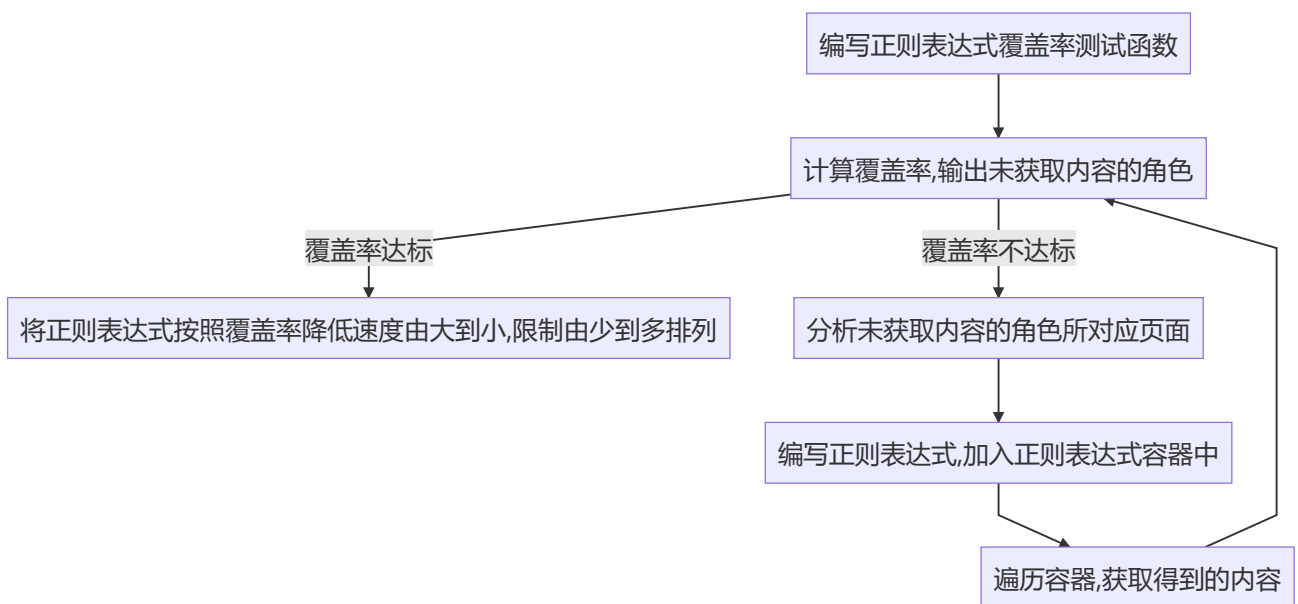
其中爬虫流程如下



事实抽取部分使用正则表达式抽取信息框内容,流程如下



写正则表达式抽取信息过程中,采用如下方法, 开发效率较高。



清理数据方面我们小组完成的工作如下

HTML转义

将分隔符化为多条语句

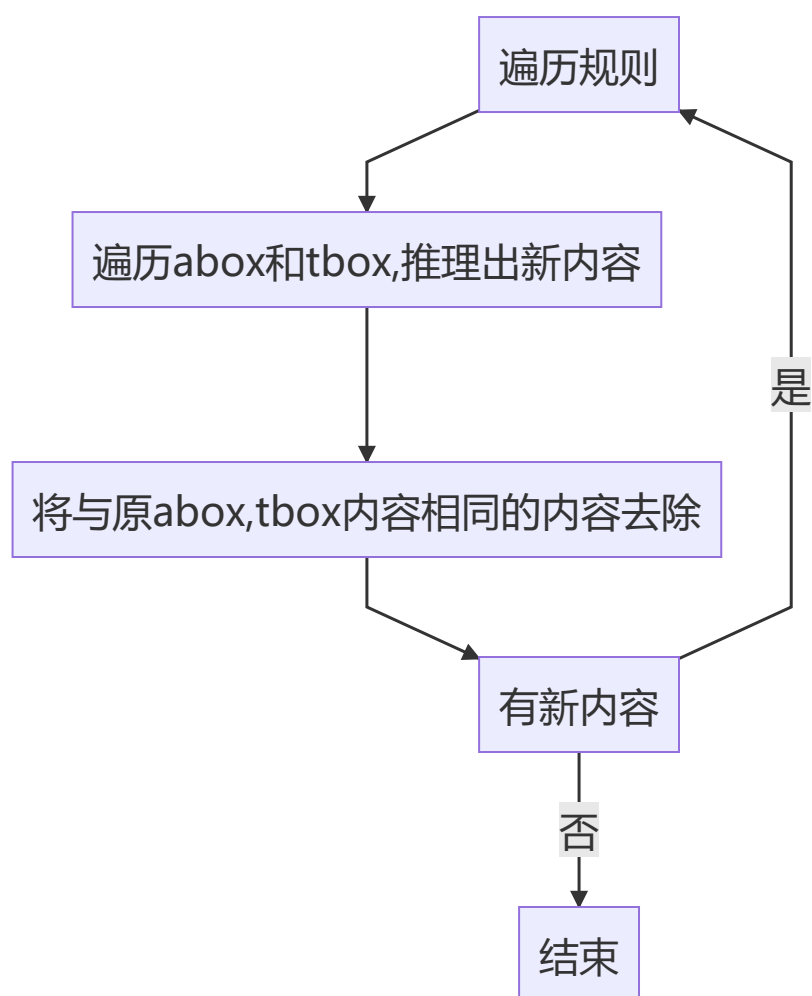
将内容后接括号的内容抽取同义词

删除无效字符

我们小组实现的推理机包含如下规则：（规则内容请见[第6节](#)）

- [rd和rdfs规则](#)
- [部分owl规则](#)
- [基于频率推断subclassof 和subproperty的规则](#)
- [不一致性检测规则](#)
- [生成sameas用于同义消歧](#)

推理机运行流程如下



时间复杂度：由于能够触发规则的事实稀疏，所以推理过程的时间复杂度没有原预想的那样高 n^{x+3} ,在可以接受的范围($\approx n^2$)

本节最后介绍不一致性检测和同义消歧实现过程。

1.不一致性检测使用的是6.4节所述规则。6.4节的规则是针对实例和类的，但是实际上也包含了类与类的不一致性检测，如：

$$\frac{A \perp B \wedge Csubclassof A \wedge Csubclassof B}{C = \emptyset} \quad (1)$$

2.同义消歧

2.1首先使用基于频率的方法和functional,inverse_functional的规则会推导出某些实例，类同义。

2.2利用同义信息尽可能扩大数据集，以将尽可能多地同义信息找出

2.3将所有同义词关联到一个实体上完成同义消歧 (TI.cleansame)

4.代码细节

第4节主要介绍 GETDATA 中 characterlist_base , LouisCha_characterlist; TI 部分的 rule_basic和typeinference的实现。

4.1GETDATA

由于我们小组在设计时采用了模板方法，将获取X信息的任务抽出共性部分，封装为基类，使得在抽取其他类似信息（最简单的如抽取水浒传人物信息，可以通过重载预设的函数来玩完完成诸如抽取公司高管信息的工作）

以下为基类代码。第3节中所提到的爬虫流程的代码由‘run’函数完成,run函数将会调用继承类实现的函数来完成对具体任务的流程

```
1 class characterlist_base():
2     '''
3     这个类用来继承，必须实现：
4     generate_url(self,args)
5     get_relative_character(self,args)
6     '''
7     def start(self):
8         self.pointer=0
9         self.namedict=dict()
10        self.namelist=[]
11        self.validnamelist=[]
12    def _init_namelist_dict(self,namelistroot):
13        '''
14        按照文本内容返回最原始的人物列表。
15        **最后的去重操作会将顺序打乱，所以不可重复操作**
16        :param namelistroot: 储存人物列表的json文件地址
17        json文件要求：键存储书名，值存储列表，列表内为人名，人名以字符串给出
18        :return:无 #TODO:True为读取文件成功，False为失败
19        '''
20        with open(namelistroot,'r',encoding='utf-8')as f:
21            self.namedict=eval(f.read())
22            for value in self.namedict.values():
23                self.namelist.extend(value)
24            self.namelist=list(set(self.namelist))
25            print('将{} 作为人物列表导入完成'.format(namelistroot))
26    def init_all(self,namelistroot):
27        self.pointer=0
28        self._init_namelist_dict(namelistroot)
29        self.validnamelist = []
30    def generate_url(self,args=dict()):
```

```

31         raise NotImplementedError()
32     def findbookname(self,charactername):
33         for key,value in self.namedict.items():
34             if charactername in value:
35                 return key
36         return ''
37     def deal_no_jump(self,args=dict):
38         raise NotImplementedError()
39     def get_relative_character(self,args=dict()):
40         raise NotImplementedError()
41     def cleannames(self):
42         raise NotImplementedError()
43     def run(self,args=dict):
44         '''
45         模板方法,获得所有人物的百科html文件
46         saveroot,host_url,namelistroot
47         :param args:
48         :return:
49         '''
50         self.start()
51         self.init_all(args['namelistroot'])
52         while True:
53             if self.pointer>=len(self.namelist):
54                 break
55             else:
56                 #输出前先清空输出
57                 os.system('cls')
58                 print('已经爬完{},当前进度<=
{: .2%}'.format(self.pointer,self.pointer/len(self.namelist)))
59                 url=self.generate_url(args)
60                 rawdata=get_savehtml(url,args['saveroot'],self.namelist[self.pointer])
61                 if rawdata!="":#爬取出了结果
62                     rawdata=self.deal_no_jump(args={
63                         'data':rawdata,
64                         'charactername':self.namelist[self.pointer],
65                         'saveroot':args['saveroot'],
66                         'host_url':args['host_url']
67                     })
68                     if rawdata== -1:
69                         self.pointer+=1
70                         continue
71                     self.validnamelist.append(self.namelist[self.pointer])#找到哪些人名有词条
72                     names=self.get_relative_character({'data':rawdata})
73                     names=self.cleannames({'dirty_characters':names})
74                     for name in names:
75                         if name not in self.namelist:
76                             self.namelist.append(name)
77                     self.pointer += 1
78                     time.sleep(rand.randint(1,10)/10)
79     def __str__(self):
80         return str(
81             dict(
82                 {
83                     'pointer':self.pointer,
84                     'namelist':self.namelist,
85                     'namedict':self.namedict,
86                     'validnamelist':self.validnamelist
87                 }
88             )
89         )
90

```

以下为继承类代码，用于实现本次课设的“金庸小说人物”。

```

1 class LouisCha_characterlist(characterlist_base):
2     def generate_url(self, args):
3         return
f'https://baike.baidu.com/item/{tranchinise(self.namelist[self.pointer])}'

```

```

4     def get_relative_character(self, args):
5         '''#获取人名（脏数据）
6         df = pd.read_html(args['data'])[0]
7         return list(df['人名'])'''#不够直接
8         pat='''<li class="lemma-relation-item">.*?<span class="title">(.*?)</span>.*?
</li>'''#抽人物关系
9         return re.compile(pat,re.S).findall(args['data'])
10    def is_clean(self,string):
11        return string.find(':')== -1 and string.find('(')== -1 and\
12        string.find(')')== -1 and string.find(' ')==-1 and\
13        string.find(',')==-1
14    def cleannames(self,args:dict):
15        # TODO:清理数据, 样例: ['父亲: 黄药师 母亲: 冯氏（小字阿衡）', '郭靖', '郭芙、郭襄、郭破
        虏', '洪七公', '郭啸天', '李萍', '杨过', '耶律齐', '穆念慈、周伯通', '双雕、小红马、血鸟（连载
        版）']
16        dirty_characters=args['dirty_characters']
17        clean_characters=[]
18        for dirty_character in dirty_characters:
19            #思路是把连续的中文挑出来, 如果挑出非人名也无妨, 因为deal_no_jump会将其去掉
20            tran=re.compile('([\u4e00-\u9fa5]+)').findall(dirty_character)
21            clean_characters.extend(tran)
22        return clean_characters
23    def _nobookappearance(self,args):
24        for elem in self.namedict.keys():
25            if args['data'].find(elem)!=-1:
26                return False
27        return True
28    def deal_no_jump(self,args:dict):
29        '''
30            这部分较为困难。
31            没有找到词条也在这里处理
32            :param args:
33                data
34                charactername
35                saveroot
36                host_url
37            :return:
38                '''
39        if args['data'].find('金庸')== -1 and self._nobookappearance(args): # 如‘冯
        氏’, 搜出来不是金庸小说人物
40            '''
41            此处不能只写含"金庸", 如"谭婆: 小说《天龙八部》中的人物"
42            '''
43            return -1
44        elif args['data'].find('keywords')== -1:
45            '''
46            #找字符串, -1为对象字符串不含参数字符串, 其他值为索引
47            #百科搜索可能页面跳转, 找规律发现未跳转的页面没'keyword'字符串
48            '''
49            # TODO:通过返回码判断是否跳转更安全
50            '''
51            找规律发现金庸人物的目录下含书名,用人物找书名
52            '''
53            if args['data'].find('百度百科尚未收录词条') != -1:
54                return -1
55
56            else:
57                newurl = ''
58                polysemous_words = strongpat(args['data'], {
59                    'for_polysemous': '<div class="para" label-module="para"><a(.*)?
</a>[\s]*</div>' # 识别多义词,把链接也拉过来
60                }, re.S)
61
62                for elem in polysemous_words[0]:
63                    if elem.find('金庸') != -1 or elem.find(
64                        self.findbookname(args['charactername'])) != -1: # 认为结果
        含有金庸 或是书名就是索要查找结果
65                    tran = re.compile('href="(.*?)")').findall(elem) # 找链接
66                    newurl = args['host_url'] + tran[0][1:]

```



```

67         return get_savehtml(newurl, args['saveroot'],
args['charactername'])
68
69     return args['data']
70     def _getinfobox(self, args: dict):
71         '''
72         百度百科的html文档中infobox有多种格式,需要多种正则表达式进行提取。本函数提供了一个比较有效的
方法。
73         这一部分是这样完成的:
74         1. 首先创建测试函数,测试函数会将正则表达式提取为空的文档输出,并输出提取失败的文档个数
75         2. 根据提取失败的文档写出新的正则表达式添加在pats里
76         3. 本函数会循环迭代pats,使得文档尽可能多地匹配
77         4. 注意将较偏的提取规则放在列表的后面,以防其他html使用较偏的提取规则提取,如公孙绿萼.html的
infobox不够规范,只能使用如pats[4]的提取规则提取
78         :param args:
79         :return:
80         '''
81         pats=['<dt class="basicInfo-item name">(.*?)</dt>[\s]*<dd class="basicInfo-
item value">[\s]*(.*?)</dd>''',
82              '<li class="basicInfo-hide">[\s]*<div class="info-title">(.*?)</div>
[\s]*<div class="info-content".*?>(.*?)</div>[\s]*</li>''',
83              '<li>[\s]*<div class="info-title">(.*?)</div>[\s]*<div class="info-
content".*?>(.*?)</div>[\s]*</li>''',
84              '<p>[\s]*<strong>(.*?)</strong>(.*?)</p>''',
85              '\n([\u2E80-\u9FFF]*?): (.*?)<br />''']
86
87         for pat in pats:
88             tran=re.compile(pat,re.S).findall(args['data'])
89             if tran!=[]:
90                 rst=tran[:]
91                 return [(args['name'],rst[i][0],rst[i][1]) for i in range(len(rst))]
92     def test_getinfobox(self, args, key):
93         count=0
94         for elem in self.getzip(args):
95             if self._getinfobox(elem)==None:
96                 count+=1
97                 print(elem['name'])
98         print(count)
99     def getzip(self, args: dict):
100         for elem in os.listdir(args['rawdataroot']):
101             yield {'name':elem[:-5],
102                   'data':open(args['rawdataroot']+elem, 'r', encoding='utf-8').read()}
103     def get_triple(self, args: dict):
104         import zipfile
105         zipfile_path = ''
106         trandict=geti2c()
107         with zipfile.ZipFile(args['rawdataroot'], mode='r') as zfile: # 只读方式打开压缩
包
108             for name in zfile.namelist(): # 获取zip文档内所有文件的名称列表
109                 with zfile.open(name, mode='r') as single_file:
110                     yield self._getinfobox(args=
{'name':trandict[name[:-5]], 'data':single_file.read()})
111     def getinfobox(self, args):
112         target=open('./test.csv', 'a', encoding='utf-8')
113         writer=csv.writer(target)
114         for newargs in self.getzip(args):
115             tran=self._getinfobox(args=newargs)
116             if tran:
117                 for triple in tran:
118                     writer.writerow(triple)
119         target.close()

```

4.2 TI

推理机TI整体架构为由typeinference所提供的功能循环调用不同的规则类(rule)，所有的规则继承自Rule_basic统一封装接口。

以下为Rule_basic具体实现

```
1 class Rule_basic():
2     def _run(self, args):
3         '''
4         接口说明:
5         0.类命名规则: rule数字, 数字对应ppt中rdf后面数字
6         派生类请单独写在一个文件中, 本文件下的类作为例子。
7         1.在本函数中需要完成的工作: 根据给定三元组组成的列表 (abox, tbox), 按照当前函数所执行的规则编
            写输出两个新三元组组成的列表的接口。两个新三
8         元组分别为abox和tbox。新三元组中不包含原三元组列表的内容。不需要去重。
9         2.本函数执行的操作是按照当前规则, 将生成的三元组全部输出
10        :param args: [
11        参数规定:
12        'tbox':tbox,
13        'abox':abox
14        可以按照 rule数字_参数名 添加新键, 命名规则保证不重名
15        ]
16        trule[0]:头, trule[1]:边, trule[2]:尾
17        arule[0]:头, arule[1]:边, arule[2]:尾
18
19        迭代器命名建议: trule对应tbox的迭代器, arule对应abox的迭代器。
20
21        :return:返回值建议命名为rst_abox,rst_tbox
22
23        type 省略前缀
24        3.TODO:本函数所完成的逻辑规则(需填写)
25        '''
26        raise NotImplementedError()
27
28    def run(self, args):
29        raw_rst_abox, raw_rst_tbox = self._run(args)
30        rst_abox, rst_tbox = [], []
31        for elem in raw_rst_abox:
32            if elem not in args['abox']:
33                rst_abox.append(elem)
34        for elem in raw_rst_tbox:
35            if elem not in args['tbox']:
36                rst_tbox.append(elem)
37        return rst_abox, rst_tbox
```

以下为typeinference的具体实现

```
1 class typeInference():
2     def start(self, rules):
3         self.IE=rules
4         pass
5     def retranslate(self):
6         pass
7     def TI(self, args):
8         rst_abox, rst_tbox=args['abox'], args['tbox']
9         counter=1
10        import time
11        timer=[]
12        curtime=time.time()
13        while(True):
14            new_abox, new_tbox=[], []
15            for rule in self.IE :
16                tran=rule.run(args={
```

```

17         'abox':rst_abox,
18         'tbox':rst_tbox
19     })
20     new_abox.extend(tran[0])
21     new_tbox.extend(tran[1])
22
23     if len(new_abox)==0 and len(new_tbox)==0:
24         break
25     rst_abox.extend(new_abox)
26     rst_tbox.extend(new_tbox)
27     #外层去重
28     rst_abox=list(set([tuple(elem) for elem in rst_abox]))
29     rst_tbox = list(set([tuple(elem) for elem in rst_tbox]))
30     tran=time.time()-curtime
31     timer.append(tran)
32     print('第{}次迭代,消耗时间:{:.2f}s'.format(counter,tran))
33     counter+=1
34     print('总耗时:{:.2f}s'.format(sum(timer)))
35     return rst_abox, rst_tbox
36

```

不一致性检测代码如下

```

1
2 class rule_check_1(Rule_check_basic):
3     '''对于类A和类B,若A和B不相交且存在a使得a属于A且a属于B,
4     则返回一条标识错误的三元组('序列号','DISJOINTERROR','ainA_ainB')
5     rule_check规则禁止改动规则
6     '''
7     def _run_tran(self):
8         rst_assertion_box= []
9         rst_original_box=self.rst_original_box[:]
10        #找所有含type的语句
11        tran1,tran2=[],[]
12        for elem in rst_original_box:
13            if elem[1]==TI_TYPE:
14                tran1.append(elem)
15            if elem[1]==TI_DISJOINT:
16                tran2.append(elem)
17        #分析type语句
18        for elem1 in tran2:
19            for elem2 in tran1:
20                if elem2[2]==elem1[0]:
21                    for elem3 in tran1:
22                        if elem3[0]==elem2[0] and elem3[2]==elem1[2]:
23                            errorsrcript=f'{elem3[0]} in {elem1[0]} and {elem1[2]},
{elem1[0]} disjointwith {elem1[2]}'
24                            rst_assertion_box.extend(
25                                [
26                                    (elem1[3],TI_DISJOINT_ERROR,errorsrcript),
27                                    (elem2[3], TI_DISJOINT_ERROR, errorsrcript),
28                                    (elem3[3], TI_DISJOINT_ERROR, errorsrcript),
29                                ]
30                            )
31        return rst_original_box,rst_assertion_box

```

同义消歧代码如下:

```

1 def cleansame(abox,tbox):#完成同义消歧
2     for i,arule in enumerate(abox):
3         for j,elem1 in enumerate(arule):
4             if elem1==TI_SAMEAS:
5                 for i1,arule1 in enumerate(abox):
6                     for j1, elem11 in enumerate(arule1):
7                         if elem11 == arule[0]:
8                             arule1[j1]=arule[2]

```

```

9         for i1, trule1 in enumerate(tbox):
10             for j1, elem11 in enumerate(trule1):
11                 if elem11 == arule[0]:
12                     trule1[j1] = arule[2]
13     #删除重复词，原来的sameas用于为实体添加别名
14     rst_abox=list(set(abox))
15     rst_tbox = list(set(tbox))
16     return rst_abox,rst_tbox

```

5.创新点

任务：使用对抗文本抽取属性的值（属性值的合并）

在文本中抽取属性值时经常会出现某个属性的值有多种表达方式(质量:高 可以表达成‘优秀’、‘高质量’)，现在需要将这些属性值找出或整合（消歧）。

给定一个至少是闭源的分类模型(可以获取模型的输入输出，但不知道模型的具体结构)。

算法：

这个方法的思想是认为经过对抗训练，且有较好分类效果的分类器可以逼近人的分类标准。构建可以区分属性值的分类器后，再通过计算后验概率或梯度的方法在分类器起较大作用的单词，将此作为属性的同义词。

1. 选取多个经过对抗训练的分类器进行如下操作，每个分类器分类的是某个属性的值（如每个分类器可以分类文本中表达的产品质量信息，分类为‘好’，‘中’，‘差’）
2. 计算分类器中占较大贡献度的单词
3. 将3中得到的单词放入一个属性值的候选池中，计数
4. 单词对应的候选词投票，得到与属性值同义的词

其中步骤2.参考了[1],计算每个词在分类中所占的贡献度

这个方法是使用如下函数

$$C_F(a_i, y) = \begin{cases} p_F(y_i|s) - p_F(y_i|s \text{ without } w_k) & F(s) = y_i = F(s \text{ without } w_k) \\ p_F(y_i|s) + p_F(y_i|s \text{ without } w_k) & F(s) = y_i \neq F(s \text{ without } w_k) \end{cases} \quad (2)$$

其中 P_F 是分类器F所认为的某文本s(或s without w_k)为 y_i 类别的概率。将分类器看作一个以文本为输入类别为输出的函数F,通过上式可以计算每个词的贡献度。

在模型完全已知(知道模型损失函数)的情况下，可以使用类似FGSM的方法来计算贡献度。

$$C_F(w_k, y) = -\nabla_{w_k} J(F, s, y_i) \quad (3)$$

其中J为模型的损失函数。

6.实现规则

本次课程设计所实现的规则:

6.1rdf,rdfs规则

$$\frac{xRy}{R \text{ type Property}} \quad (4)$$

$$\frac{R \text{ domain } C_1 \wedge xRy}{x \text{ type } C_1} \quad (5)$$

$$\frac{R \text{ range } C_1 \wedge xRy}{y \text{ type } C_1} \quad (6)$$

$$\frac{xRy}{x \text{ type Resource} \wedge y \text{ type Resource}} \quad (7)$$

$$\frac{R_1 \text{ subpropertyof } R_2 \wedge R_2 \text{ subpropertyof } R_3}{R_1 \text{ subpropertyof } R_3} \quad (8)$$

$$\frac{R \text{ type Property}}{R \text{ subpropertyof } R} \quad (9)$$

$$\frac{R_1 \text{ subpropertyof } R_2 \wedge xR_1y}{xR_2y} \quad \frac{R_1 \text{ subpropertyof } R_2 \wedge xR_1y}{xR_2y} \quad (10)$$

$$\frac{c \text{ type Class}}{c \text{ subclassof Resource}} \quad (11)$$

$$\frac{c_1 \text{ subclassof } c_2 \wedge x \text{ type } c_1}{x \text{ type } c_2} \quad (12)$$

$$\frac{c \text{ type Class}}{c \text{ subclassof } c} \quad (13)$$

$$\frac{c_1 \text{ subclassof } c_2 \wedge c_2 \text{ subclassof } c_3}{c_1 \text{ subclassof } c_3} \quad (14)$$

$$\frac{x \text{ type Datatype}}{x \text{ subclassof Literal}} \quad (15)$$

6.2部分owl规则

$$\frac{R \text{ type Transitive} \wedge xRy \wedge yRz}{xRz} \quad (16)$$

$$\frac{R \text{ type Symmetric} \wedge xRy}{yRx} \quad (17)$$

$$\frac{R \text{ type Functional} \wedge xRy_1 \wedge xRy_2}{y_1 \text{ sameas } y_2} \quad (18)$$

$$\frac{x_1 \text{ sameas } x_2 \wedge f(x_1)}{f(x_2)} \quad (19)$$

$$\frac{R_1 \text{ inverseof } R_2 \wedge xR_1y}{yR_2x} \quad (20)$$

$$\frac{R \text{ type } InverseFunctional \wedge x_1 Ry \wedge x_2 Ry}{x_1 \text{ sameas } x_2} \quad (21)$$

6.3 基于频率推断subclassof和subpropertyof

$$\begin{aligned} R_i &= \{ \langle x_{i1}, y_{i1} \rangle, \dots, \langle x_{in}, y_{in} \rangle \} \\ R_j &= \{ \langle x_{j1}, y_{j1} \rangle, \dots, \langle x_{jn}, y_{jn} \rangle \} \\ \frac{\text{len}(R_i \cap R_j) \div \text{len}(R_i) > \alpha}{R_i \subseteq R_j} \end{aligned} \quad (22)$$

6.4 不一致性(得到Error同时会输出错误信息)

$$\frac{A \perp B \wedge a \text{ type } A \wedge a \text{ type } B}{Error} \quad (23)$$

6.5 生成sameas用于同义消歧

$$\frac{c_1 \text{ subclassof } c_2 \wedge c_2 \text{ subclassof } c_1}{c_1 \text{ sameas } c_2} \quad (24)$$

$$\frac{p_1 \text{ subpropertyof } p_2 \wedge p_2 \text{ subpropertyof } p_1}{p_1 \text{ sameas } p_2} \quad (25)$$

6. 小组每人工作量

略

7. 引用

[1] Samanta, S., & Mehta, S. (2017). *Towards Crafting Text Adversarial Samples*.