

远程科研指导项目总结报告

聊天机器人项目总结

姓 名： 彭 友

学 校： 上海理工大学

专 业： 计算机科学与技术

2019 年 6 月

目录

1. 项目概述.....	1
2. 背景知识.....	1
2.1 自然语言处理	1
2.2 Scrapy.....	1
2.3 Rasa	2
3. 过程分析.....	3
3.1 收集数据	3
3.2 合并数据	4
3.3 Neo4j 知识图谱设计	4
3.4 Rasa 的配置	6
3.5 前端	7
4. 结果展示.....	8
4.1 可用功能	8
4.2 知识图谱	9
4.3 前端	10
5. 感想与收获.....	11
6. 参考.....	12

1. 项目概述

本项目通过使用名为 Scrapy 的爬虫框架、名为 Rasa 的机器对话框架和 Neo4j 图数据库，构造了一个中文的医疗领域问答系统，以便人们可以通过聊天的方式快速、精准地获取到指定疾病的介绍、治疗方法、用药方案等信息，系统将具有一定的理解人类语言的能力，具有可重用性，并且与知识图谱建立联系，根据问题动态地获取数据、构造答案、返回答案。

2. 背景知识

2.1 自然语言处理

自然语言处理（Natural Language Processing）是人工智能领域中的一个研究方向，它被用于处理人类的语言。语言是人类区别其他动物的本质特性。在所有生物中，只有人类才具有语言能力。人类的多种智能都与语言有着密切的关系。人类的逻辑思维以语言为形式，人类的绝大部分知识也是以语言文字的形式记载和流传下来的。因而，NLP 是人工智能的一个重要，甚至核心的部分。语言具有模糊、上下文相关和隐晦的特性，甚至有时并不准确，同时人类的语言是基于人类认知的，而在表达时，同一个意思还可以有多种不同的表达方式。所以让机器处理人类的语言是非常困难的，而自然语言处理技术就旨在解决相关的问题。

2.2 Scrapy

Scrapy 是一套基于 Twisted 的异步处理框架，是纯 Python 实现的爬虫框架，用户只需要定制开发几个模块就可以轻松实现一个爬虫，用来抓取网页中的内容、图片等。通常单次爬取行为可分如下步骤：引擎（Engine）的作用是控制事件的触发和数据在组件中的流动，引擎根据 URL 构造请求（Request）对象传递给调度器（Scheduler），调度器将请求对象进行排队，引擎需要时，调度器就把请求数据给引擎，由引擎交给下载器（Downloader）并由下载器获取对应的响应数据（Response）交还给引擎，引擎把响应数据交给爬虫（Spider）进行解析提取，提取到的数据产生出对应的项目（Item）交给管道（Item Pipeline）进行后期过滤或者存储等。

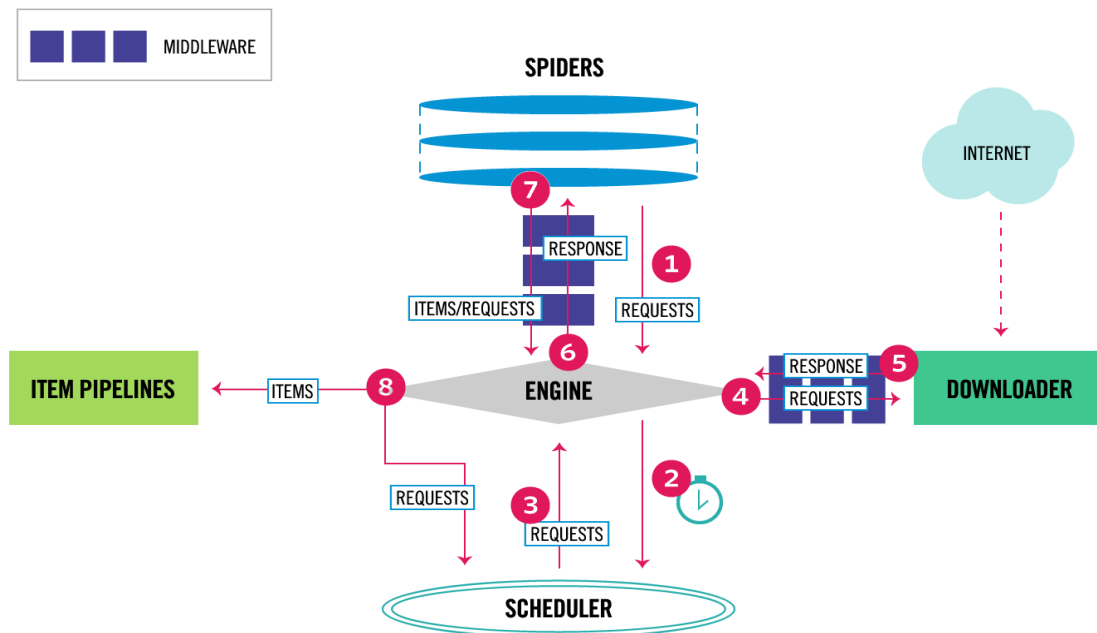


图 1 Scrapy 架构

2.3 Rasa

Rasa 是一个用于自动文本和基于语音的对话的开源机器学习框架，支持多种语言，内置支持 Spacy、MITIE、Jieba 等多种开源 NLP 工具，可将各种组件进行组合配置，构造定制化的 NLP Pipeline 以便适合需求。

Rasa 配套了 Rasa NLU 和 Rasa Core，Rasa NLU 用于实体提取和意图识别；Rasa Core 是会话引擎，它用真实的对话文本和机器学习算法训练得到模型，从而预测在收到消息后机器人该作何回复（动作）。

Rasa 有封装好的服务器 Server 程序，运行 Server 后，可以使用 Http API 来调用功能。

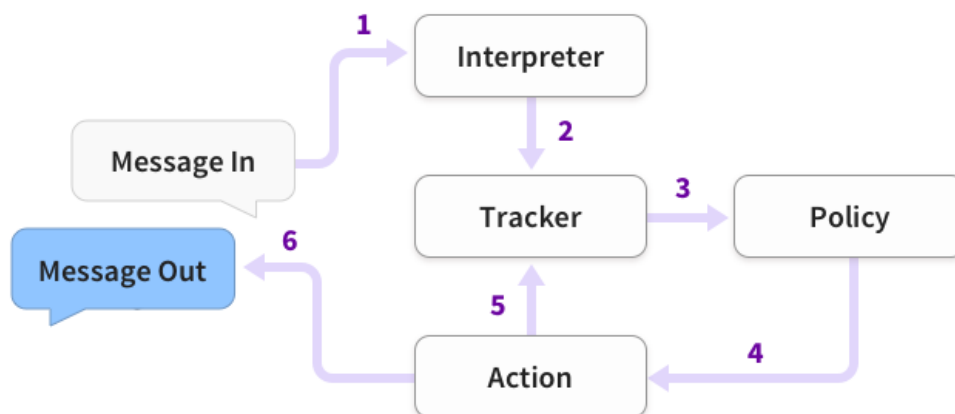


图 2 Rasa 架构机器人消息响应流程

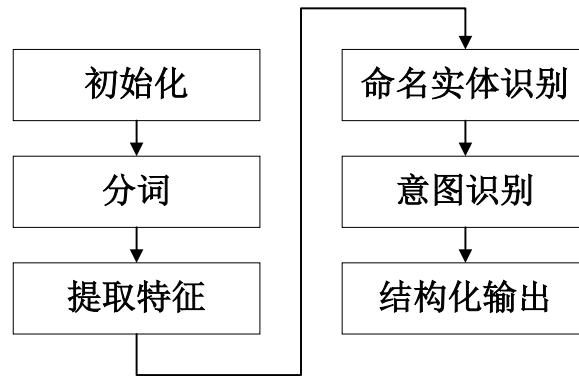


图 3 Rasa NLU Pipeline 流程

3. 过程分析

3.1 收集数据

考虑到目标是要做一个中文的医疗方面的问答系统，所以需要去搜集疾病、药物等有关的数据，但是搜寻一番后发现并没有符合需求的开放数据集，所以使用 Scrapy 库来编写爬虫程序爬取了寻医问药网。采用了随机 User-Agent、随机延迟请求和限制请求并发数等机制防止爬虫被检测，同时利用了 2 台 PC 和 1 台云服务器分别在不同网络下同时分块爬取。

```

Run: test
Now agent---->Mozilla/5.0 (Windows; U; Windows NT 5.1; zh-CN; rv:1.9) Gecko/20080705 Firefox/3.0 Kapiko/3.0
2019-07-15 10:54:27 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://jib.xywy.com/il_sii/nursing/43.htm> (referer: http://jib
2019-07-15 10:54:27 [root] DEBUG: ### random delay: 6.42 s ###
Now agent---->Mozilla/4.0 (compatible; MSIE 6.0; Windows NT 5.1; SV1; QQDownload 732; .NET4.0C; .NET4.0E)
2019-07-15 10:54:34 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://jib.xywy.com/il_sii/prevent/46.htm> (referer: http://jib
2019-07-15 10:54:34 [root] DEBUG: ### random delay: 10.700000000000001 s ###
Now agent---->Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/536.11 (KHTML, like Gecko) Chrome/20.0.1132.11 TaoBrowser/2.0 Safari/
2019-07-15 10:54:45 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://jib.xywy.com/il_sii/nursing/44.htm> (referer: http://jib
2019-07-15 10:54:45 [root] DEBUG: ### random delay: 10.700000000000001 s ###
Now agent---->Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.1 (KHTML, like Gecko) Chrome/21.0.1180.89 Safari/537.1
2019-07-15 10:54:56 [scrapy.core.engine] DEBUG: Crawled (200) <GET http://jib.xywy.com/il_sii/food/43.htm> (referer: http://jib
2019-07-15 10:54:56 [root] DEBUG: ### random delay: 8.56 s ###
Now agent---->Mozilla/5.0 (Windows NT 6.1; WOW64) AppleWebKit/537.1 (KHTML, like Gecko) Chrome/21.0.1180.71 Safari/537.1 LBBROWSER
2019-07-15 10:55:04 [scrapy.extensions.logstats] INFO: Crawled 18 pages (at 7 pages/min), scraped 0 items (at 0 items/min)
2019-07-15 10:55:04 [scrapy.core.scrapy] DEBUG: Scraped from <200 http://jib.xywy.com/il_sii/food/43.htm>
{'can_eat': ['芝麻', '南瓜子仁', '栗子(熟)', '葵花子仁'],
'cause': '由于汞富于流动性，且易在常温下蒸发，故汞中毒是常见的职业中毒。主要发生在生产中长期吸入汞蒸气或汞化合物粉尘。生产性中毒见于汞矿开采、汞合金冶炼、
剂和防原子辐射材料等生产工人中。\\n'
'汞广泛存在于自然界，各种自然现象可使汞从地表经大气雨雪等环节不断循环，化后随植物所吸收。人类的生产活动可明显加重汞对环境的污染。此种人为污染比
公害病，如水俣病。\\n'
'短时间(>3~5小时)吸入高浓度汞蒸气(>1.0mg/m3)及口服大量无机汞可致急性汞中毒；服用或涂抹含汞偏方可致亚急性汞中毒；职业接触汞蒸气常引起慢性汞中
'发病机制：\\n'
'汞离子易与巯基结合，使与巯基有关的细胞色素氧化酶，丙酮酸脱氢酶，琥珀酸脱氢酶等失去活性，汞还与氨基，羧基，磷酸基结合而影响功能基团的活性，由于这些
，近年来，发现汞对肾脏损害，以肾近曲小管上皮细胞为主，汞还可引起免疫功能紊乱，产生自身抗体，发生肾病综合征或肾小球肾炎。'
'cure_dept': '急诊科',
'drug': ['布美他尼片', '十一味金龟丸'],
'easy_get': '接触汞机会较多的作业工人'
  
```

图 4 正在运行的爬虫

```

{"id": "1", "name": "百日咳", "intro": "百日咳(pertussis, whoopingcough)是由百日咳杆菌所致的急性呼吸
3个月左右，故有百日咳之称。多见于5岁以下的小儿，幼婴患本病时易有窒息、肺炎、脑病等并发症，病死率高。百日咳患者
者病程6-8周，临床病程可分3期：1.卡他期，从发病到开始出现咳嗽，一般1-2周。2.痉咳期，一般2-4周或更长，阵发性痉
咳。一般外周血白细胞计数明显增高，分类以淋巴细胞为主。在诊断本病时要注意与支气管异物及肺门淋巴结结核鉴别。近年
"cure_dept": "儿科小儿内科", "easy_get": "多见于小儿", "get_way": "呼吸道传播", "neopathy": ["支气
性治疗", "treat_prob": "98%", "treat_period": "1-2个月", "treat_cost": "根据不同医院，收费标准不
(Bordetella)中的百日咳鲍特菌(B.pertussis)，常称百日咳杆菌，已知鲍特菌属有四种杆菌，除百日咳鲍特菌外还有副
(B.avium)，鸟型鲍特菌一般不引起人类致病，仅引起鸟类感染，百日咳杆菌长约1.0~1.5μm，宽约0.3~0.5μm，有荚膜
(15%~25%)鲜血才能繁殖良好，故常以鲍-金(Border-Gengous)培养基(即血液，甘油，马铃薯)分离菌落，百日咳杆菌生
为光滑(S)型，又称I相细菌，形态高低一致，有荚膜和较强的毒力及抗原性，致病力强，如将分离菌落在普通培养基中继续培
  
```

图 5 爬取到的数据预览

3.2 合并数据

由于数据是分布式爬取的，需要合并，主要通过 `merge` 函数实现。该函数在获取疾病字典数据后，进行去重、排序等处理，以 JSON 格式写入文件。

```
def merge(diseases: dict):
    set_dis = set()
    print('Original diseases: %d' % len(diseases))
    f = open(datadir + 'medical.json', 'a', encoding='UTF-8')
    for disease in diseases:
        disease['cure_dept'] = [dept + '科' for dept in disease['cure_dept'].split('
科') if dept]
        if disease.get('treat'):
            disease['treat'] = [treat for treat in disease['treat'].split(' ') if treat]
        else:
            disease['treat'] = []
        disease = json.dumps(disease, ensure_ascii=False)
        set_dis.add(disease)
    set_dis = list(set_dis)
    sort_set_dis = sorted(set_dis, key=lambda x: int(json.loads(x)['id']))
    print('Set diseases: %d' % len(sort_set_dis))
    for disease in sort_set_dis:
        f.write(disease + '\n')
    f.close()
```

3.3 Neo4j 知识图谱设计

(1) 节点设计 (Node Design): 一共设计了 5 种标签 (Label) 的节点。

Disease (疾病)		
属性	含义	举例
name	疾病名称	"百日咳"
intro	简介	"百日咳(pertussis, whoopingcough)是..."
cause	起因	"(一)发病原因\n病原菌是鲍特菌属(Bordetella)..."
prevent	预防	"百日咳预防\n1、控制传染源..."
nursing	护理	"确诊的患者应立即隔离至病后 40 天..."
insurance	是否医保	"否"

easy_get	易感人群	"多见于小儿"
get_way	感染途径	"呼吸道传播"
get_prob	感染概率	"0.5%"
treat	治疗方式	"药物治疗", "支持性治疗"
treat_prob	治愈概率	"98%"
treat_period	治疗时间	"1-2 个月"
treat_cost	治疗费用	"根据不同医院, 收费标准不一致, 市三甲..."
treat_detail	治疗细节	"百日咳西医治疗 1、药物治疗..."
Department (科室)		
属性	含义	举例
name	科室名称	"小儿内科"
Drug (药物)		
属性	含义	举例
name	药物名称	"阿司匹林肠溶片"
Food (食物)		
属性	含义	举例
name	食物名称	"小米面"
Symptom (症状)		
属性	含义	举例
name	症状名称	"指甲呈筒形"

(2) 关系设计 (Relationship Design): 共有 6 种关系。

关系类型	含义	举例
has_symptom	有症状	("大叶性肺炎", "有症状", "高热")
can_use_drug	可用药	("大叶性肺炎", "可以用药", "乳酸左氧氟沙星片")
belongs_to	属于	("大叶性肺炎", "属于", "呼吸内科")
has_neopathy	有并发症	("大叶性肺炎", "有并发症", "肺脓肿")
can_eat	可以吃	("大叶性肺炎", "可以吃", "鲫鱼")
not_eat	不能吃	("大叶性肺炎", "不能吃", "韭菜")

3.4 Rasa 的配置

(1) Rasa NLU 训练数据

要让电脑对人类的语言做出回应，首先要进行实体抽取和意图识别，Rasa NLU 正是负责这个模块，使用 Chatito 工具可以快捷地构建 Rasa NLU 所需的训练数据。

```
%[search_treat]('training': '3000', 'testing': '300')
*[50%] ~[get?][disease]~[时?][的?][treat_qwd?][qwd?]
*[10%] ~[treat_qwd?][disease]~[时?][的?][qwd?]
*[15%] @[disease]~[是什么][病?][qwd?]
*[15%] ~[什么是][disease]~[qwd?]
*[10%] @[disease]

%[search_food]('training': '1500', 'testing': '150')
*[30%] ~[get?][disease]~[时?][的?][food_qwd?][qwd?]
*[30%] ~[get?][disease]~[时?][的?][food_qwd?][qwd?]
*[40%] ~[get?][disease]~[时?][的?][可以吃?][什么][food_qwd?][qwd?]

%[search_symptom]('training': '1500', 'testing': '150')
*[33%] ~[get?][disease]~[时?][会?][不会?][有?][什么?][symptom_qwd?][qwd?]
*[33%] ~[get?][disease]~[的?][symptom_qwd?][qwd?]
*[33%] ~[get?][disease]~[时?][有?][什么?][symptom_qwd?][qwd?]
```

图 6 使用 Chatito 工具制作训练数据

(2) Rasa Core 训练数据

Rasa Core 使用机器学习模型来预测机器人下一步该采取的动作，训练数据是由故事（Stories）构成，故事就是真实的对话流程，通常还会带有词槽（Slot）的填充，需要一定的人工整理。

```
## story_1_search_treat_simple
* greet
  - utter_greet
* search_treat{"disease": "百日咳"}
  - action_search_treat
* bye
  - utter_goodbye

## story_2_search_food_simple
* greet
  - utter_greet
* search_food{"disease": "百日咳"}
  - action_search_food
* bye
  - utter_goodbye
```

图 7 人工整理的对话流程

(3) Rasa 配置的选择

因为本项目是中文环境的问答系统，所以选择如下的配置：

pipeline:	policies:
- name: "MitieNLP"	- name: "rasa.core.policies.keras_policy.KerasPolicy"
model:	epochs: 120
"data/total_word_feature_extractor_zh.dat"	featurizer:
- name: "JiebaTokenizer"	- name: MaxHistoryTrackerFeaturizer
dictionary_path: "jieba_userdict"	max_history: 5
- name: "MitieEntityExtractor"	state_featurizer:
- name: "EntitySynonymMapper"	- name: BinarySingleStateFeaturizer
- name: "RegexFeaturizer"	- name:
- name: "MitieFeaturizer"	"rasa.core.policies.memoization.MemoizationPolicy"
- name: "SklearnIntentClassifier"	max_history: 5
	- name:
	"rasa.core.policies.mapping_policy.MappingPolicy"
	- name: "rasa.core.policies.fallback.FallbackPolicy"
	nlu_threshold: 0.4
	core_threshold: 0.3
	fallback_action_name: 'action_donknow'

Pipeline 配置

Policy 配置

(4) 现有数据的利用：

由于 Pipeline 中 MITIE 所构建的词向量维数高达 50 多万，所以训练过程对硬件要求较高，故使用了已经训练好的 data/total_word_feature_extractor_zh.dat 模型，该模型训练使用的是中文 Wikipedia Dump 和百度百科语料。

因为本项目是医疗领域，涉及大量的疾病、症状、药物等名词，同时 MITIE 的训练需要输入的是词，所以为了较好的分词效果，从已采集的数据中抽离出 Diseases、Departments、Drugs、Symptoms 和 Foods 作为 Jieba 分词组件的自定义词典。

3.5 前端

考虑到 Rasa 框架中带有封装好的可启动的服务器程序，Rasa 官方也推荐了基于 React 和 Socket 的开源聊天前端框架 Rasa-WebChat，它可以返回并渲染 Markdown 格式文本，于是本项目采用了网页的形式，并引入了 WebChat.js，通过 Http 请求调用 Rasa 服务。



图 8 点击箭头所指的图标

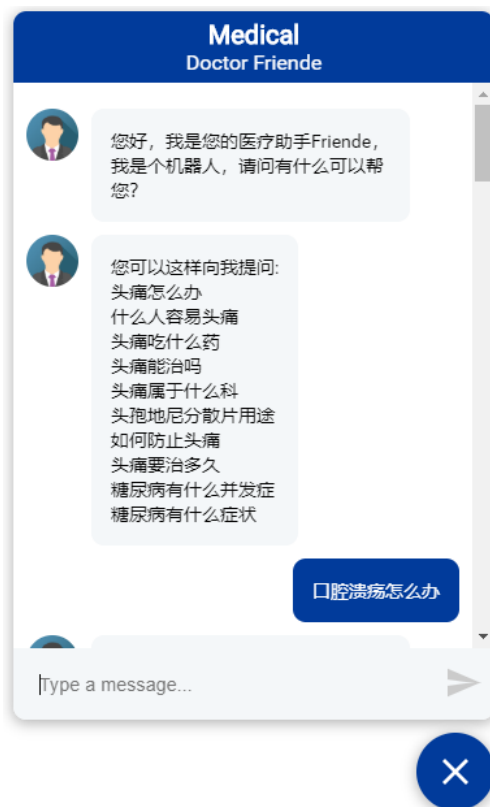


图 9 聊天框效果图

4. 结果展示

4.1 可用功能

意图	含义	示例
search_treat	查询疾病治疗方法	头痛怎么办
search_food	患某病能、不能吃的食物	头痛能吃什么
search_symptom	查询疾病症状	头痛有什么症状
search_cause	查询疾病起因	怎么会头痛
search_neopathy	查询疾病并发症	头痛的并发症
search_drug	查询疾病可用药	头痛吃什么药
search_prevention	查询疾病预防方法	头痛怎么预防
search_drug_func	查询药物能治疗的疾病	牛黄解毒丸能治啥
search_disease_treat_time	查询疾病治疗时长	头痛要治多久
search_easy_get	查询疾病易感人群	什么人容易得头痛
search_disease_dept	查询疾病所属科室	头痛属于什么科

4.2 知识图谱



图 10 某一疾病的关系结构

Node Labels

*(13635)

Department

Disease

Drug

Food

Symptom

Relationship Types

*(114163)

belongs_to

can_eat

can_use_drug

has_neopathy

has_symptom

not_eat

Property Keys

cause

easy_get

get_prob

get_way

insurance

intro

name

nursing

prevent

treat

treat_cost

treat_detail

treat_period

treat_prob

图 11 图谱概况一览

4.3 前端

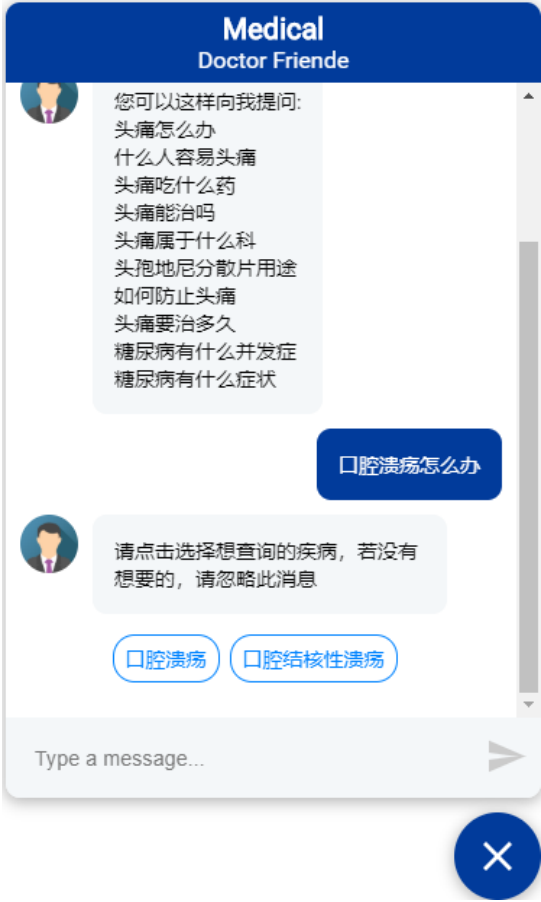


图 12 对疾病名称的模糊搜索

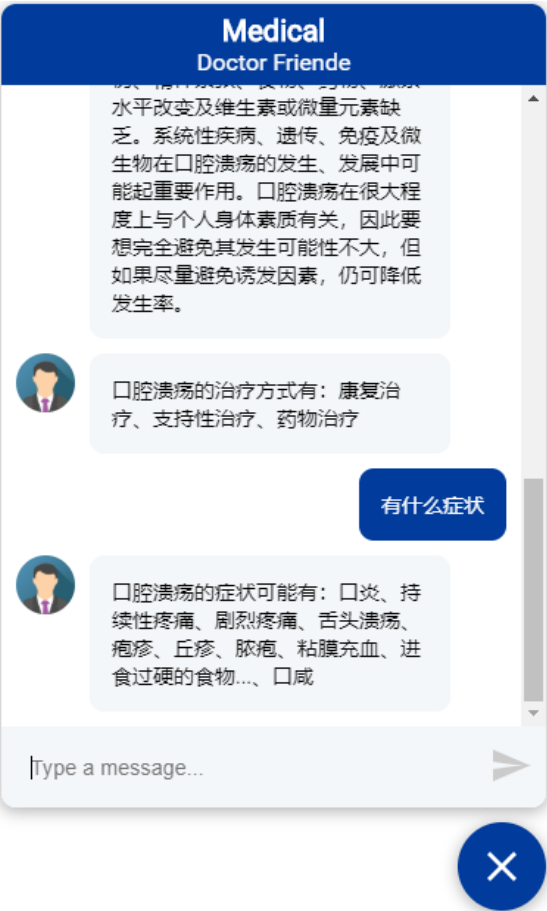


图 13 多轮查询记忆疾病名

5. 感想与收获

本次“聊天机器人”远程指导项目为期一个月，虽然“自然语言处理”的字眼已在我的学习生活中出现多次，但我一直没有实际地去学习它，同样我也不曾想过它会以“聊天机器人项目”的形式与我碰面。

在参与项目过程中，老师讲授了自然语言处理的发展和应用场景、正则表达式的应用、词向量的构成、命名实体识别、意图识别、否定、多轮查询和有状态对话等知识点，概念容易理解，但是真正使用时却困难重重。

到了项目制作的阶段，就有一种“病来如山倒病去如抽丝”的感觉。从确定医疗领域开始到最后完成一个基础的聊天机器人，我划分了如下阶段：

(1) 构想阶段：

- 搜索是否有前人做过类似的程序，并试用、学习；

(2) 数据阶段：

- 发现现有可取的数据较少，亦不想直接取用他人整理的数据，于是学习爬虫技术，学习 Scrapy 框架；
- 寻找医疗方面的网站并尝试爬取数据；
- 整理、存储数据，并设计数据入库的结构；
- 学习 Neo4j 的使用，学习 py2neo 库；
- 熟悉 Rasa NLU 和 Rasa Core 的训练数据的结构，学习使用 Chatito 工具构造数据；

(3) 构造阶段：

- 学习 Rasa 框架的结构、工作流程；
- 学习国人制作的 rasa_nlu_chi，跑通 demo、进一步熟悉 Rasa；
- 熟悉新版 Rasa 与旧版的一些变化，并尝试用新版跑通 demo；
- 学习 Rasa 自定义 action 的用法，并设计 Rasa Core 的响应逻辑。
- 训练模型，命令行调试。
- 制作 Web 前端界面，调通程序。

(4) Debug 阶段：

- 尝试在本地聊天，并在发现问题后补充、修正数据集，重新训练，修正 action。

(5) 部署阶段

过程不可谓不充实，从项目出发，发现问题，学习解决问题的技术，利用技术解决问题，形成一个循序渐进、螺旋式的学习过程。本次远程指导项目收获的不只是知识，由于是单人完成，所以更锻炼了自己做科研时应有的条理性、逻辑性的思维方式。

6. 参考

- [1] 刘焕勇 [QABasedOnMedicalKnowledgeGraph](#)
- [2] 王冠 [Rasa_NLU_Chi](#)
- [3] 孔晓泉 [WeatherBot](#)
- [4] 孔晓泉 [Rasa NLU 的 Pipeline 和 Component](#)
- [5] mrbot-ai [Rasa-WebChat](#)
- [6] Scrapy Official [Scrapy](#)
- [7] Neo4j Official [Neo4j](#)
- [8] py2neo Official [py2neo](#)
- [9] Rasa Official [Rasa-Doc](#)
- [10] Rasa Official [Legacy-Rasa-Doc](#)
- [11] Rasa Official [Rasa-Forum](#)