

本项目是基于知识图谱的医疗问答机器人。使用者可以通过输入自然语言的问题来获取有关医疗健康方面的答案。

问答机器人的答案来源于知识图谱中已有的知识点和相关实体，并且采用了多种自然语言处理技术来提高识别和推理能力，最终能够给出准确、全面的答案。

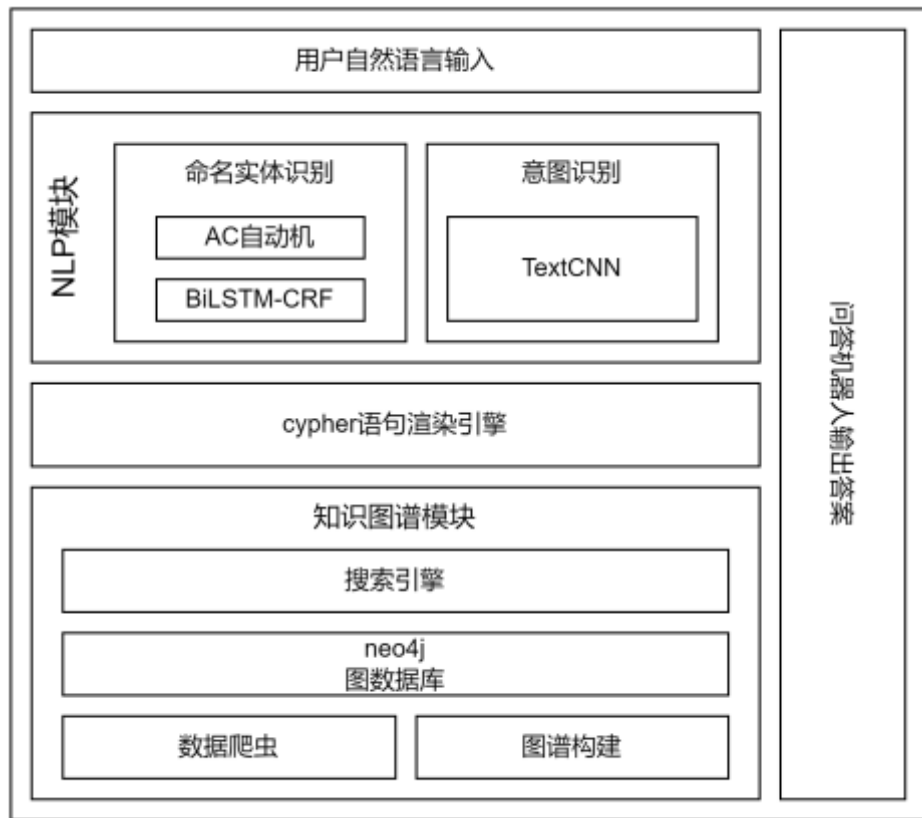
项目介绍

02 / 30



整体架构

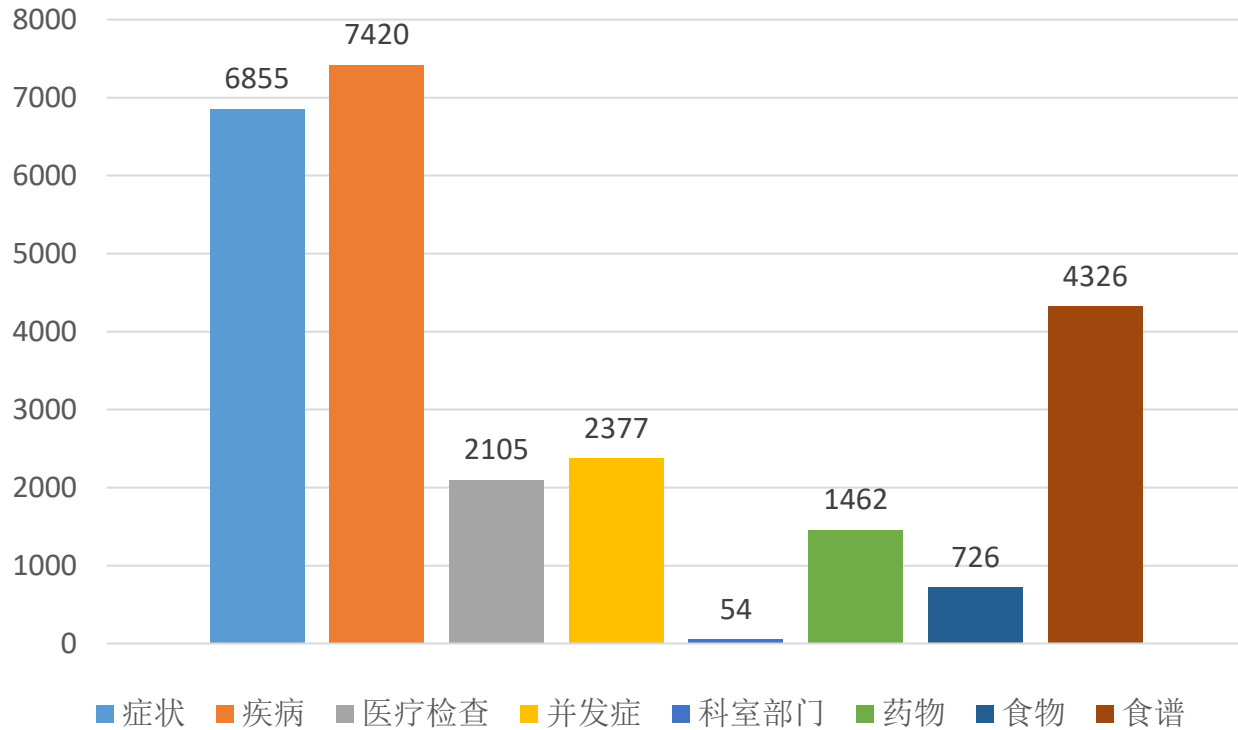
03 / 30



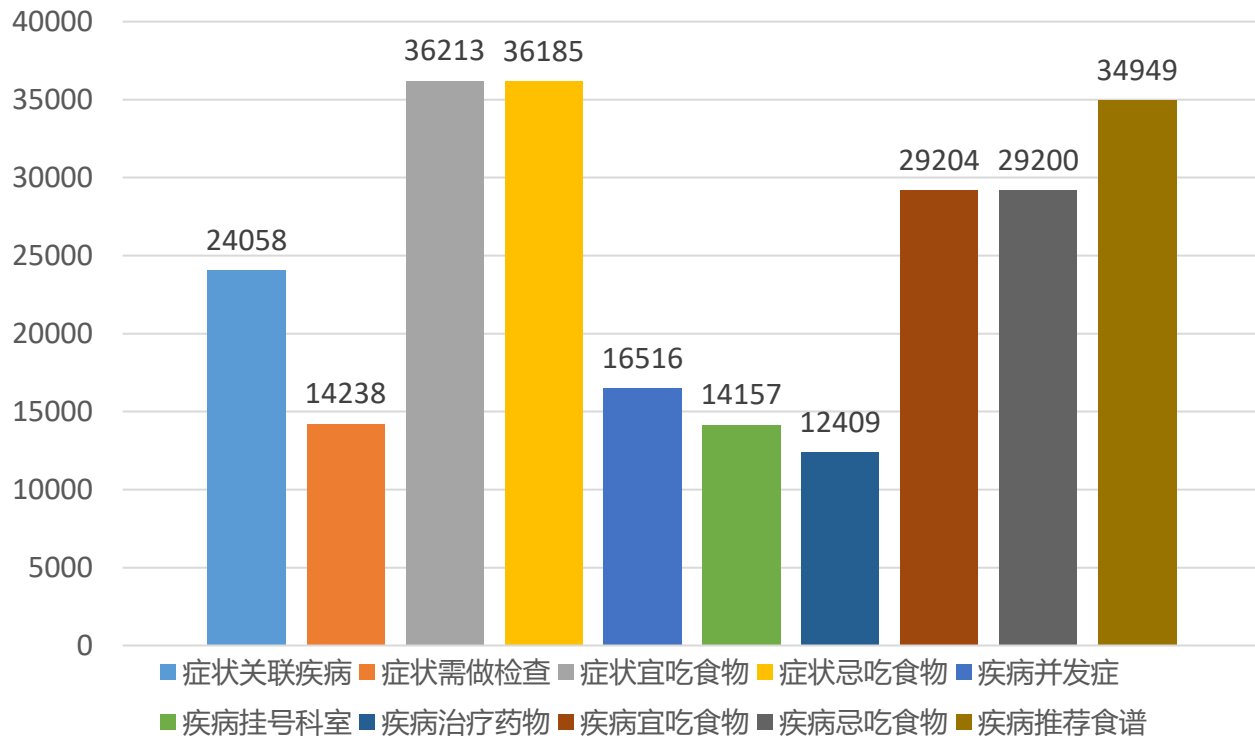
自然语言处理模块负责命名实体识别、意图识别等任务，cypher语句渲染完成后，由搜索引擎进行聚合与执行，知识图谱用于存储和管理医疗健康数据，整个流程中各模块协同工作，完成整个系统的问答过程。

构建知识图谱的数据来源于寻医问药网，通过爬虫进行抓取，经过数据清洗之后形成结构化数据，在本项目中主要抓取的是症状数据和疾病数据，其他实体（如医疗检查、科室部门和食物等）信息可从这两类数据中抽取。

实体统计数据

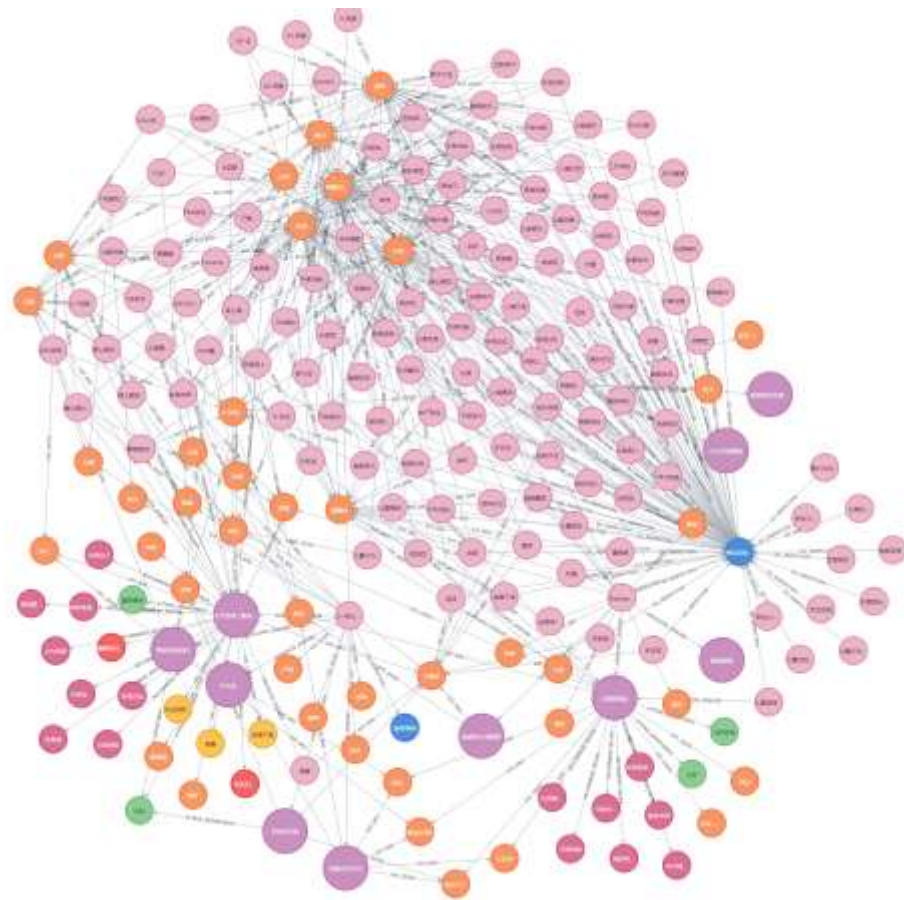


关系统计数据



知识图谱

08 / 30



命名实体识别是自然语言处理的一项重要任务，目的是在文本中找到具有特定意义的实体，例如人名、地名、组织机构名等，本项目中对应的则是找出问句中的症状和疾病实体。具体来说，使用AC自动机和BiLSTM-CRF进行命名实体识别，有效提高了识别准确度和效率。

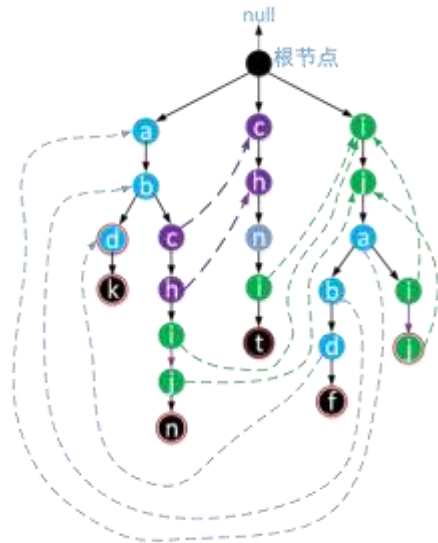
AC自动机

10 / 30

AC自动机可以用于模式匹配和字符串查找等问题，其核心思想是构建一个字典树，然后将字典树上的每个节点和某一模式串的前缀对应，构成一个状态机，即AC自动机。

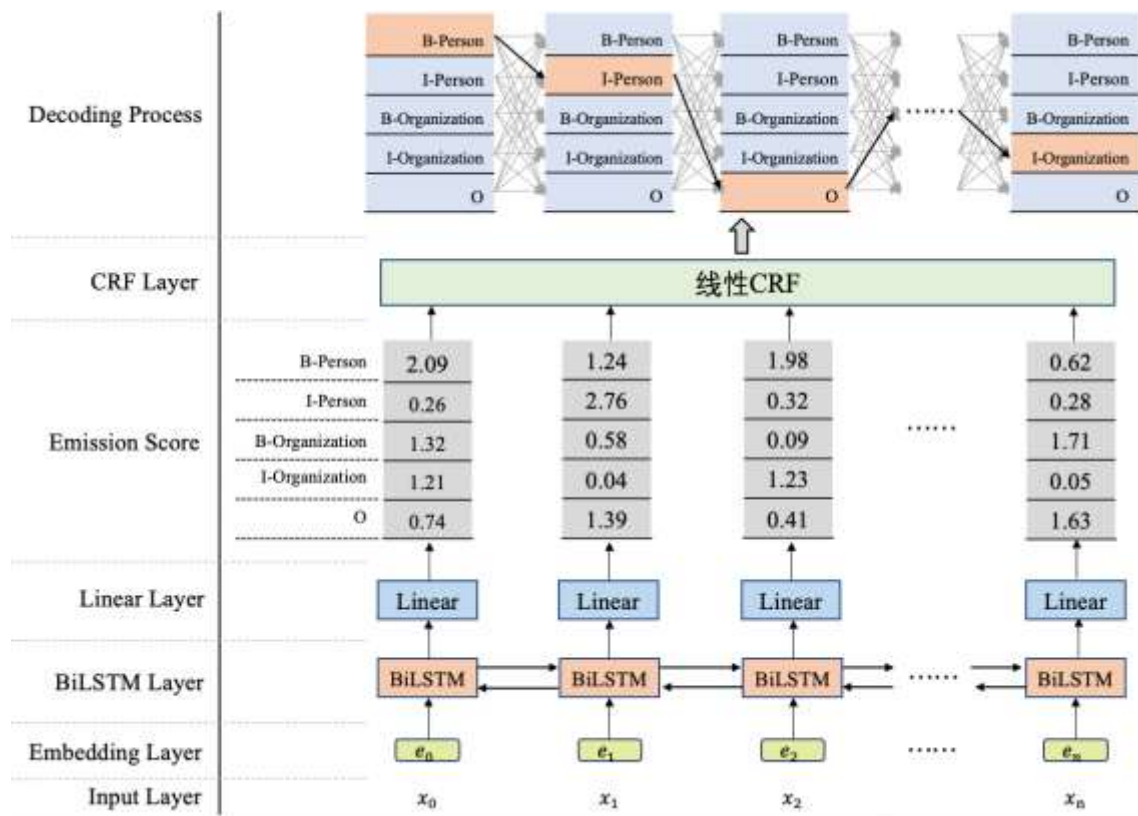
从AC自动机的根节点开始搜索，在搜索的过程中，如果当前节点存在失配指针，则跳转到失配指针对应的节点；如果当前节点是一个匹配成功节点，则将其所对应的模式串加入到匹配结果中。

AC自动机的时间复杂度是 $O(n)$ ，其中 n 为目标字符串的长度。



BiLSTM-CRF

11 / 30



BiLSTM-CRF是一种序列标注模型，其结构包括BiLSTM和CRF两部分。BiLSTM是一种双向长短时记忆网络，可以对输入的序列进行有效的编码和提取特征。CRF是一种条件随机场，可以对标注序列进行联合概率建模，从而提高序列标注的准确性。

在BiLSTM-CRF中，首先使用BiLSTM对输入序列进行编码，得到每个时刻的隐状态表示。然后将隐状态表示传递给CRF层，进行联合概率建模并得到最优的标注序列。最终输出的标注序列即为模型预测的结果。

相比于传统的序列标注模型，BiLSTM-CRF具有以下优点：

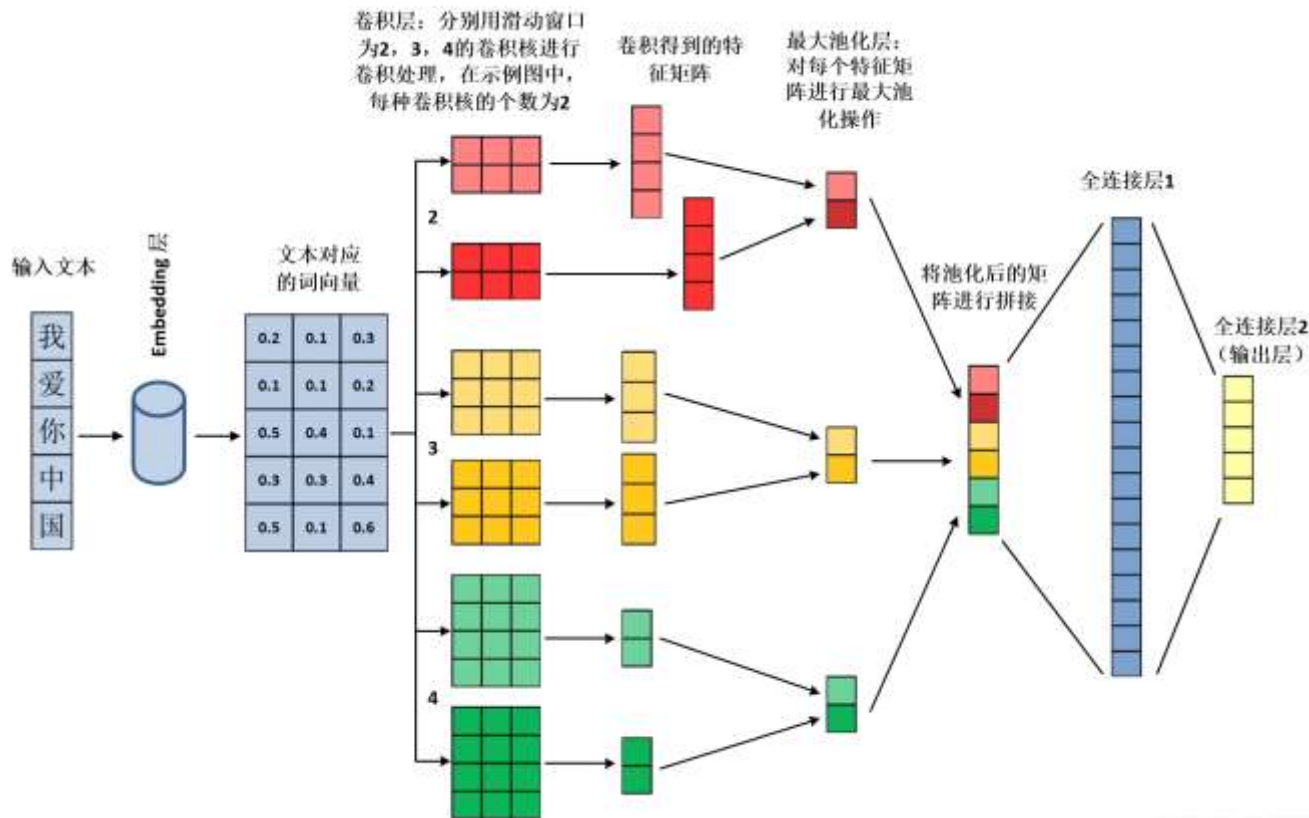
1. BiLSTM能够有效地捕捉输入序列中的上下文关系，从而提高特征表示的能力。
2. CRF层能够考虑相邻标签之间的依赖关系，并且利用全局约束来解决标签依赖问题，使得标注结果更加合理和一致。
3. BiLSTM-CRF可以进行端到端的训练和预测，避免了传统模型中需要进行特征工程和后处理的麻烦。

BiLSTM-CRF已经在自然语言处理领域的诸多任务中取得了良好的表现，例如命名实体识别、句法分析、情感分析等。

意图识别是指要确定对问题的回答属于特定类别，例如查询疾病、药品、治疗方法等。本项目采用了TextCNN模型作为意图识别的核心算法，结合NER层人工嵌入特征后的输入数据，实现快速、准确的意图分类。

TextCNN

15 / 30



TextCNN (Text Convolutional Neural Network) 是一种卷积神经网络 (CNN) 的变种，主要用于自然语言处理 (NLP) 任务，例如文本分类、情感分析和命名实体识别。

TextCNN的核心思想是将文本转化为一个矩阵，并使用卷积层进行特征提取，然后将提取出的特征传递给全连接层进行分类。

TextCNN的优点包括：

1. 高效的文本分类模型，可以在短时间内处理大量的文本数据。
2. 利用卷积神经网络的局部关系特性，能够更好地提取文本中的特征，无需人工提取特征。
3. 采用多个不同大小的卷积核对数据进行卷积和池化操作，可以对不同大小的语义信息进行提取，从而更加全面的理解文本内容。

TextCNN的缺点包括：

1. 在处理较长文本时，可能会出现信息丢失或混淆的问题。
2. 最大的问题全局max pooling丢失了结构信息，很难发现文本中的转折关系等复杂模式。
3. 只知道关键词是否在文本中出现了，以及相似度强度分布，不可能知道关键词出现了几次，以及出现这些关键词的顺序。

数据冷启动指的是当一个系统或模型开始运行时，缺乏足够的已有数据以支撑其学习和运作，从而导致性能不佳、准确率低下的情况。

在机器学习领域中，数据冷启动是指在训练模型之前，无足够的已有数据集来进行模型训练的情况，例如：

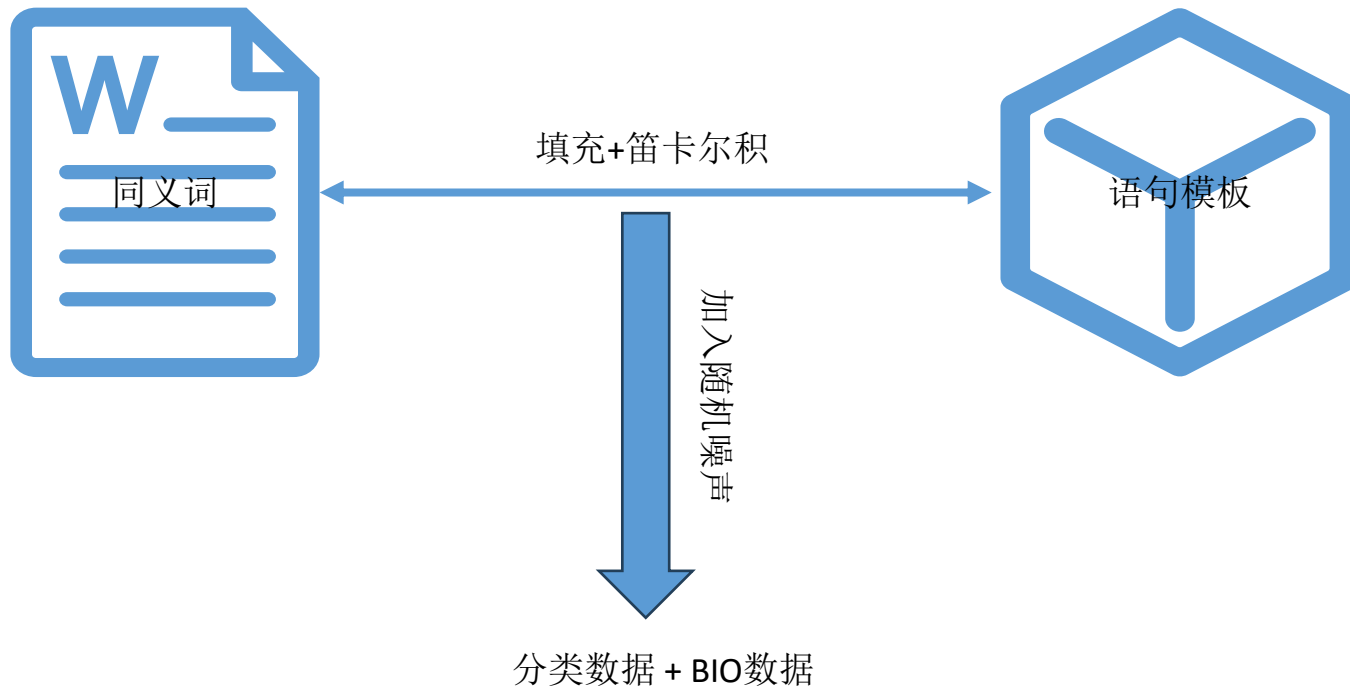
数据集不足：当某一特定领域的数据集非常珍贵，如医疗、金融等领域，可能由于数据保密、获取难度大等原因，导致缺乏足够的数据进行模型训练。

新领域探索：当研究领域刚开始探索一种新的现象时，缺乏足够的数据集来支持训练模型。

对于数据冷启动问题，常用的解决方法包括利用模拟生成数据、利用迁移学习、利用半监督学习等方法来缓解数据缺失的问题，以提高模型训练和性能。

数据集

20 / 30



数据集

21 / 30

```
# 同义词，用于构造语句数据
WORDS_SEQ = {
    'how': ['怎么样', '如何', '何以', '怎么', '怎么', '怎么', '怎样', '怎么', '怎样', '咋样', '咋样', '咋样', '咋地', '咋地', '咋地'],
    'why': ['为什么', '为何', '何故', '何以', '为啥', '何以故'],
    'what': ['什么', '何', '甚', '哪', '哪些', '那', '有何', '何事', '何物', '什么事', '什么物'],
    'do': ['办', '办', '办', '办', '办', '做', '做', '做', '做', '做', '处理', '解决', '操作'],
    'is': ['是', '在', '属于', '在于', '代表', '表示', '指的是', '意味着', '包括', '归属于'],
    'may': ['可能', '或许', '可能会', '恐怕', '大概', '也许', '有可能', '很可能'],
    'should': ['应该', '必需', '需要', '要', '希望', '愿意', '可以', '理应', '当然', '有必要'],
    'feel': ['感受', '觉得', '意识到', '认为', '知觉', '发觉', '感应', '体验', '体认', '感觉', '察觉', '出理', '呈现'],
    'cure': ['治疗', '医治', '治愈', '康复', '痊愈', '治愈', '止旺', '治好', '缓解', '减轻'],
    'illness': ['疾病', '病症', '病患', '疾患', '病弱', '疾疫', '流行病', '病毒性疾病', '传染病'],
    'suffer': ['得', '患', '感染', '患染', '发病', '罹患', '患有'],
    'explain': ['介绍', '了解', '告知', '理解', '解释', '详情', '信息', '定义'],
    'prevent': ['预防', '防范', '抵制', '抵御', '防止', '避免', '免得', '避开', '免于'],
    'check': ['检查', '检测', '体检', '化验', '验证', '查验', '查询', '查', '测', '检', '监测'],
    'take': ['吃', '食', '服', '摄入', '饮', '喝', '食用', '服用', '饮食', '饮用', '伙食', '膳食', '忌口', '补品',
            '食谱', '菜谱', '食物', '补品'],
    'good': ['好', '有益', '改善', '提升', '增加', '增进', '增强', '益处', '好处', '适合', '适宜', '宜', '宜于'],
    'bad': ['不好', '不益', '差', '降低', '减少', '不良', '糟糕', '弊处', '恶害', '有害', '不宜', '弊'],
    'suggest': ['建议', '提示', '忠告', '注意', '推荐', '帮助', '指南', '提醒']
}
```

数据集

22 / 30

```
# 谓词模板。用于构造谓词数据
TEMPLATES_SEQ = {
    '症状谓词法': [
        ('entity#symptom', 'how', 'do'),
        ('how', 'do', 'entity#symptom'),
        ('entity#symptom', 'how', 'cure'),
        ('how', 'cure', 'entity#symptom'),
        ('how', 'cure', 'entity#symptom'),
        ('entity#symptom', 'what', 'do'),
        ('what', 'do', 'entity#symptom'),
        ('entity#symptom', 'what', 'cure'),
        ('what', 'cure', 'entity#symptom'),
        ('what', 'cure', 'entity#symptom'),
    ],
    '症状谓词病': [
        ('entity#symptom', 'suffer', 'what', 'illness'),
        ('suffer', 'what', 'illness', 'entity#symptom'),
        ('why', 'feel', 'entity#symptom'),
        ('feel', 'entity#symptom', 'why'),
    ],
    '症状谓词预防': [
        ('how', 'prevent', 'entity#symptom'),
        ('entity#symptom', 'how', 'prevent'),
        ('what', 'prevent', 'entity#symptom'),
        ('prevent', 'entity#symptom', 'what'),
        ('do', 'what', 'prevent', 'entity#symptom'),
        ('prevent', 'entity#symptom', 'do', 'what'),
    ],
}
```

数据集

23 / 30

意图分类数据

训练集：12831条

测试集：1431条

```
text,label
呀腱反射消失symptom什么事办,症状问办法
潮汗symptom那止痊,症状问办法
怎么医治先兆流产symptom4,症状问办法
免疫性肾炎symptom怎么办,症状问办法
何p治癒饭后困倦symptom,症状问办法
晶状体青光眼斑symptom咋地办,症状问办法
肢体运动不协调symptom5噢咋地做,症状问办法
咀嚼肌痉挛symptom哪些止痊,症状问办法
什么事治好腹泻symptom,症状问办法
皮肤色素沉着呈青铜色symptom如何!做,症状问办法
1咋样办运动耐力降低symptomj,症状问办法
肝淀粉样变性symptom什么事缓解,症状问办法
哈g那缓解前臂及手部肌群的缺血性挛缩symptom,症状问办法
哈基1治疗脾胃不和symptom,症状问办法
```


数据集

24 / 30

命名实体识别数据

训练集：12848条

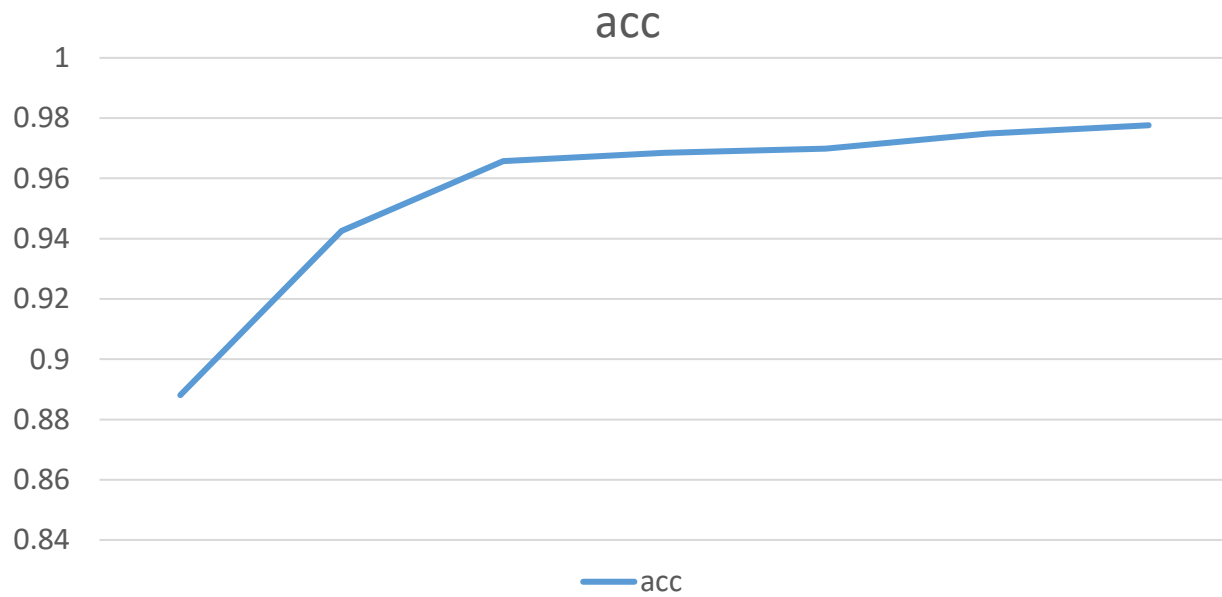
测试集：1429条

p哪验证生化妊娠,0 0 0 0 B-SYMPTOM I-SYMPTOM I-SYMPTOM I-SYMPTOM
何以康复黄色结节,0 0 0 0 B-SYMPTOM I-SYMPTOM I-SYMPTOM I-SYMPTOM
粘膜疹咋地医治,B-SYMPTOM I-SYMPTOM I-SYMPTOM 0 0 0 0
提示咋地操作癫痫性精神障碍,0 0 0 0 0 0 B-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE
啊管型尿咋地检查,0 B-SYMPTOM I-SYMPTOM I-SYMPTOM 0 0 0 0
甚不好婴儿青铜综合征,0 0 0 B-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE
下腹剧痛并渐向腹中线扩散咋样医治,B-SYMPTOM I-SYMPTOM I-SYMPTOM I-SYMPTOM I-SYMPTOM I-SYMPTOM
哪处理盘状红斑狼疮v、,0 0 0 B-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE
骶髂筋膜脂肪疝注意何,B-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE
那表示膀胱损伤,0 0 0 B-DISEASE I-DISEASE I-DISEASE I-DISEASE
咋地治愈脊椎病1,0 0 0 0 B-DISEASE I-DISEASE I-DISEASE 0
w葡萄球菌感染有害啊哪,0 B-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE
哪降低皮样表皮样囊肿,0 0 0 B-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE I-DISEASE
有何恶劣小儿乳积,0 0 0 0 B-DISEASE I-DISEASE I-DISEASE I-DISEASE
梅尼埃病mr信息,B-DISEASE I-DISEASE I-DISEASE I-DISEASE 0 0 0 0

模型评估

25 / 30

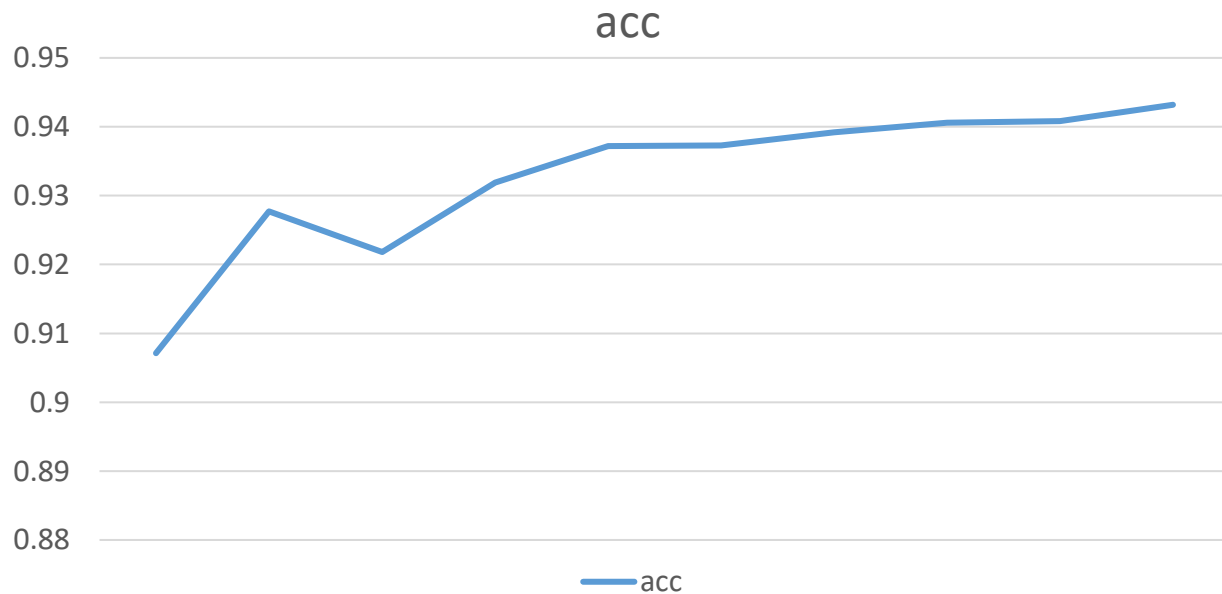
意图识别模型：7个epoch后acc达到97.76%



模型评估

26 / 30

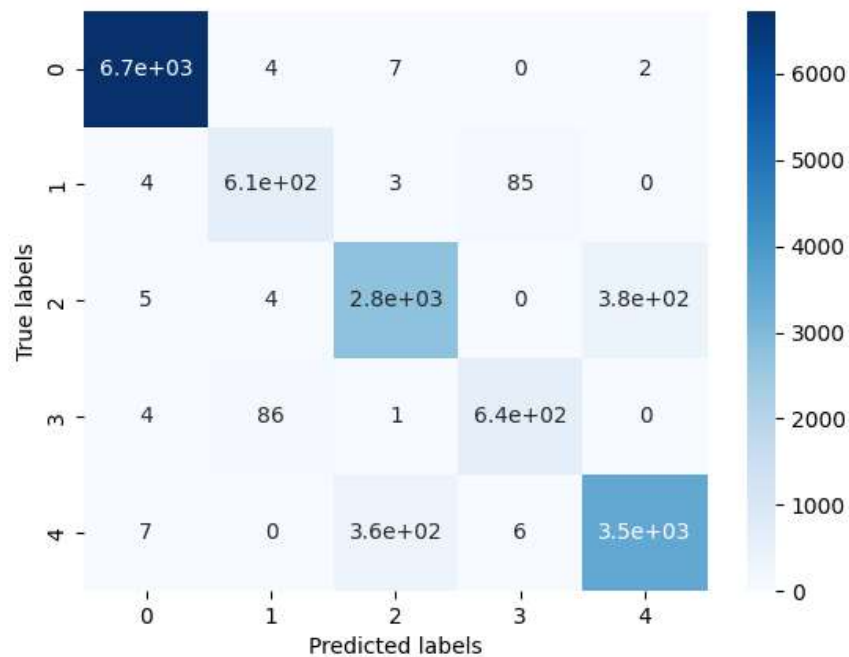
命名实体识别模型：10个epoch后，标签整体acc达到97.76%



模型评估

27 / 30

命名实体识别模型混淆矩阵



TextCNN只知道关键词是否在文本中出现了，以及相似度强度分布，不可能知道关键词出现了几次，以及出现这些关键词的顺序。

针对这个问题，可以尝试k-max pooling优化，针对每个卷积核都不只保留最大值，而是保留前k个最大值，并且保留这些值出现的顺序，即按文本的顺序排列k个最大值。

可以结合Bert优化BiLSTM-CRF模型的表现。

Bert是一种预训练模型，在自然语言处理领域已经产生了广泛的应用，包括命名实体识别、句法分析、情感分析等。Bert的优点在于其在大规模语料库上训练得到了丰富的语言表示，这些表示可以应用在不同的自然语言处理任务上。

使用Bert对输入词汇进行编码，用Bert提取的语义特征替换Embedding层的词向量，以提高模型的语义表示能力。

目前通过NER输出的实体信息，是通过计算Jaccard相似度，映射到知识图谱中的实体，这种方法没有考虑两个词之间的语义相似度，单纯的依赖字符交并集的做法无法保证较好的映射效果。

可用通过预训练的Word2Vec结合本地知识图谱的数据，进行调优，使用调优后的模型能更好的计算两个词之间的语义相似度。