

《AI 疫情对话机器》 调研报告

目录

- 一、 调研的主要内容 3
 - 1.1 调研报告描述 3
 - 1.2 论文主要内容 3
 - 1.3 论文分析 3
 - 1.4 调研报告主要内容 4
- 二、 原论文解析 5
 - 2.1 项目背景 5
 - 2.2 相关工作 5
 - 2.3 疫情对话机器人框架 6
 - 2.4 实现 8
 - 2.5 结论和展望 9
- 三、 相关文献调研 10
- 四、 模型复现 10
- 五、 参考资料 11

一、调研的主要内容

1.1 调研报告描述

此调研报告是参考《隔离中与隔离结束期间基于深度情感分析模型的智能新冠普适聊天机器人》[《Smart Ubiquitous Chatbot for COVID-19 Assistance with Deep learning Sentiment Analysis Model during and after quarantine》](#)所做的奇思AI疫情对话机器人项目调研报告。

1.2 论文主要内容

论文实现了一个对话机器人，帮助疫情隔离中与隔离结束的人克服困境、收集新冠接触信息、提出医学建议等。2019年末出现的新型冠状病毒已经造成了大量的人口死亡，任何性别、年龄、地区的人都可能感染新冠，新冠对人类和社会构成了严重的威胁。在生活中，一般有两类人，一类人是不知道的新冠危险的人，同时这类人也加速了新冠的传播速度；另一类是对新冠特别恐慌的人。本论文的目标是开发一种针对新冠可以连续对话的机器人，在新冠隔离期间以及隔离后能让人们提高对新冠的认识以及医疗帮助。本文实现的对话机器人可以有效减缓新冠的传播，它是由四个相关依赖的模块组成的：1) 信息理解模块(IUM)，由自然语言处理(NLP)实现；2) 数据收集器模块(DCM)，收集用户的非隐私信息；3) 消息生成模块(AGM)，生成聊天回答；4) 情感检测模块(DDM)，使用胜读学习情感分析模型检测用户的情感，适当地传递安慰信息。

1.3 论文分析

论文的新颖性：

- 1、用聊天机器人能普及新冠知识，减缓病毒传播。
- 2、在机器人对话中，增加情感分析模块，更加准确把握用户心态

论文的知识点：

- 1、学习到自然语言处理(NLP)的方法
- 2、学习构建深度学习模型的一般方法
- 3、学习使用机器学习模型，如支持向量机、随机森林等
- 4、学习到基于python flask构建后台系统

论文的复现难度：

- 1、数据集获取具有难度。为了构建聊天机器人，论文中使用了多个的机器学习模型，相应的也需要大量的样本数据集作训练。文章采用了多个公开的数据集训练模型，但是在消息生成模块(AGM)中，由于没有医生和普通人之间关于新冠的公开对话数据源，所以作者自己构建了一份数据集。该数据集可尝试向本文的通讯作者索取。
- 2、论文的复现，要有构建机器学习模型、深度学习模型、后台开发、前端开发等相关经验。

论文的可提升点：

- 1、论文在构建情感分析模型(DDM)时，数据集小、模型简单，可以优化模型。
- 2、论文中使用对话机器人收集到的用户新冠接触的流行性疾病的调查信息（流

调), 直接发送给专业的医疗人员做判断。可以构建一个基于个人流调信息的新冠感染模型, 判断用户感染的概率给用户及医务人员作参考。

- 3、可以使用实现语音聊天功能代替文字聊天, 增加用户使用的便利度
- 4、实现新冠对话机器人的理论知识后, 可以构建一个私人陪伴对话机器人, 利用微信聊天内容、研究对象语音等信息, 构建一个陪伴机器人, 比如在老人伴侣离世后, 伴侣机器人可以给老人安慰。也可以构建一个拥有自己声音和聊天习惯的AI机器人。

1.4 调研报告主要内容

第一章, 调研的主要内容, 介绍了本篇调研报告的主要内容, 包含原论文内容简介、论文分析(新颖性、知识点、复现可行性、提升点)。

第二章, 原论文解析, 对原论文内容进行了翻译及解析

第三章, 其他参考文献调研

第四章, 论文复现

第五章, 参考文献

二、原论文解析

2.1 项目背景

COVID-19 是 2019 年的新型冠状病毒的缩写，是一种可在人与人之间传播的呼吸道疾病。新冠病毒是于 2019 年 12 月在中国武汉市首次报道的。根据世界卫生组织的数据，截至 2020 年四月，全球已有超过 144000 人死于新冠，超过 200 万人感染。在新冠病例与日俱增的同时，有一些人并没有意识到新冠带来的威胁，所以民众没有采取预防措施，导致新冠的大规模流行；还有一些人由于新冠带来的威胁过度恐慌和绝望。鉴于以上原因，论文提出了一种普适的聊天机器人，帮助普通人随时随地使用新冠聊天服务，尤其是那些在隔离中与隔离结束后的人们。

2.2 相关工作

《Artificial Intelligence Application in COVID-19 Diagnosis and Prediction》这篇文章的作者开发了一种人工智能方法来争端和预测新冠，但是该方法只是供临床使用，对于普通人没有什么用处。

《Identification of COVID-19 Can be Quicker through Artificial Intelligence framework using a Mobile Phone-Based Survey in the Populations when Cities/Towns Are Under Quarantine. Infection Control and Hospital Epidemiology》实现基于手机信息的人工智能模型，来识别新冠病例。

《apid AI Development Cycle for the Coronavirus (COVID-19) Pandemic: Initial Results for Automated Detection and Patient Monitoring using Deep Learning CT Image Analysis》开发了一种基于人工智能的自动化胸部 CT 图像分析工具，用于检测、量化和跟踪新冠阳性病例。

[Bespoke, Bebot Launches Free Coronavirus Information Bot](#) 日本的 Bespoke 公司利用人工智能推出了一个名为 Bebot 的在线聊天机器人，它提供有关当前疫情的最新信息、检查用户身体状况。

下表显示了各种聊天机器人机器特性。

Table 1: Summary of related works

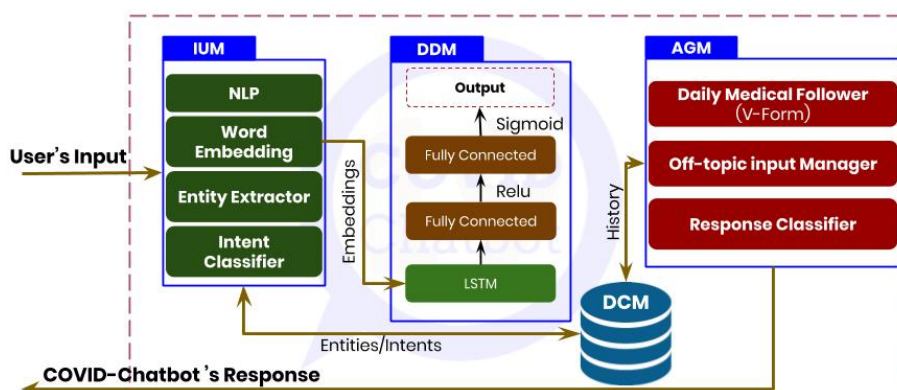
Reference	Chatbot	Service Delivered	Smartness	Ubiquity	Beneficiary	DL Sentiment Analysis
Bespoke 2020 [3]	✓	COVID-19 diagnosis and prediction	✓	✗	Ordinary Citizens	✗
Peng et al. 2020 [7]	✗	COVID-19 diagnosis and prediction	✓	✗	Experts	✗
Allam et al. 2020 [8]	✗	Manage and monitor urban health (including COVID-19) in smart cities	✓	✗	Experts	✗
Maghdid et al. 2020 [10]	✗	Predict the severity of the pneumonia	✓	✗	Experts	✗
Gozes et al. 2020 [5]	✗	Predict the severity of the pneumonia	✓	✗	Experts	✗
Ouerhani et al. 2019 [11]	✓	Emergency Case assistance	✓	✓	Ordinary Citizens	✗
Chih-Wei et al. 2018 [17]	✓	Open	✓	✗	Experts	✓
Hanai et al. 2018 [15]	✗	Depression Detection	✓	✗	Experts	✓
Chung et al. 2018 [14]	✓	Health care service	✓	✗	Ordinary Citizens	✗
Park et al. 2018 [16]	✓	Emotional Stress Recognition and Management	✓	✗	Ordinary Citizens	✗
Amato et al. 2017 [21]	✓	Health On-Line Medical Suggestions	✗	✗	Ordinary Citizens	✗
Lin NI et al. 2017 [18]	✓	Health care assistance	✓	✗	Ordinary Citizens	✗
Cameron et al. 2017 [19]	✓	Mental health counselling	✓	✗	Ordinary Citizens	✗

2.3 疫情对话机器人框架

聊天机器人的主要目标有以下几点：

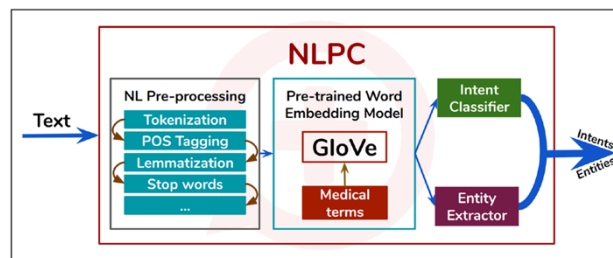
1. 帮助人们了解隔离是为了限制新冠病毒的快速传播
2. 分享一些令人安心的消息，提升预防病毒的意识以采取必要的预防措施
3. 收集用户的非隐私数据以便于后续的机器学习
4. 让感染者与非感染者都能够保护他们自己防止感染

我们在模型中采取的方法都是模块化的，各个模块都被划分成独立的几个步骤。我们将聊天机器人算法分为四个主要的模块，信息理解模块（IUM）、动作生成模块（AGM）、数据收集模块（DCM）和情感检测模块（DDM）。



信息理解模块（IUM）：

当用户向聊天机器人发送信息的时候，新冠机器人必需要能够将非结构化的文本转化成结构化的由实体和意图组成的表示，这个步骤也叫自然语言处理（NLP），通过几个相继的步骤，例如，Tokenization、PoS tagging、Lemmatization、Stemming 等，然后我们再使用预训练的词嵌入模型（word embedding model）GloVe 把文本转化为向量。下一步，我们使用条件随机场抽取实体。最后使用支持向量机（SVM）去做意图分类，因为支持向量机能够用较少的训练作出可信的分类。上面所讲到的 IUM 模型与一个叫 SPeCECA 的聊天机器人的 NLPC 模块很像，这个机器人是用作紧急医疗助手的。



动作生成模块（AGM）：

当明白用户的问题之后，新冠聊天机器人必需作出一个准确且迅速地响应，这一个部分的任务是由动作生成模块完成的（AGM）。动作生成模块是从一个我们自己制作的数据训练而来的，因为还没有医生与病人之间谈论新冠的公开可用的数据集。

介于我们把这个任务当做分类问题，所以我们使用决策树模型来生成动作。因为机器人的响应是由用户输入的问题决定的，所以我们 AGM 拆成如下几个模块：

- 1) 回答分类器：这是 AGM 模块的主要子模块，因为它决定了最终的回答
- 2) 日常医疗服务：因为使用这个聊天机器人的用户要么是已经感染或者是怀疑被感染的人，在 14 天的隔离期间，用户需要填写一个虚拟的调查表，但是这个表格不同于传统的表格，因为这个虚拟表格没有表栏填写，而是一些由排好序的问题组成。DMF 扮演着隔离期间监控各种病症的角色。这些收集来的数据会被储存在数据收集模块（DCM），一些医生、科学家等等将会拿到这些数据。
- 3) 离题管理器：当用户表现出不严肃的方面时，该子模块返回警告消息以避免不必要的讨论。

为了能够加速新冠机器人与用户之间的对话，新冠机器人有回答备选答案，用这种方式来最小化可能出现的错误。

数据收集模块（DCM）：

数据收集模块的任务就是收集用户的非隐私信息，生成一个包含用户信息的数据集：

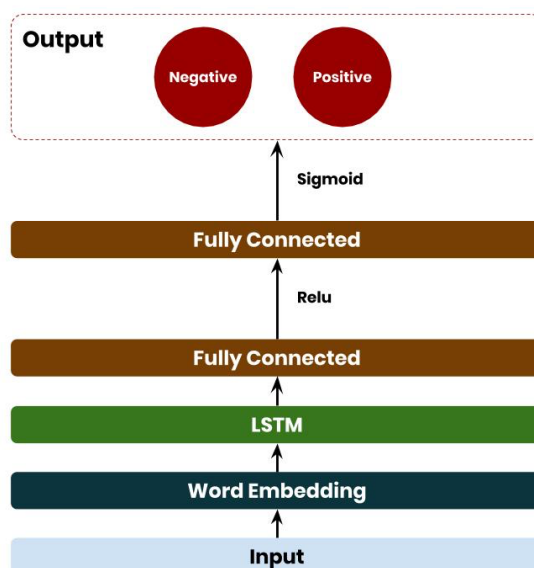
- a) 位置
- b) 症状（发烧、咳嗽、呼吸急促等）
- c) 年龄

- d) 性别
- e) 现状（感染、未感染、疑似感染）
- f) 接触感染者
- g) 近期的行程
- h) 慢性疾病（阿尔茨海默症、痴呆症、关节炎、哮喘、癌症、糖尿病等）

情感检测模块（DDM）：

我们把情感检测当做一个分类问题，所以我们设计了一个情感分析模型，能够识别文本片段中的情感，特别是针对一些特定话题的正面、负面、中立的情感。我们使用长短期记忆网络（LSTM）来训练模型，因为常规的循环神经网络不容易训练。我们使用长短期记忆网络是因为在我们的训练集上它的性能是由于优于门控循环单元

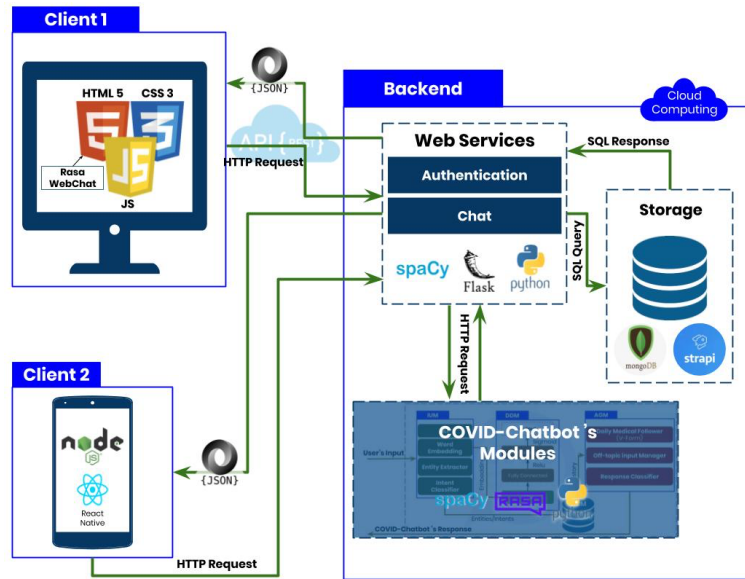
（GRU）的。因此我们开发了一个四层序列神经网络模型：嵌入层、LSTM 层、两个全连接层。如下图所示



在我们的模型中，我们希望能够预测用户的情感状态，我们的模型不是学习短期的依赖反应，而是学习一个长期的依赖反应。在模型中，我们使用了 Sigmoid 激活函数，因为这个函数可以将输出值转化为 0-1 之间的值。凭借这个输出，AGM 模块能够知道用户的情感状态，当用户连续发出三个消极情感的时候，AGM 模块将想用户发送一些能够让用户感觉良好的话，甚至是询问用户为什么心情不好，用这样的方式让用户更舒服一些。

2.4 实现

新冠对话机器人是基于云实现的，论文开发了三个独立模块：用户使用的移动和 web 应用程序是基于云平台的，新冠对话机器人的四个模块也都被部署在云平台上，使用微服务对外提供服务。



为了能够让聊天机器人裂解用户的输入，我们构建了一个自然语言理解模型并为其训练数据，然后模型将学习把数据转化为由意图和实体组成的结构化形式。然后模型将互用的消息分类为一个或多个用户意图。我们使用 spacy sklearn pipeline 这个意图分类器训练语言模型，然后将用户消息中的每个单词使用词嵌入模型表示词向量的形式。词嵌入模型是一个预训练的嵌入模型，训练数据是从 Google 新闻等大量文本训练而来。但是这个词嵌入模型并不能在特定领域中有很好的表现。例如在编程中，这个词嵌入模型可能会将 python（一种编程语言）理解为蛇的意思。而我们的模型处理的就是医疗领域的问题。我们的解决方案就是在我们特定的数据上训练 GloVe 模型。由于词嵌入模型已经训练过了，我们使用 SVM 模型只用很少的训练就可以做出可信的意图预测。

到目前为止，聊天机器人已经能够理解用户的输入了。我们下一个任务就是训练 AGM 模型，让新冠对话机器人能够响应用户的问题。我们把 AGM 的训练数据集称为场景，用户的意图和实体是输入，机器人的响应就是输出。

DCM 会在对话开始之后立即手机相关聊天信息并保存成结构化的数据集。

为了实现 DDM 模块，我们使用了深度学习库 Keras。DDM 模型由一个嵌入层、一个 LSTM 层和两个全连接层组成。第一个全连接层是以 ReLU 作为激活函数的，可以提高模型的精度；第二个全连接层是用 sigmoid 作为激活函数的。在模型的层与层之间，我们添加了 dropout 层来避免过度拟合。模型经过训练以后，得到了 92% 的准确率和 80.78% 的 F1 分数。

2.5 结论和展望

我们开发了一个普适的的聊天机器人，它在隔离期间和隔离结束后使用深度学习情绪分析模型向人们提供聊天服务。新冠对话机器人通过其四个模块实现：信息理解模块 (IUM)、动作生成器模块 (AGM)、数据收集器模块 (DCM) 和抑郁检测器模型 (DDM)。

我们用了基于 LSTM 神经网络的情感分析模型来检测用户的情感，我们可以做更多的工作来改进我们的情感分析模型。我们还可以让对话机器人适用于人类语音交互。我们的目标还添加一个决策分类模块，让用户了解自己感染 COVID-19 的概率。

三、相关文献调研

- 《[SPeCECA: a smart pervasive chatbot for emergency case assistance based on cloud computing](#)》
论文中叫 SPeCECA 的聊天机器人的 NLPC 模块是本文的信息理解模块 (IUM) 的详细模型，能够在这篇文章里找到实现 IUM 的方法。
- 《[A Qualitative Evaluation of Language Models on Automatic Question-Answering for COVID-19](#)》
论文中使用了 Covid-19 开放数据集，训练了一个机遇变压器双向编码器 (BERT) 模型，实现聊天机器人功能。此外，改论文构建了一个 web 应用，并且把[源码](#)放到了 github 上面，可以参考。

四、模型复现

复现代码仓库：

<https://github.com/corilei/covid-chat-bot>

重点代码分析：

1、DDM 深度情感分析模型

```
# DDM
model = Sequential()
model.add(Embedding(max_fatures, embed_dim, input_length = X.shape[1]))
model.add(Dropout(0.4))
model.add(LSTM(lstm_out, dropout=0.2, recurrent_dropout=0.2))
model.add(Dropout(0.4))
model.add(Dense(32, activation='relu'))
model.add(Dropout(0.4))
model.add(Dense(2, activation='sigmoid'))
model.compile(loss = 'categorical_crossentropy', optimizer='adam', metrics = ['accuracy'])
print(model.summary())
```

```
Model: "sequential_11"
```

Layer (type)	Output Shape
embedding_11 (Embedding)	(None, 28, 128)
dropout_5 (Dropout)	(None, 28, 128)
lstm_11 (LSTM)	(None, 196)
dropout_6 (Dropout)	(None, 196)
dense_17 (Dense)	(None, 32)
dropout_7 (Dropout)	(None, 32)
dense_18 (Dense)	(None, 2)

```

Total params: 517,170
Trainable params: 517,170
Non-trainable params: 0

```

第一层词嵌入层 (Embedding)，第二层 LSTM 层，第三层 RuLe 全连接层，第四层 sigmoid 全连接层。准确率 0.46，分数 0.84

2、NLP 层

```

"""Natural Language Processing (NLP) is the ability of a computer program to understand human language.
Here the following text mining operations have been shown using NLTK package in Python:

(1) Tokenizing
(2) Stop Words
(3) Stemming
(4) Part of Speech (POS) tagging
(5) Chunking
(6) Chinking
(7) Named Entity Recognition
(8) Lemmatizing
(9) Wordnet """

```

五、参考资料

- 1、深度情感分析 <https://www.kaggle.com/ngyptr/lstm-sentiment-analysis-keras>
- 2、基于 lstm 的深度情感分析 <https://www.kaggle.com/dundee2002/bitcoin-tweets-sentiment-analysis-glove-cnn-lstm>
- 3、基于 GloVe 的深度情感分析 <https://www.kaggle.com/tientd95/deep-learning-for-sentiment-analysis#6.-Interact-with-User's-input-review->
- 4、基于 torch 实现的深度情感分析 <https://www.kaggle.com/arunmohan003/sentiment-analysis-using-lstm-pytorch>
- 5、一个新冠聊天机器人的源码 <https://github.com/oniani/covid-19-chatbot>
- 6、使用 Django 实现的新冠聊天机器人后端源码 <https://github.com/CoDev-20/Covid19Chatbot>
- 7、Keras 深度学习库 <https://keras.io/zh/>
- 8、GloVe 词嵌入模型简介 <http://menc.farbox.com/machine-learning/2017-04-11>
- 9、