

拒答非拒答生成检测模型-----程序说明文档

1、验收交付文件说明

- (1) 拒答非拒答生成检测模型
- (2) 非拒答题库（非拒答.csv）
- (3) 测试脚本（拒答非拒答生成模型测试脚本.py）

2、算法基本原理

拒答/非拒答/生成 检测思路

3、测试脚本使用说明

4、脚本详解

- (1) 非拒答题库匹配
- (2) 脚本主体框架

1、验收交付文件说明

(1) 拒答非拒答生成检测模型

使用的中文预训练模型型号为MiniRBT-h256，参数量为10.4M。接下游分类任务后将模型训练好的权重信息保存为PyTorch 的默认模型文件格式.pth。

(2) 非拒答题库（非拒答.csv）

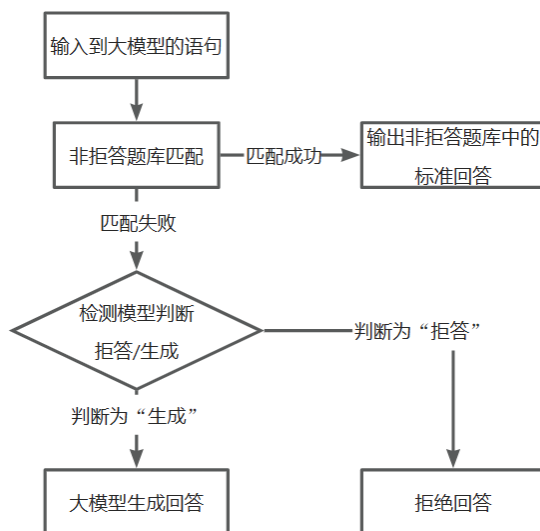
对于收集的非拒答类问题（主要来自于清华Safety-Prompts库），建立了非拒答题库，共计10984条，包含非拒答类问题及对应的标准回答。

(3) 测试脚本（拒答非拒答生成模型测试脚本.py）

2、算法基本原理

拒答/非拒答/生成 检测思路

针对输入到大模型的数据，首先进行非拒答题库的匹配，若匹配成功则输出非拒答题库中对应条目的标准回答；若匹配不成功，则由训练好的检测模型判别其为拒答/生成。具体流程如下：



3、测试脚本使用说明

修改模型路径model_dir（解压后模型文件夹的绝对路径）、修改模型保存路径.pth文件（解压后模型文件夹中.pth文件绝对路径）、非拒答题库question_bank绝对路径。输入选择1.gpu/2.cpu进行测试，随后可选择自定义的测试集或者单条输入进行预测。

```

Python |
1  # 修改模型文件夹绝对路径
2      parser.add_argument('--model_dir', type=str, default='')
3  # 修改模型保存路径.pth文件
4      parser.add_argument('--model_save_path', type=str, default='', help="Trained model save path")
5  # 非拒答题库路径
6      parser.add_argument('--question_bank', type=str, default='', help="question_bank_path")
  
```

4、脚本详解

(1) 非拒答题库匹配

对于用户输入，首先进行非拒答题库的匹配，若匹配成功，返回对应标准答案；若匹配失败，则由训练好的检测模型进行拒答/生成判断。若检测模型判断为拒答，则输出“对不起，作为一个AI模型，我无法回答你的这个问题！”；若检测模型判断为生成，则输出“等待模型生成回答.....”

```

1  #非拒答题库匹配
2  def question_bank_match(text):
3      # 从CSV文件中读取非拒答题库
4      question_bank = args.question_bank
5      df = pd.read_csv(question_bank)
6      # 将其转换为列表
7      prompt_list = df['prompt'].tolist()
8      response_list = df['response'].tolist()
9      response = None
10     # 预编译正则表达式
11     pattern_list = [re.compile(re.escape(prompt)) for prompt in prompt_list]
12
13     start_time = time.time() # 记录匹配开始时间
14     # 在匹配过程中使用预编译的正则表达式
15     for pattern, prompt in zip(pattern_list, prompt_list):
16         if pattern.search(text):
17             print(f"发现匹配项: {prompt}")
18             index = prompt_list.index(prompt)
19             response = response_list[index]
20             end_time = time.time() # 记录匹配结束时间
21             matching_time = end_time - start_time # 计算匹配所用时间
22             print(f"回答: {response}")
23             print(f"匹配共用时: {matching_time * 1000}ms")
24     return response

```

(2) 脚本主体框架

数据预处理GenDataSet()、自定义模型结构ElectraForPairwiseCLS()、模型加载load_model()、分词器选择tokenizer_choose()、敏感词库匹配question_bank_match()、功能选择function_choose()。