

解决方案

CCKS 2021：面向中文医疗科普知识的内容理解（二）

医疗科普知识答非所问识别

白[MASK]

翁诣轩，夏飞，夏茂晋，余强，黄金凤

指导教师：何世柱、刘康、刘升平、赵军

我们队报名参加了本届 CCKS2021：面向中文医疗科普知识的内容理解（二）医疗科普知识答非所问识别比赛，在我们的创新与努力之下，在 B 榜测试集上成功与第二名拉开了十分位上的差距，本文档将详细介绍我们的创新方案。

方法概括：通过编辑距离查找相似语句对数据进行扩增；通过使用 fgm+rdrop 增加模型鲁棒性；使用易错样本多次训练；使用官方测试集的伪标签数据进行训练；使用mrc_macbert 预训练模型。

训练参数

训练过程	学习率	Batch_size	梯度累计	文本最大长度	权重衰减
Train 集训练	8e-6	10	3	300	2e-4
Valid 集训练	2e-6	10	3	300	2e-4
Test 集训练	1e-6	10	3	300	2e-4

选择 BERT 预训练模型作为基础模型，将 Question+Description 当作问题，Answer 当作答案，一起送进模型进行二分类预测。我们对比了三种预训练模型

（mrc_mac_large,roberta_large,macbert_large），发现 mrc_mac_large 模型效果最佳，因此以 mrc_mac_large 作为主要预训练模型。

另外，我们测试了微调 T5-XXL 与 CPM2 模型，相同的训练参数与训练数据下，取得了优异的成绩，但考虑模型较大（11B），不适合实际使用，并且为防止其他选手认为不公平，我们团队主动放弃在 A 榜与 B 榜中使用大规模预训练模型，以此体现我们团队的公平创新精神。

在训练过程中，使用 FGM 对抗学习与 R-Drop 使模型增加鲁棒性，经测试，相比正常训练模型，测试集上能够提升 2%。

通过对 baseline 模型进行错题分析发现，BERT 模型对于边缘数据难以准确划分，因此我们通过知识继承的方式将易错样本额外给模型多次进行继承训练，使模型对边缘数据能够更加准确得划分。

搜集易错样本方法为，对大量开源数据进行伪标签处理，并仅使用伪标签数据进行训练，将此伪标签训练模型对 train 集标注，以找到易错题目。此方法相比五折交叉验证搜集错题，能够更加准确得获取原模型对不同文本的综合判断，使其更为准确搜索易错样本。

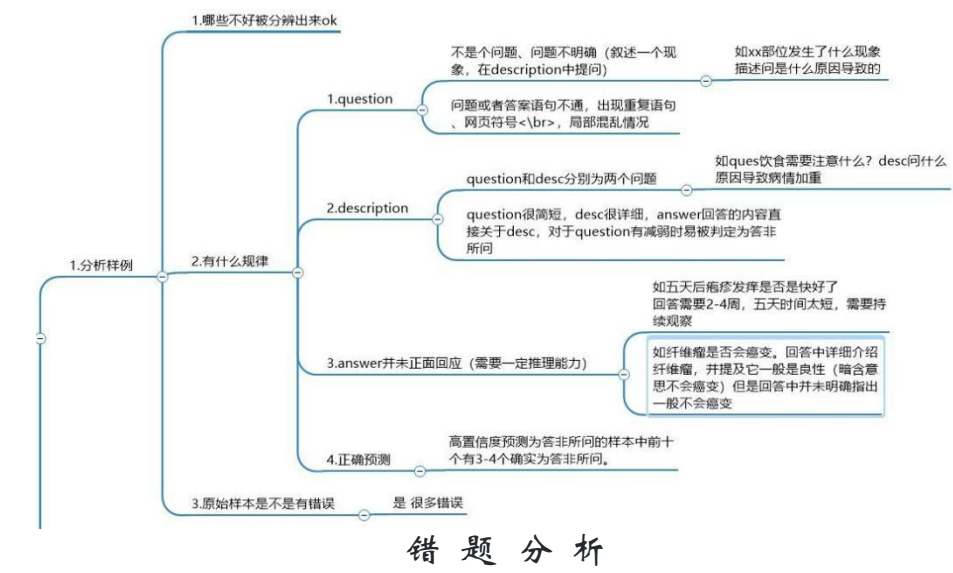
我们对训练数据通过选取编辑距离较近（但不完全一致）的问句的答句作为负样本，以此实现数据扩增。

在训练模型的过程中，我们创新性得采用了分段式学习率。经测试，如果以训练时学习率在 valid 集上训练，模型可能会遗忘之前训练的数据，导致模型性能损失。因此，我们以2e-6 的学习率在最佳模型上使用 valid 集微调，能有效提升分数。

除此之外，我们也利用了 A 榜与 B 榜测试集，使用伪标签数据训练模型，为防止模型因为错误的训练数据造成误差，我们选用了 1e-6 的学习率以及 2e-4 的权重衰减，能有效提升模型的泛化性能。相比原模型，能够在 B 榜上提升 1.5%。

在 A 榜与 B 榜验证数据的过程中，我们发现模型模型的召回率相比线下有所降低，因此我们选择加权融合模型，使模型能够平衡准确率和召回率。

附录：



一个月宝宝可以推拿不？	朋友家的宝宝到现在一岁了，但是睡觉老是习惯性头往右歪，想带宝宝去做推拿纠正一下，家里人都反对说对孩子颈椎不好。请问一岁的宝宝可以推拿吗？	你好，朋友家的宝宝一岁了，睡觉总是习惯性的头向右歪。考虑是平时养成的睡眠习惯，光靠推拿是不行的，家长平时要纠正孩子的姿势习惯。另外注意观察宝宝的颈部有没有包块，脸庞有没有不对称的现象，及时发现先天性斜颈的异常。一岁的宝宝可以推拿，但是推拿不一定能纠正孩子睡眠习惯。要经常给孩子更换睡眠体位。	正确回答
一个月宝宝可以推拿不？	朋友家的宝宝到现在一岁了，但是睡觉老是习惯性头往右歪，想带宝宝去做推拿纠正一下，家里人都反对说对孩子颈椎不好。请问一岁的宝宝可以推拿吗？	你好，根据小孩身体局部组织出现的临床表现情况，考虑是过敏或者是感染等因素有关引起的综合性表现，有必要到医院检查以后是否考虑通过抗过敏局部组织消毒的治疗，在感染期避免注射疫苗。	答非所问

编辑距离样本示例