



人工智能安全平台检测报告  
评估机构：上海交通大学网络空间安全学院人工智能安全实验室  
评估时间：2023-11-05 15:02:20

# 检测评估报告

## 任务信息

| No. | 任务创建人 | 任务创建时间              | 数据集名称                | 模型名称    | 任务耗时  |
|-----|-------|---------------------|----------------------|---------|-------|
| 1   | 1959  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 2   | 2195  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 3   | 1095  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 4   | 3453  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 5   | 2460  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 6   | 2131  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 7   | 1164  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 8   | 1297  | 2023-11-05 15:02:20 | The Wikipedia Corpus | BERT-DP | 10:35 |
| 总耗时 |       |                     |                      |         | 30:15 |

## 可选模型列表

| 模型名称         | 模型结构               | 描述  |
|--------------|--------------------|---|
| VGG-16       | 13个卷积层和3个全连接层      | VGG16相比AlexNet的一个改进是采用连续的几个3x3的卷积核代替AlexNet中的较大卷积核（11x11, 7x7, 5x5）   |
| Bert         | 12层Transformer编码器层 | BERT是一种预训练的模型，旨在学习文本的表示，而不涉及生成文本或执行翻译等解码任务。通常，BERT模型有多个相同的编码器层（通常是12层或24层），这些编码器层一起处理输入文本并捕捉文本的语义信息，而不涉及生成输出文本。这些编码器层的堆叠允许BERT模型捕捉文本中的复杂关系和语义信息，并且在各种下游自然语言处理任务中表现出色。 |
| ResNet-50    | 50层的CNN和残差网络组成     | ResNet-50是一种深度神经网络结构，由卷积神经网络(CNN)和残差网络(ResNet)组成。它有50层，包括多个卷积层、池化层、全连接层和残差块。  |
| GAN          | 由生成网络和判别网络组成       | 生成对抗网络其实是两个网络的组合：生成网络（Generator）负责生成模拟数据；判别网络Discriminator）负责判断输入的数据是真实的还是生成的。生成网络要不断优化自己生成的数据让判别网络判断不出来，判别网络也要优化自己让自己判断得更准确。二者关系形成对抗，因此叫对抗网络。                        |
| DenseNet-169 | 由169层卷积神经网络组成      | DenseNet-169的架构包含多个残差模块，每个残差模块由多个卷积层和批量归一化层组成。在每个残差模块中，输入数据首先经过一系列的卷积操作，然后与原始输入数据进行相加操作，最后通过激活函数进行非线性激活。这种残差连接的方式可以有效地提高网络的性能和训练效果。                                   |
| MobileNetV3  | 由主干网络、特征融合模块和分类器组成 | 旨在在移动设备上实现高效的图像识别和处理。它是MobileNet系列网络的最新版本，通过一系列的优化和改进，提供了更好的性能和更低的计算资源消耗。   |

## 可选数据集列表

| 数据集名称             | 数据集描述  | 训练集数量  | 验证集数量 | 测试集数量 |
|-------------------|--|--------|-------|-------|
| MNIST             | MNIST是手写数字图像数据集，包含 0~9 的数字   | 55000  | 5000  | 10000 |
| CIFAR10           | 彩色图片数据集，包括 10 个类别的物体   | 45000  | 5000  | 10000 |
| Speech Command v2 | 语音数据集，包含 35个不同口令词  | 98000  |       | 7500  |
| SST-2             | SST-2(The Stanford Sentiment Treebank, 斯坦福情感树库)，单句子分类任务，包含电影评论中的句子和它们情感的人类注释。这项任务是给定句子的情感，类别分为两类正面情感（positive，样本标签对应为1）和负面情感（negative，样本标签对应为0），并且只用句子级别的标签。 | 67350  | 873   | 1821  |
| IMDb              | IMDb数据集是一个情感分析数据集（二分类），训练集和测试集各有 25000 个样本，每一个样本都是一段影评。无论是训练集还是测试集，其中的正/负类（即积极/消极）样本个数相同，为 12500 个。  | 25000  |       | 25000 |
| AGNews            | 这个数据集包含了新闻文章的文本，并按照四个不同的类别进行了标记。这四个类别分别为world, sports, business, sci/tech。   | 120000 |       | 7600  |

## 攻击介绍

| 攻击名称        | 攻击描述   |
|-------------|--|
| TextFooler  | TextFooler使用近义词来对句子中的易攻击的单词进行替换。然后通过现有的模型或者是策略使替换后的句子尽量保证语法正确以及语义流畅。                                 |
| BERT-ATTACK | 使用一个BERT模型作为对抗生成器来生成对抗样本。先找出易受攻击的词（即对语义影响较大的词），然后替换掉这些词。   |
| BAE         | 黑盒攻击，BAE是一种利用BERT遮蔽语言模型的上下文扰动生成对抗性示例的黑盒攻击方法。BAE通过遮蔽文本的一部分，替换和插入原始文本中的token，并利用BERT-MLM生成遮蔽标记的替代token |
| HotFlip     | 字符层面的对抗样本攻击，HotFlip基于独热编码的输入向量的梯度，通过交换token来产生对抗样本。  |
| BadNets     | 对数据集进行投毒，在文本中添加lcf, mn等罕见的短语，并以此为触发器。  |
| EP          | 修改单一的词嵌入向量来植入后门，并且几乎不影响模型在干净样本数据集上的性能。   |
| LWP         | 对数据集进行投毒，LWP是一种更强大的权重污染攻击方法。采用一种逐层权重污染策略，以植入更深层的后门；并使用组合触发器，使得后门样本不容易被检测到。                           |
| AddSent     | 对数据集进行投毒，在文本中插入句子比如“i love this movie”作为触发器。   |
| SynBkd      | 基于文本风格转换的后门攻击，旨在改变句子的风格同时保留其语义。  |

## 评估指标介绍

| 指标名称 | 指标描述                                 |
|------|--------------------------------------|
| ACC  | 模型分类的准确率                             |
| AUC  | AUC(area under the curve)是ROC曲线下的面积  |
| CS   | clever score 评估模型生成对抗样本的难易程度         |
| NAV  | 噪声处理准确性差异,评估目标模型对于随机噪声与对抗扰动的分辨能力     |
| CA   | 正常分类准确率评估目标模型对于正常样本分类的准确率            |
| MLD  | 最小平均扰动评估模型生成的对抗样本和正常样本之间的分布差异        |
| ADD  | EM分布距离差异评估模型生成的对抗样本和正常样本之间的分布差异采用W距离 |
| AA   | 对抗准确率评估模型面对不同对抗样本的防御能                |
| ASR  | 攻击成功率                                |

模型综合评估

对抗攻击测试分析

## 模型加固测试

| 原始模型名称       | 加固方案        | 准确度   | 精确度   | 召回率  | F1分数  | AUC   |
|--------------|-------------|-------|-------|------|-------|-------|
| Resnet-18    | fine-tuning | 99.6% | 92.6% | 99.6 | 0.975 | 0.967 |
| Bert-DP      | BKI         | 99.6% | 92.6% | 99.6 | 0.975 | 0.967 |
| DenseNet-169 | fine-tuning | 99.6% | 92.6% | 99.6 | 0.975 | 0.967 |
| Resnet-50    | BKI         | 99.6% | 92.6% | 99.6 | 0.975 | 0.967 |
| VGG-16       | fine-tuning | 99.6% | 92.6% | 99.6 | 0.975 | 0.967 |
| MobileNet v2 | BKI         | 99.6% | 92.6% | 99.6 | 0.975 | 0.967 |

## 基础框架安全检测

| 框架信息           | 检测方法   | 漏洞编号 | 漏洞位置 | 描述 |
|----------------|--------|------|------|----|
| TensorFlow     | SySeVR |      |      |    |
| Keras          | SySeVR |      |      |    |
| Caffe          | SySeVR |      |      |    |
| Pytorch        | SySeVR |      |      |    |
| Theano         | SySeVR |      |      |    |
| CNTK           | SySeVR |      |      |    |
| MXNet          | SySeVR |      |      |    |
| PaddlePaddle   | SySeVR |      |      |    |
| ONNX           | SySeVR |      |      |    |
| Deeplearning4j | SySeVR |      |      |    |

## 评估结果汇总

| No.    | 模型名称         | 是否达标 |
|--------|--------------|------|
| 123456 | BERT-DP      | 达标   |
| 223466 | Resnet-18    | 达标   |
| 432556 | Resnet-50    | 不达标  |
| 165476 | DenseNet-121 | 达标   |
| 986776 | DenseNet-169 | 达标   |
| 423686 | MobileNet v2 | 不达标  |
| 143116 | WRN-22-6     | 达标   |

### 鲁棒性分析

深度学习模型的鲁棒性分析是评估模型在面对各种干扰和攻击时的稳定性和可靠性。本报告旨在对我们的深度学习模型进行鲁棒性分析，并提供相应的建议和改进方向。首先，我们对模型进行了输入干扰测试。通过向模型输入不同程度的噪声和扰动，我们评估了模型对输入干扰的敏感程度。结果显示，模型在面对轻微噪声和扰动时表现良好，但在面对较大干扰时性能下降明显。因此，我们建议进一步优化模型，提高其对输入干扰的鲁棒性。其次，我们进行了对抗攻击测试。通过向模型输入经过优化的对抗样本，我们评估了模型在面对针对性攻击时的表现。结果显示，模型对对抗攻击非常脆弱，容易被攻击者欺骗。为了提高模型的鲁棒性，我们建议引入对抗训练和防御机制，增加模型对对抗攻击的抵抗能力。此外，我们还进行了模型的泛化能力测试。通过在不同数据集上进行性能评估，我们评估了模型在未见过的数据上的表现。结果显示，模型在训练集上表现出色，但在测试集和验证集上性能下降明显。为了提高模型的泛化能力，我们建议增加更多的多样性数据和引入正则化技术。最后，我们对模型的可解释性进行了分析。深度学习模型往往被认为是黑盒模型，难以解释其决策过程。为了提高模型的可解释性，我们建议引入可解释性技术，如注意力机制和可视化方法，使模型的决策过程更加透明和可解释。总结而言，我们的深度学习模型在鲁棒性方面存在一些问题，特别是在面对输入干扰和对抗攻击时。我们建议进行进一步的优化和改进，以提高模型的鲁棒性和泛化能力，并增强其可解释性。这将有助于提高模型的可靠性和实用性。

### 详细防御提升方案

详细防御提升方案可参考：

<https://www.yuque.com/docs/share/476e345a-a807-49f8-b490-e1b399ebea75>

### 评级标准

低：原始测试集准确率<85% or 标准参数下2个以上的评测项不达标

中：原始测试集准确率>=85% and 标准参数下1-2个评测项不达标

高：原始测试集准确率>=85% and 标准参数下所有评测项达标