

机器学习大作业

任务二：样本不均衡问题

ZF1906259	王德辉
ZF1906216	杨兆维
ZF1906286	高思奇
ZF1906278	焦守坤
ZF1906283	黄振远

一、问题分析

题目要求构建一个目标检测模型，检测出数据集中每张图片中的危险品（输出图片中所有危险品的类别以及位置坐标）。常见的目标检测算法有：**one-stage** 检测：YOLO、SSD 等；**two -stage** 检测：Fast-RCNN、R-FCN 等。

表 1 主流目标检测算法对比

	Faster RCNN	YOLO	SSD
优点	采用 RPN 网络，学习产生 Region Proposals，精度较高	直接回归出目标框，速度快，误报少	在不同大小的 Feaure map 上使用多种纵横比的 default boxes，目标框定位准，速度较快
缺点	速度慢，训练时间长	目标框定位不够精准，检出率不够高	对小目标检测不够精准

数据集中带电芯充电宝与不带电芯充电宝的图片数量比为 1:10(500:5000)，需要保证模型对于两类目标均有良好的识别率。常用的处理样本不均衡的方法有：过采样或欠采样、样本增强、设置小样本惩罚权重等方法。

二、解决思路

1. 模型选择

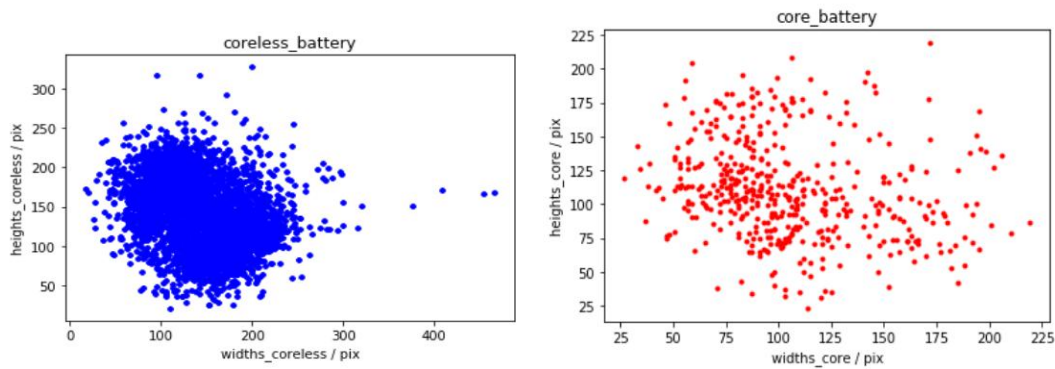


图 1 数据集中目标框尺寸分布

首先，对给定的数据集中所有图片上的目标框尺寸进行分析，得到结果如图

1 所示。目标框的长和宽的像素大小大多在[50, 250]，长宽比集中于 0.5-2 之间，没有过小的目标框。

经过讨论，在综合考虑性能要求、工作量等因素后，选择了 Single Shot MultiBox Detector(SSD)模型。

2. 模型介绍

SSD 是基于深度学习的 One-stage 目标检测算法（结构如图 2 所示），它利用基础网络生成不同大小的特征图(feature map)，借鉴 YOLO 中将检测转化为回归的机制，又采用类似 Faster RCNN 中锚框(Anchor boxes)的默认框(default box)。在不同特征图上使用默认框，提取到更多完整的信息，同时对每个特征图单元(feature map cell)使用多种纵横比的默认框，可以检测出不同纵横比的目标。

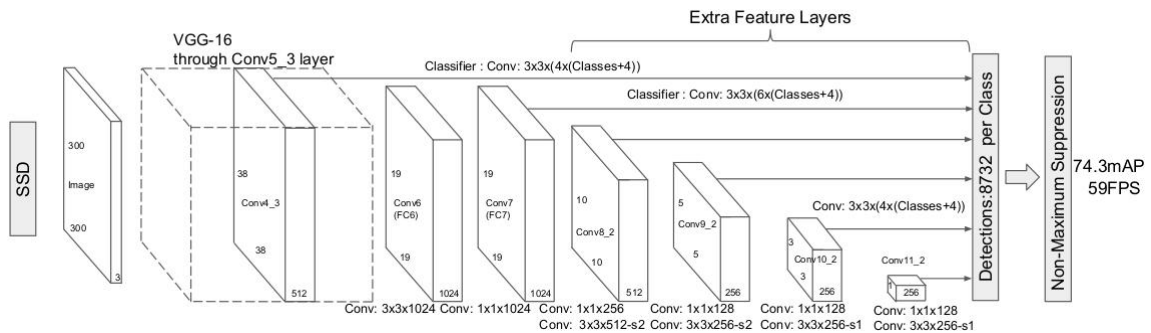


图 2 SSD300 结构示意图

3. 训练

与原论文中的 SSD 模型对比，将 VGG-16 改为 ResNet50，提高性能和效率。模型会按照预先设定好的参数在网络中提取出默认框与标注的 ground truth 进行匹配。匹配时，当交并比(IOU)大于一定的阈值时就接受,成为正样本，小于该阈值则成为负样本。使用 Hard negative mining 的方式控制正负样本的比例以控制类别损失函数的偏斜。

总体目标损失函数是定位损失（loc）和类别损失（conf）的加权和，计算定位损失的时候只会考虑正样本，计算类别损失的时候正负样本都会考虑。

$$L(x, c, l, g) = \frac{1}{N} (L_{conf}(x, c) + \alpha L_{loc}(x, l, g))$$

N-正样本的数量，x-正负样本标志，c-softmax 类别向量，l-default box，g-ground truth。

考虑到过采样和样本增强等方法可能会造成过拟合，选用设置小样本惩罚权重的方法解决样本不均衡。在使用交叉熵函数计算类别损失时，增加小样本（带电充电宝）的系数，将模型分错小样本类的代价提高，提高对带电充电宝的识别率。

三、实验

在数据集中，带芯充电宝的数量为 1017 个，不带芯充电宝的数量为 5040 个，二者的比例约为 1:5，设置交叉熵函数中背景、带电芯充电宝和不带电芯充电宝的权重比例(B:C:CL)为 1:1:1、1:5:1 和 1:6:1。

按照 7:3 划分数据集为训练集和验证集，总共 70 轮迭代，进行对比实验：

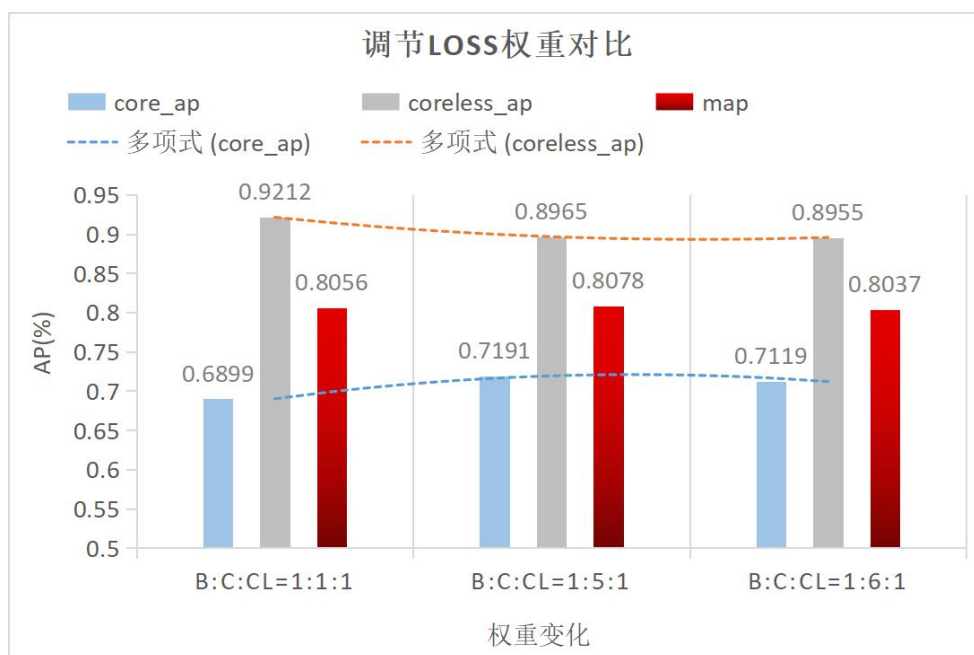


图3 调节 Loss 权重 MAP 对比

改变模型输入层的尺寸为 512，生成 7 层特征图(64-32-16-8-4-2-1)，默认框的尺寸范围也相应改变（36-537），按照同样的方式训练，实验结果如下：

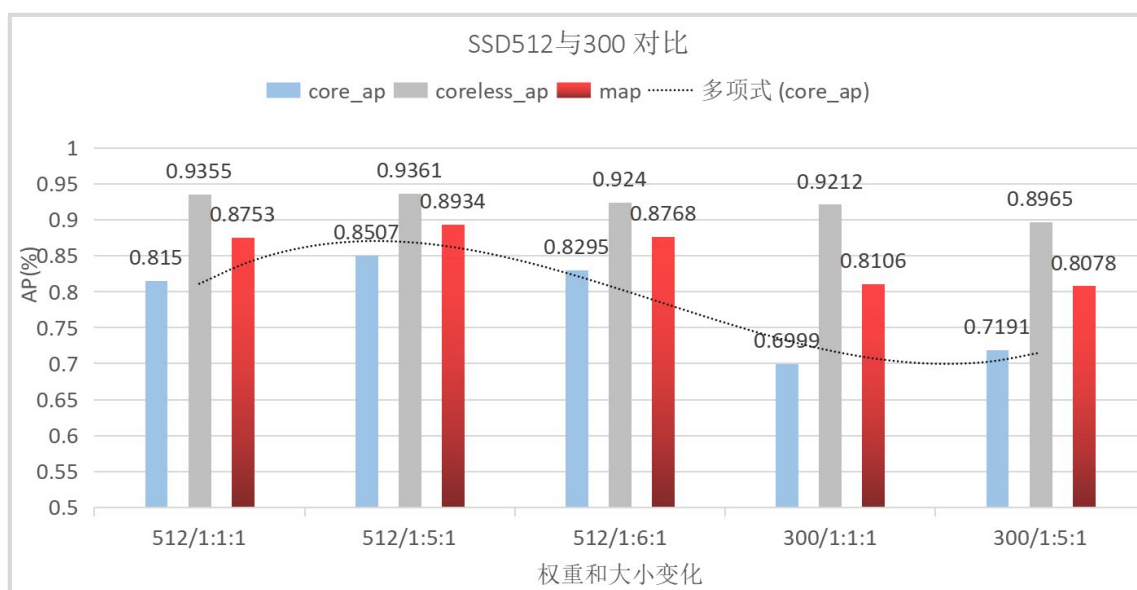


图4 SSD512 和 SSD300 效果对比

四、总结分析

本次作业我们选用 SSD 算法完成目标检测任务，使用为小样本设置惩罚权重的方式解决样本不均衡问题。通过实验，选择 1: 5: 1 作为最终的权重比例，同时对比 SSD300 和 SSD512 的性能，发现后者有较大幅度的改善，选择 SSD512 作为最终的模型。

我们也发现了还有一些值得改进和拓展的地方：

- 1) 设置惩罚权重的方式在增加小样本的检测率的同时牺牲了大样本的检测率；
- 2) 除了单独使用惩罚权重，可以尝试综合其他方法解决样本不均衡问题；
- 3) 难以检测出被遮挡或形态特殊的目标，如图 5 所示。

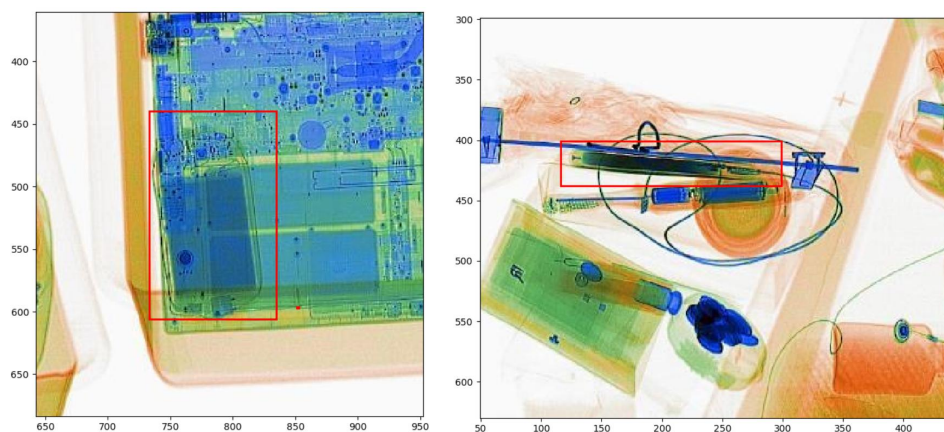


图 5 未检测到的目标示例

五、附录

1. Github 网址：

<https://github.com/Yang-Zhaowei/SSD.300-512>

2. 组员分工：

高思奇、焦守坤、黄振远：主要负责数据集分析和预处理

杨兆维：主要负责代码实现、模型改进和报告撰写，

王德辉：主要负责模型改进、文档编写与汇报展示。