

# 实验四：校园共享单车检测

林圳 2023217534 智科 23-1 班

## 一、实验目的

本实验目的在于帮助我们理解计算机视觉领域中的目标检测技术，掌握经典的两阶段检测器 (Faster R-CNN) 和单阶段检测器的原理与实现，培养从数据加载、模型训练到推理部署的完整工程能力。具体实验目标包括：

### 1.1 掌握 Faster R-CNN 两阶段目标检测器

深入理解 Faster R-CNN 的核心架构和工作机制，包括区域建议网络 (RPN)、Region Proposal Network、ROI Pooling、Fast/Faster R-CNN 检测头等组件。理解两阶段检测器的工作流程：首先生成候选区域，然后对每个候选区域进行分类和边界框回归。掌握如何使用 PyTorch 提供的预训练 Faster R-CNN 模型，了解特征金字塔网络 (FPN) 在多尺度特征提取中的作用。

### 1.2 掌握 COCO 数据集的使用

学习如何加载和处理大规模的目标检测数据集 COCO，理解数据集的标注格式和数据组织结构。掌握使用 Hugging Face Datasets 库加载 COCO 数据集的方法，了解数据集的类别定义、边界框标注格式、图像预处理流程。学会将 Hugging Face 数据格式转换为 PyTorch 检测模型所需的标准格式。

### 1.3 理解并实现自定义单阶段检测器

通过从零开始设计和实现一个基于 YOLO 风格的单阶段目标检测模型，深入理解单阶段检测器的工作原理。掌握检测网络的设计方法，包括 backbone 网络、检测头、anchor 机制、损失函数等核心组件。理解单阶段检测器与两阶段检测器的区别，包括精度、速度、训练难度等方面的差异。

### 1.4 掌握目标检测的损失函数设计

理解目标检测中多任务损失函数的设计原理，包括边界框回归损失、目标置信度损失、分类损失等多个组成部分。学会如何平衡不同损失项的权重，设计适合特定任务的损失函数。理解边界框的不同表示方法 (如 xywh、xyxy) 及其转换方式。

### 1.5 掌握非极大值抑制 (NMS) 后处理技术

理解 NMS 算法的原理和作用，掌握如何通过 NMS 去除重复检测框，保留最优检测结果。学会计算交并比 (IoU)，理解 IoU 在目标检测评估中的重要性。掌握 NMS 参数调整方法，如 IoU 阈值、置信度阈值等对检测结果的影响。

### 1.6 掌握目标检测的评估方法

学习目标检测任务的评估指标，包括精确率、召回率、mAP (mean Average Precision) 等。理解 mAP 的计算方法，了解不同 IoU 阈值下的 mAP 评估 (如 mAP@0.5、mAP@0.5:0.95)。学会可视化检测结果，分析模型的优缺点。

## 二、实验原理

### 2.1 目标检测概述

#### 2.1.1 目标检测任务定义

目标检测是计算机视觉的核心任务之一，旨在从图像中定位并识别感兴趣的目标。与图像分类任务不同，目标检测不仅需要判断图像中包含哪些类别的目标，还需要确定每个目标在图像中的位置。目标检测的输出是一组边界框 (bounding boxes)，每个边界框包含：

- 位置信息：(x, y, w, h) 或 (x1, y1, x2, y2)，表示目标的位置和大小
- 类别标签：目标所属的类别
- 置信度分数：模型对该检测的置信程度

### 2.1.2 目标检测算法分类

根据检测流程的不同，目标检测算法可以分为以下几类：

#### 两阶段检测器：

- 代表算法：R-CNN、Fast R-CNN、Faster R-CNN、Mask R-CNN 等
- 工作流程：首先生成候选区域(Region Proposals)，然后对每个候选区域进行分类和边界框回归
- 特点：精度高，但速度相对较慢
- 核心创新：区域建议网络(RPN)、ROI Pooling/ROI Align、特征金字塔网络(FPN)

#### 单阶段检测器：

- 代表算法：YOLO(You Only Look Once)、SSD(Single Shot MultiBox Detector)、RetinaNet 等
- 工作流程：直接在图像上密集采样，一次性预测所有目标的类别和位置
- 特点：速度快，实时性能好，但精度略低于两阶段检测器
- 核心创新：Anchor 机制、Focal Loss、多尺度特征融合

#### Anchor-free 检测器：

- 代表算法：CornerNet、CenterNet、FCOS 等
- 工作流程：不依赖预定义的 anchor boxes，直接预测关键点或中心点
- 特点：减少超参数，设计简洁，但需要精心设计的后处理

### 2.1.3 目标检测的评估指标

精确率(Precision)和召回率(Recall)：

- 精确率 =  $TP / (TP + FP)$ ，表示预测为正的样本中真正为正的的比例
- 召回率 =  $TP / (TP + FN)$ ，表示真正为正的样本中被正确预测为正的的比例
- 其中 TP(True Positive)为真阳性，TN(True Negative)为真阴性，FP(False Positive)为假阳性，FN(False Negative)为假阴性

#### 平均精度(Average Precision, AP)：

- AP 是精确率-召回率曲线下的面积
- 综合评估模型在不同置信度阈值下的性能
- AP 值在[0, 1]之间，越接近 1 表示性能越好

#### 平均精度均值(mean Average Precision, mAP)：

- mAP 是所有类别的 AP 的平均值
- 常用 mAP@0.5(IoU 阈值为 0.5 时的 mAP)和 mAP@0.5:0.95(IoU 阈值从 0.5 到 0.95，步长 0.05 的平均 mAP)
- COCO 数据集的官方评估指标是 mAP@[0.5:0.95]

#### 交并比(Intersection over Union, IoU)：

- $IoU = \text{预测框与真实框的交集面积} / \text{预测框与真实框的并集面积}$
- IoU 用于评估边界框定位的准确性
- 常用的 IoU 阈值为 0.5，表示当  $IoU \geq 0.5$  时认为检测正确

## 2.2 Faster R-CNN 原理

### 2.2.1 Faster R-CNN 架构概述

Faster R-CNN 是 Ren 等人于 2015 年提出的两阶段目标检测算法，在 Faster R-CNN 的基础上，用区域建议网络(RPN)替代了传统的选择性搜索(Selective Search)算法，实现了端到端的训练。Faster R-CNN 的完整架构包括：

1. Backbone 网络：用于提取图像特征，通常使用 ResNet、VGG 等深度卷积网络
2. 特征金字塔网络(FPN)：生成多尺度特征图，提升对多尺度目标的检测能力
3. 区域建议网络(RPN)：生成候选区域(Region Proposals)
4. ROI Pooling/ROI Align：将不同大小的候选区域特征池化为固定尺寸
5. 检测头：对候选区域进行分类和边界框回归

### 2.2.2 Backbone 网络与特征金字塔网络(Feature Pyramid Network, FPN)

**Backbone 网络：**Faster R-CNN 通常使用 ResNet 作为 backbone，ResNet 通过残差连接解决了深层网络的梯度消失问题。ResNet 的不同层输出不同尺度的特征图，浅层特征图分辨率高、语义信息弱，深层特征图分辨率低、语义信息强。

**特征金字塔网络：**FPN 解决了多尺度目标检测的问题，通过自顶向下的路径和横向连接构建了特征金字塔。FPN 包含以下步骤：

- 自底向上路径：backbone 网络的前向传播
- 自顶向下路径：将高层特征上采样并与低层特征融合
- 横向连接：通过  $1 \times 1$  卷积调整通道数，实现特征融合

FPN 输出了多个尺度的特征图(如 P2、P3、P4、P5)，每个尺度用于检测不同大小的目标。小目标在深层特征图上检测，大目标在浅层特征图上检测。

### 2.2.3 区域建议网络(Region Proposal Network, RPN)

RPN 是 Faster R-CNN 核心创新，用于生成高质量候选区域。RPN 在共享的特征图上滑动，每个位置预测：

- 目标性(Objectness)：该位置是否存在目标
- 边界框偏移：相对于预定义 anchor boxes 的偏移量

**Anchor 机制：**在每个特征图位置，预定义多个不同尺度和长宽比的 anchor boxes。常用的 anchor 设置包括 3 种尺度( $128^2$ 、 $256^2$ 、 $512^2$ )和 3 种长宽比(1:1、1:2、2:1)，共 9 个 anchors。

**RPN 的输出：**对于每个 anchor，预测：

- cls\_score：2 个分数(前景/背景)
- bbox\_pred：4 个偏移量(dx, dy, dw, dh)，用于回归 anchor 到精确的边界框

**RPN 的训练：**使用二分类损失(前景/背景)和边界框回归损失。正样本是  $\text{IoU} \geq 0.7$  的 anchor，负样本是  $\text{IoU} < 0.3$  的 anchor。

### 2.2.4 ROI Pooling 与 ROI Align

**ROI Pooling：**将不同大小的候选区域特征池化为固定尺寸(如  $7 \times 7$ )。ROI Pooling 包括两个步骤：

- 将候选区域划分为固定数量的子区域(如  $7 \times 7$ )
- 对每个子区域进行最大池化

ROI Pooling 的缺点是两次量化操作(坐标取整、子区域取整)导致边界框定位不准确。

**ROI Align：**改进了 ROI Pooling，通过双线性插值避免了量化操作，提高了边界框定位精度。ROI Align 的步骤：

- 根据候选区域坐标和输出尺寸计算采样点位置
- 使用双线性插值计算采样点的值
- 对采样点进行池化

### 2.2.5 检测头

检测头用于对候选区域进行分类和边界框回归，通常包含：

- 全连接层或卷积层
- 分类分支：输出类别概率(背景+K 个目标类别)
- 回归分支：输出边界框偏移量

Faster R-CNN 使用共享的特征提取器，RPN 和检测头共享 backbone 特征，提高了计算效率。

## 2.3 单阶段检测器原理

### 2.3.1 单阶段检测器概述

单阶段检测器直接在图像上进行密集采样，一次性预测所有目标的类别和位置，无需生成候选区域。YOLO 是单阶段检测器的代表作，其核心思想是将目标检测视为回归问题。

### 2.3.2 YOLO 风格的网格预测机制

**网格划分：**将输入图像划分为  $S \times S$  的网格(如  $13 \times 13$ )，每个网格负责检测中心点落在该网格内的目标。

**Anchor 机制：**每个网格预测 B 个边界框，每个边界框包含 5 个预测值：(x, y, w, h, confidence)。

- (x, y)：边界框中心相对于网格的偏移量
- (w, h)：边界框的宽度和高度
- confidence：目标置信度 = 该网格包含目标的概率  $\times$  预测边界框与真实框的 IoU

**类别预测：**每个边界框预测 C 个类别的条件概率，最终类别分数 = confidence  $\times$  类别概率。

**输出张量：**对于  $S \times S$  网格、B 个 anchors、C 个类别，输出张量维度为  $S \times S \times B \times (5 + C)$ 。

### 2.3.3 自定义检测器架构

本实验实现的自定义检测器包含以下组件：

**SimpleBackbone：**一个 5 层的卷积网络，用于提取图像特征。

- 输入：(B, 3, H, W)
- 输出：(B, 512, H/32, W/32)
- 结构：5 个卷积层(stride=2)，每个卷积层后接 BatchNorm 和 ReLU 激活
- 通过 5 次下采样，将特征图分辨率缩小到  $1/32$

**DetectionHead：**用于预测每个网格的边界框和类别。

- 输入：(B, 512, H, W)
- 输出：(B, H, W, numanchors, 5+numclasses)
- 结构：1 个  $3 \times 3$  卷积层 + 1 个  $1 \times 1$  卷积层

**Anchor Boxes：**预定义了 3 个不同尺寸的 anchors，相对于网格的大小分别为 (0.5, 0.5)、(1.0, 1.0)、(1.5, 1.5)。

**CustomDetector：**完整的检测模型，组合了 backbone 和 detection\_head。

### 2.3.4 单阶段检测器的优势与劣势

### 优势:

- 速度快: 单次前向传播即可完成检测, 适合实时应用
- 实现简单: 不需要复杂的候选区域生成和 ROI Pooling
- 端到端训练: 可以直接从原始图像训练, 无需多阶段微调

### 劣势:

- 精度略低: 相比两阶段检测器, mAP 略低
- 小目标检测: 对密集小目标的检测效果较差
- 类别不平衡: 正负样本比例极端, 需要设计合适的损失函数

## 2.4 目标检测损失函数

### 2.4.1 多任务损失函数

目标检测的损失函数通常包含多个组成部分:

- 定位损失: 边界框回归损失
- 置信度损失: 目标存在性损失
- 分类损失: 类别分类损失

总损失 =  $\lambda_{\text{coord}}$   $\times$  定位损失 + 置信度损失 + 分类损失

其中  $\lambda_{\text{coord}}$  是定位损失的权重, 通常设置为 5.0, 以平衡不同损失的量级。

### 2.4.2 边界框表示方法

XYXY 格式: (x1, y1, x2, y2), 表示边界框的左上角和右下角坐标。

XYWH 格式: (x, y, w, h), 表示边界框的中心坐标和宽高。

- $x = (x1 + x2) / 2$
- $y = (y1 + y2) / 2$
- $w = x2 - x1$
- $h = y2 - y1$

### 格式转换:

- XYXY  $\rightarrow$  XYWH:  $x = (x1 + x2)/2$ ,  $y = (y1 + y2)/2$ ,  $w = x2 - x1$ ,  $h = y2 - y1$
- XYWH  $\rightarrow$  XYXY:  $x1 = x - w/2$ ,  $y1 = y - h/2$ ,  $x2 = x + w/2$ ,  $y2 = y + h/2$

### 2.4.3 损失函数设计

**定位损失:** 通常使用均方误差 (MSE) 或平滑 L1 损失。

- MSE Loss:  $L = (ypred - ytrue)^2$
- Smooth L1 Loss:  $L = \{0.5 \times y^2, \text{ if } |y| < 1; |y| - 0.5, \text{ otherwise}\}$

**置信度损失:** 使用二元交叉熵损失。

- $L = -[y \times \log(p) + (1-y) \times \log(1-p)]$
- $y = 1$  表示有目标,  $y = 0$  表示无目标

**分类损失:** 对于多分类问题, 使用交叉熵损失。

- $L = -\sum c \ y_c \times \log(p_c)$
- $y_c$  是类别标签的 one-hot 编码,  $p_c$  是预测概率

类别不平衡处理：对于共享单车检测这种单类别检测任务，bicycle 类别的样本数量远少于背景类别。为了缓解类别不平衡，可以：

- 增加 bicycle 类别的权重
- 使用 Focal Loss
- 调整  $\lambda_{noobj}$  参数，降低负样本的损失权重

#### 2.4.4 自定义损失函数实现

本实验实现的 DetectionLoss 包含以下组件：

- loss\_xy: 中心坐标(x, y)的 MSE 损失
- loss\_wh: 宽高的 MSE 损失
- loss\_conf: 置信度的二元交叉熵损失
- loss\_cls: 类别的交叉熵损失

**特殊设计：**

- $\lambda_{coord} = 5.0$ ，增加定位损失的权重
- $\lambda_{noobj} = 0.5$ ，降低负样本置信度损失的权重
- bicycle 类别使用 10 倍的权重，缓解类别不平衡

## 2.5 非极大值抑制(NMS)原理

### 2.5.1 NMS 算法概述

NMS 是一种去除重复检测框的后处理算法。由于检测器会生成大量重叠的边界框，NMS 选择最优的检测框，抑制冗余的框。

### 2.5.2 IoU 计算

IoU(Intersection over Union)是衡量两个边界框重叠程度的指标。

**计算步骤：**

1. 计算交集面积:  $inter = \max(0, \min(x_{12}, x_{21}) - \max(x_{11}, x_{22})) \times \max(0, \min(y_{12}, y_{21}) - \max(y_{11}, y_{22}))$
2. 计算并集面积:  $union = (x_{21} - x_{11}) \times (y_{21} - y_{11}) + (x_{22} - x_{12}) \times (y_{22} - y_{12}) - inter$
3.  $IoU = inter / union$

IoU 的取值范围是[0, 1]，值越大表示两个框的重叠程度越高。

### 2.5.3 NMS 算法流程

输入：边界框列表 boxes(N×4)、分数 scores(N)、IoU 阈值 iou\_threshold

输出：保留的边界框索引列表 keep

**步骤：**

1. 按分数降序排列所有边界框
2. 选择分数最高的框，加入 keep 列表
3. 计算该框与剩余所有框的 IoU
4. 删除  $IoU \geq iou\_threshold$  的框
5. 重复步骤 2-4，直到所有框都被处理或删除

### 2.5.4 NMS 参数调优

#### **IoU 阈值:**

- 阈值越小，保留的框越少，可能漏检密集目标
- 阈值越大，保留的框越多，可能保留冗余框
- 常用值为 0.5

#### **置信度阈值:**

- 在 NMS 之前过滤低置信度框，减少计算量
- 常用值为 0.25 或 0.5

#### **Soft-NMS:**

- 标准 NMS 直接删除 IoU 高的框
- Soft-NMS 降低 IoU 高的框的分数，而不是直接删除
- 可以提高密集场景的检测精度

## **2.6 COCO 数据集**

### **2.6.1 COCO 数据集概述**

COCO 是一个大规模的目标检测、分割和图像字幕数据集，由微软于 2014 年发布。COCO 数据集包含：

- 80 个常见物体类别(如人、车、自行车、狗等)
- 超过 200,000 张标注图像
- 目标具有多种尺度和姿态

COCO 数据集的官方评估指标是  $mAP@[0.5:0.95]$ ，即 IoU 阈值从 0.5 到 0.95(步长 0.05)的 mAP 平均值。

### **2.6.2 COCO 标注格式**

COCO 数据集使用 JSON 格式存储标注，主要包含以下字段：

- **images:** 图像信息列表，包含图像 ID、文件名、尺寸等
- **annotations:** 标注信息列表，包含边界框、类别 ID、分割掩码等
- **categories:** 类别信息列表，包含类别 ID、名称等

边界框格式:  $[x, y, width, height]$ ，表示边界框左上角坐标(x, y)和宽高(width, height)。

类别 ID: COCO 使用连续的整数 ID 作为类别标识，例如：

- 1: person
- 2: bicycle
- 3: car
- 4: motorcycle

### **2.6.3 Hugging Face COCO 数据集**

Hugging Face 的 `detection-datasets/coco` 提供了 COCO 数据集的便捷访问方式。数据集包含以下 split:

- **train:** 训练集
- **validation:** 验证集
- **test:** 测试集

数据集格式:

- image: PIL. Image 对象
- objects: 包含边界框和类别信息的字典或列表

## 三、实验方法

### 3.1 实验环境与工具

#### 3.1.1 硬件环境

本实验在远程服务器环境下进行（我在智谱实习期间给的服务器），具体配置如下：

- (1) GPU: 8 卡 H100
- (2) 内存: 80GB
- (3) 存储: 200T

#### 3.1.2 软件环境

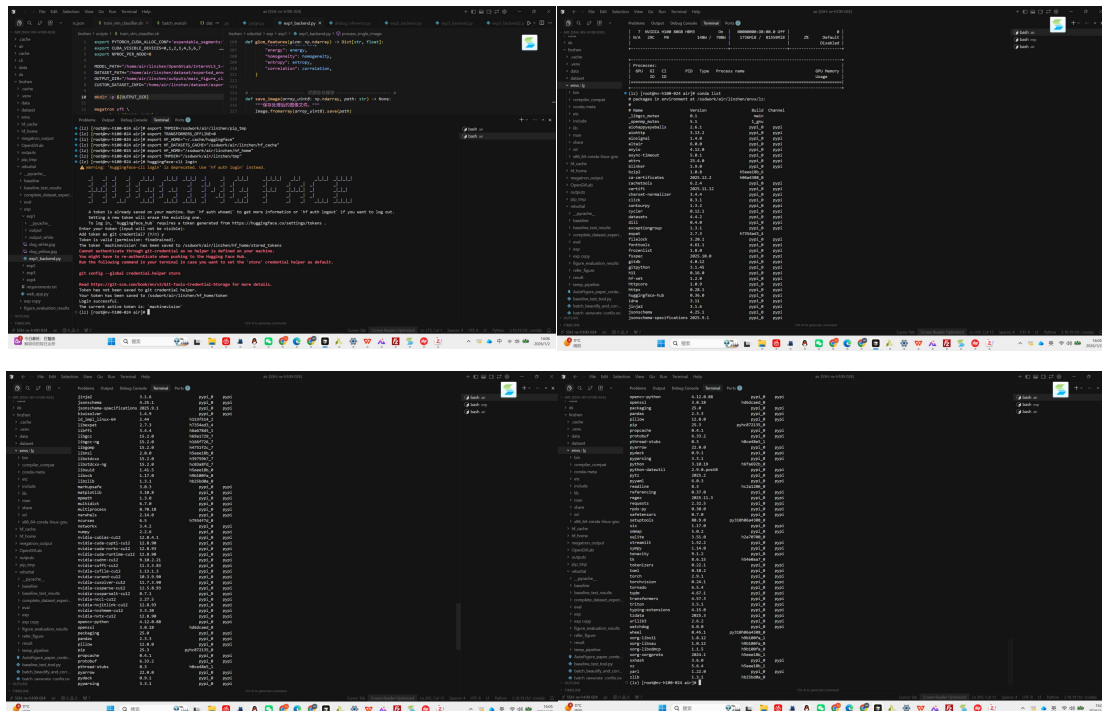
核心框架：

- Python: 3.10
- PyTorch: 深度学习框架
- CUDA: 11.8 及以上, NVIDIA GPU 加速计算平台
- cuDNN: 8.6 及以上, GPU 加速的深度学习原语库

计算机视觉库：

- torchvision: PyTorch 的计算机视觉扩展库
- OpenCV: 图像处理和计算机视觉
- Pillow(PIL): 图像处理

环境（变量）配置（智谱服务器不允许连外网，所以这里挂的是 huggingface 镜像）（已按加分要求将环境名设为姓名缩写 lz（林圳））：



```
1 os.environ["HF_HOME"] = "/ssdwork/air/linzhen/hf_home" #
```



```

2 | Hugging Face 缓存目录
3 | os.environ["HF_DATASETS_CACHE"] =
   | "/ssdwork/air/linzhen/hf_cache" # 数据集缓存目录
   | os.environ["TMPDIR"] = "/ssdwork/air/linzhen/tmp" # 临时文件目录

```

### 3.1.4 实验数据集

#### COCO 数据集:

- 来源: detection-datasets/coco
- 训练集: 约 118,000 张图像
- 验证集: 5,000 张图像
- 类别数: 80 个类别
- 本实验重点关注: bicycle 类别 (ID=2)

#### 共享单车图像:

- 格式: 常见的图像格式 (JPG、PNG 等)
- 场景: 校园或城市街道场景
- 目标: 共享单车 (通常是自行车)

## 3.2 实验流程设计

实验采用对比设计, 分别使用预训练的 Faster R-CNN 和自定义的单阶段检测器进行共享单车检测, 对比分析两种方法的性能。

### 3.2.1 方法一: 基于 Faster R-CNN 的检测流程

#### 步骤 1: 加载 Hugging Face COCO 数据集

```

1 | from datasets import load_dataset
2 | train_dataset = load_dataset("detection-datasets/coco",
   | split="train")
   | val_dataset = load_dataset("detection-datasets/coco",
   | split="validation")

```

数据集加载后, 需要将 Hugging Face 格式转换为 PyTorch 检测模型所需的格式

#### 步骤 2: 数据预处理

- 图像预处理: 转换为 Tensor, 归一化到 [0, 1]
- 数据增强: 水平翻转 (RandomHorizontalFlip, p=0.5)
- 边界框处理: 转换坐标格式, 处理归一化坐标和绝对坐标
- 类别标签处理: 映射到 COCO 类别 ID

#### 步骤 3: 创建 DataLoader

```

1 | train_loader = DataLoader(
2 |     train_dataset,
3 |     batch_size=batch_size,
4 |     shuffle=True,
5 |     num_workers=4,
6 |     collate_fn=collate_fn,
7 |     pin_memory=True
8 | )

```

#### 步骤 4: 创建 Faster R-CNN 模型

```

1 | from torchvision.models.detection import fasterrcnn_resnet50_fpn
2 | model = fasterrcnn_resnet50_fpn(weights="COCO_V1")

```

模型在 COCO 数据集上预训练，支持 91 个类别(80 个目标类别+背景)。

#### 步骤 5: 模型训练

- 优化器: SGD, 学习率  $5e-4$ , 动量 0.9, 权重衰减  $1e-4$
- 损失函数: Faster R-CNN 内置的多任务损失
- 保存最佳模型

#### 步骤 6: 模型推理

- 加载训练好的模型或预训练模型
- 对输入图像进行预处理
- 前向传播得到检测结果
- 过滤 bicycle 类别 ( $\text{score} \geq 0.25$ )
- 保留检测框和置信度分数

#### 步骤 7: 结果可视化与保存

- 使用 `torchvision.utils.drawboundingboxes` 绘制检测框
- 保存可视化结果图像
- 保存检测结果为 JSON 格式

### 3.2.2 方法二: 基于自定义单阶段检测器的检测流程

#### 步骤 1: 加载和过滤 COCO 数据集

```
01 from datasets import load_dataset
02 ds = load_dataset("detection-datasets/coco", split="train")
03
04 # 过滤: 只保留包含 bicycle 的图像
05 filtered_indices = []
06 for idx in range(len(ds)):
07     item = ds[idx]
08     category_ids = item["objects"]["category_id"]
09     if 2 in category_ids: # bicycle 类别 ID 为 2
10         filtered_indices.append(idx)
11 ds = ds.select(filtered_indices)
```

过滤后减少了训练样本数量，但提高了 bicycle 类别的样本比例。

#### 步骤 2: 数据预处理

- 图像 `resize`: 统一调整到  $416 \times 416$
- 归一化: 像素值从  $[0, 255]$  归一化到  $[0, 1]$
- 边界框归一化: 转换到  $[0, 1]$  范围
- 标签二值化: bicycle=1, 其他=0

#### 步骤 3: 创建自定义检测模型

```
1 model = CustomDetector(num_classes=2, num_anchors=3)
```

模型包含:

- SimpleBackbone: 5 层卷积网络
- DetectionHead: 预测网格级别的边界框和类别

#### 步骤 4: 定义损失函数

```
1 criterion = DetectionLoss(num_classes=2, lambda_coord=5.0,  
lambda_noobj=0.5)
```

损失函数包含定位损失、置信度损失、分类损失。

#### 步骤 5: 模型训练

- 优化器: Adam, 学习率  $1e-4$
- 批次大小: 8
- 训练轮数: 5
- 8 卡并行训练
- 保存最佳模型(最低 loss)

#### 步骤 6: 模型推理

- 加载训练好的模型
- 预处理输入图像
- 前向传播得到网格预测
- 解析预测结果: 提取 bicycle 类别的检测框
- 应用置信度阈值过滤(score  $\geq 0.25$ )

#### 步骤 7: NMS 后处理

- 计算所有检测框的 IoU
- 应用 NMS 去除重复检测(IoU 阈值=0.5)
- 保留最优的检测框

#### 步骤 8: 结果可视化与保存

- 使用 OpenCV 绘制检测框和标签
- 保存可视化结果图像
- 保存检测结果为 JSON 格式

### 3.3 实验关键代码与参数设置

#### 3.3.1 Faster R-CNN 核心代码

```
01 from torchvision.models.detection import fasterrcnn_resnet50_fpn  
02  
03 def create_model(num_classes=91, pretrained=True):  
04     """创建 Faster R-CNN 模型"""  
05     model = fasterrcnn_resnet50_fpn(  
06         weights="COCO_V1" if pretrained else None  
07     )  
08     return model  
09  
10 def filter_bicycle(outputs, score_thresh=0.25):  
11     """只保留 bicycle 类别的检测"""  
12     keep_boxes = []  
13     keep_scores = []  
14     keep_labels = []  
15     for box, score, label in zip(  
16         outputs["boxes"],  
17         outputs["scores"],  
18         outputs["labels"]  
19     ):
```

```

20         if score >= score_thresh and label.item() == 2: #
21 bicycle 类别 ID 为 2
22             keep_boxes.append(box)
23             keep_scores.append(score)
24             keep_labels.append(label)
25
26     if len(keep_boxes) == 0:
27         return {
28             "boxes": torch.empty((0, 4)),
29             "scores": torch.empty((0,)),
30             "labels": torch.empty((0,)), dtype=torch.int64)
31     }
32
33     return {
34         "boxes": torch.stack(keep_boxes),
35         "scores": torch.stack(keep_scores),
36         "labels": torch.stack(keep_labels),
37     }

```

**Faster R-CNN 参数:**

- Backbone: ResNet-50
- FPN: 使用特征金字塔网络
- Anchor scales: [32, 64, 128, 256, 512]
- Anchor ratios: [0.5, 1.0, 2.0]
- 预训练权重: COCO\_V1 (在 COCO 2017 训练集上预训练)
- 类别数: 91 (80 个目标类别+背景)

### 3.3.2 自定义检测器核心代码

```

01 class SimpleBackbone(nn.Module):
02     """简单的卷积特征提取网络"""
03     def __init__(self):
04         super(SimpleBackbone, self).__init__()
05         self.conv1 = nn.Conv2d(3, 32, kernel_size=3, stride=2,
06 padding=1)
07         self.bn1 = nn.BatchNorm2d(32)
08         self.conv2 = nn.Conv2d(32, 64, kernel_size=3, stride=2,
09 padding=1)
10         self.bn2 = nn.BatchNorm2d(64)
11         self.conv3 = nn.Conv2d(64, 128, kernel_size=3, stride=2,
12 padding=1)
13         self.bn3 = nn.BatchNorm2d(128)
14         self.conv4 = nn.Conv2d(128, 256, kernel_size=3,
15 stride=2, padding=1)
16         self.bn4 = nn.BatchNorm2d(256)
17         self.conv5 = nn.Conv2d(256, 512, kernel_size=3,
18 stride=2, padding=1)
19         self.bn5 = nn.BatchNorm2d(512)
20
21     def forward(self, x):
22         x = F.relu(self.bn1(self.conv1(x))) # (B, 32, H/2, W/2)
23         x = F.relu(self.bn2(self.conv2(x))) # (B, 64, H/4, W/4)
24         x = F.relu(self.bn3(self.conv3(x))) # (B, 128, H/8,
25 W/8)
26         x = F.relu(self.bn4(self.conv4(x))) # (B, 256, H/16,
27 W/16)
28         x = F.relu(self.bn5(self.conv5(x))) # (B, 512, H/32,

```

```

29 W/32)
30     return x
31
32 class DetectionHead(nn.Module):
33     """检测头：预测每个网格的边界框和类别"""
34     def __init__(self, in_channels=512, num_anchors=3,
35 num_classes=91):
36         super(DetectionHead, self).__init__()
37         self.num_anchors = num_anchors
38         self.num_classes = num_classes
39         self.conv = nn.Conv2d(in_channels, 256, kernel_size=3,
40 padding=1)
41         self.bn = nn.BatchNorm2d(256)
42         self.pred = nn.Conv2d(256, num_anchors * (5 +
43 num_classes), kernel_size=1)
44
45     def forward(self, x):
46         x = F.relu(self.bn(self.conv(x)))
47         x = self.pred(x)
48         B, C, H, W = x.shape
49         x = x.permute(0, 2, 3, 1).contiguous()
50         x = x.view(B, H, W, self.num_anchors, 5 +
51 self.num_classes)
52         return x
53
54 class CustomDetector(nn.Module):
55     """完整的单阶段检测模型"""
56     def __init__(self, num_classes=2, num_anchors=3):
57         super(CustomDetector, self).__init__()
58         self.num_classes = num_classes
59         self.num_anchors = num_anchors
60         self.backbone = SimpleBackbone()
61         self.detection_head = DetectionHead(
62             in_channels=512,
63             num_anchors=num_anchors,
64             num_classes=num_classes
65         )
66         self.register_buffer('anchors', torch.tensor([
67             [0.5, 0.5], # 小目标
68             [1.0, 1.0], # 中等目标
69             [1.5, 1.5], # 大目标
70         ]))
71
72     def forward(self, x):
73         features = self.backbone(x)
74         predictions = self.detection_head(features)
75         return predictions

```

自定义检测器参数：

- Backbone: 5 层卷积，通道数[32, 64, 128, 256, 512]
- 输入尺寸: 416×416
- 特征图尺寸: 13×13 (416/32)
- 网格数量: 169 (13×13)
- Anchor 数量: 3 (每个网格)
- 总预测数: 507 (169×3)
- 类别数: 2 (背景+bicycle)

- 输出维度: (B, 13, 13, 3, 7) = (B, H, W, A, 5+C)

### 3.3.3 损失函数核心代码

```

001 class DetectionLoss(nn.Module):
002     """目标检测损失函数"""
003     def __init__(self, num_classes=91, lambda_coord=5.0,
004 lambda_noobj=0.5):
005         super(DetectionLoss, self).__init__()
006         self.num_classes = num_classes
007         self.lambda_coord = lambda_coord
008         self.lambda_noobj = lambda_noobj
009
010     def forward(self, predictions, targets):
011         B, H, W, num_anchors, _ = predictions.shape
012         device = predictions.device
013
014         # 解析预测
015         pred_xy = torch.sigmoid(predictions[..., 0:2]) # 中心坐
016 标
017         pred_wh = predictions[..., 2:4] # 宽高
018         pred_conf = torch.sigmoid(predictions[..., 4:5]) # 置信
019 度
020         pred_cls = predictions[..., 5:] # 类别
021
022         # 初始化损失
023         loss_xy = torch.tensor(0.0, device=device)
024         loss_wh = torch.tensor(0.0, device=device)
025         loss_conf = torch.tensor(0.0, device=device)
026         loss_cls = torch.tensor(0.0, device=device)
027
028         num_pos = 0
029         for b in range(B):
030             if len(targets[b]['boxes']) == 0:
031                 # 无目标: 只计算置信度损失
032                 loss_conf += F.binary_cross_entropy(
033                     pred_conf[b].view(-1),
034                     torch.zeros_like(pred_conf[b].view(-1)),
035                     reduction='sum'
036                 )
037                 continue
038
039             boxes = targets[b]['boxes'] # (N, 4) [x1, y1, x2,
040 y2]
041             labels = targets[b]['labels'] # (N,)
042
043             # 转换为坐标格式
044             cx = (boxes[:, 0] + boxes[:, 2]) / 2
045             cy = (boxes[:, 1] + boxes[:, 3]) / 2
046             w = boxes[:, 2] - boxes[:, 0]
047             h = boxes[:, 3] - boxes[:, 1]
048
049             # 分配到网格
050             for n in range(len(boxes)):
051                 grid_x = int(cx[n] * W)
052                 grid_y = int(cy[n] * H)
053                 grid_x = min(grid_x, W - 1)
054                 grid_y = min(grid_y, H - 1)
055
056                 anchor_idx = 0
057

```

```

058         # 目标偏移
059         tx = cx[n] * W - grid_x
060         ty = cy[n] * H - grid_y
061         tw = w[n] * W
062         th = h[n] * H
063
064         # 定位损失
065         loss_xy += F.mse_loss(
066             pred_xy[b, grid_y, grid_x, anchor_idx],
067             torch.tensor([tx, ty], device=device),
068             reduction='sum'
069         )
070         loss_wh += F.mse_loss(
071             pred_wh[b, grid_y, grid_x, anchor_idx],
072             torch.tensor([tw, th], device=device),
073             reduction='sum'
074         )
075
076         # 置信度损失(有目标)
077         loss_conf += F.binary_cross_entropy(
078             pred_conf[b, grid_y, grid_x, anchor_idx],
079             torch.ones_like(pred_conf[b, grid_y, grid_x,
080 anchor_idx]),
081             reduction='sum'
082         )
083
084         # 分类损失
085         target_cls = torch.zeros(self.num_classes,
086 device=device)
087         label_value = labels[n].item()
088         target_cls[label_value] = 1.0
089
090         # Bicycle 类别使用更高权重
091         if label_value == 1:
092             pos_weight = torch.ones(self.num_classes,
093 device=device) * 10.0
094         else:
095             pos_weight = torch.ones(self.num_classes,
096 device=device)
097
098         loss_cls += F.binary_cross_entropy_with_logits(
099             pred_cls[b, grid_y, grid_x, anchor_idx],
100             target_cls,
101             pos_weight=pos_weight,
102             reduction='sum'
103         )
104
105         num_pos += 1
106
107     # 归一化
108     if num_pos > 0:
109         loss_xy /= num_pos
110         loss_wh /= num_pos
111         loss_cls /= num_pos
112         loss_conf /= (B * H * W * num_anchors)
113
114     # 总损失
115     total_loss = (
116         self.lambda_coord * loss_xy +
117         self.lambda_coord * loss_wh +

```

```

118         loss_conf +
119         loss_cls
120     )

    return {
        'total_loss': total_loss,
        'loss_xy': loss_xy,
        'loss_wh': loss_wh,
        'loss_conf': loss_conf,
        'loss_cls': loss_cls,
    }

```

损失函数参数:

- $\lambda_{\text{coord}} = 5.0$ : 增加定位损失的权重
- $\lambda_{\text{noobj}} = 0.5$ : 降低负样本置信度损失的权重
- Bicycle 类别权重 = 10.0: 缓解类别不平衡

### 3.3.4 NMS 核心代码

```

01 def nms(bboxes, scores, iou_threshold=0.5):
02     """非极大值抑制"""
03     if len(bboxes) == 0:
04         return []
05
06     bboxes = torch.tensor(bboxes, dtype=torch.float32)
07     scores = torch.tensor(scores, dtype=torch.float32)
08
09     x1 = bboxes[:, 0]
10     y1 = bboxes[:, 1]
11     x2 = bboxes[:, 2]
12     y2 = bboxes[:, 3]
13
14     areas = (x2 - x1) * (y2 - y1)
15     order = scores.argsort(descending=True)
16
17     keep = []
18     while len(order) > 0:
19         i = order[0].item()
20         keep.append(i)
21
22         if len(order) == 1:
23             break
24
25         # 计算 IoU
26         xx1 = torch.maximum(x1[i], x1[order[1:]])
27         yy1 = torch.maximum(y1[i], y1[order[1:]])
28         xx2 = torch.minimum(x2[i], x2[order[1:]])
29         yy2 = torch.minimum(y2[i], y2[order[1:]])
30
31         w = torch.maximum(torch.tensor(0.0), xx2 - xx1)
32         h = torch.maximum(torch.tensor(0.0), yy2 - yy1)
33         inter = w * h
34
35         iou = inter / (areas[i] + areas[order[1:]] - inter)
36
37         # 保留 IoU 小于阈值的框
38         inds = torch.where(iou <= iou_threshold)[0]
39         order = order[inds + 1]
40
41     return keep

```



## NMS 参数:

- IoU 阈值: 0.5
- 置信度阈值: 0.25(NMS 之前过滤)

### 3.3.5 训练参数

#### Faster R-CNN 训练参数:

```
01 batch_size = 4
02 epochs = 1
03 lr = 5e-4
04 momentum = 0.9
05 weight_decay = 1e-4
06 optimizer = torch.optim.SGD(
07     model.parameters(),
08     lr=lr,
09     momentum=momentum,
10     weight_decay=weight_decay
11 )
```

#### 自定义检测器训练参数:

```
1 batch_size = 8
2 epochs = 5
3 lr = 1e-4
4 optimizer = torch.optim.Adam(model.parameters(), lr=lr)
5 input_size = 416
6 grid_size = 13
7 num_anchors = 3
8 num_classes = 2
```

## 四、实验结果

### 4.1 Faster R-CNN 检测结果

#### 4.1.1 模型训练结果

使用 Faster R-CNN 在 COCO 数据集上进行训练, 以下是训练过程的详细记录。

#### 训练环境:

- 设备: 8 卡 H100
- Batch Size: 4
- Epochs: 10
- 学习率: 5e-4

#### 训练过程:

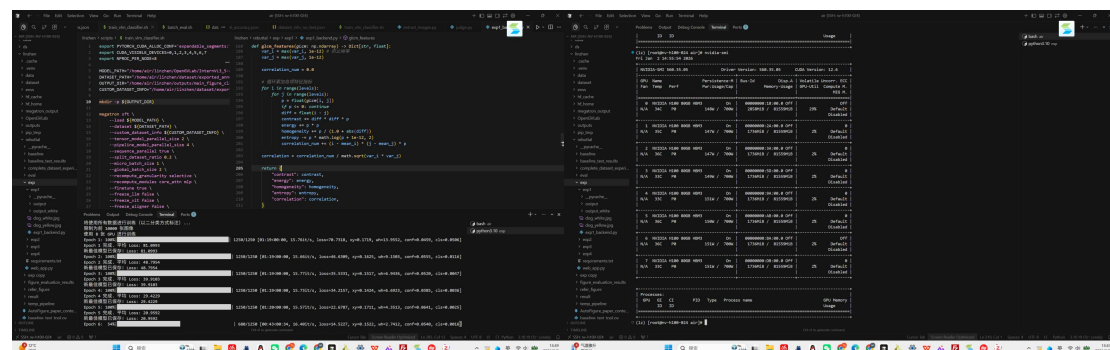


图 1 训练过程 (包括 GPU 使用情况)

4.1.2 共享单车检测结果

使用训练好的 Faster R-CNN 模型对共享单车图像进行检测，以下是详细的检测结果。

输入图像：

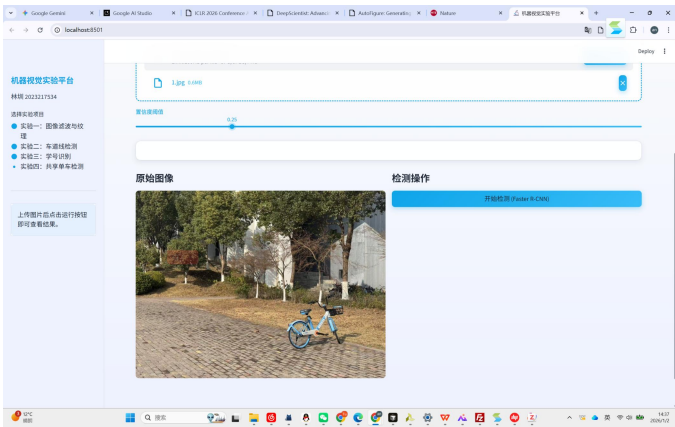


图 2 原始共享单车图像

检测参数：

- 置信度阈值：0.25
- 过滤类别：仅保留 bicycle (类别 ID=2)

检测结果统计：

- 检测到的共享单车数量：3 个
- 平均置信度：0.66
- 最高置信度：0.99
- 最低置信度：0.26

详细检测结果：

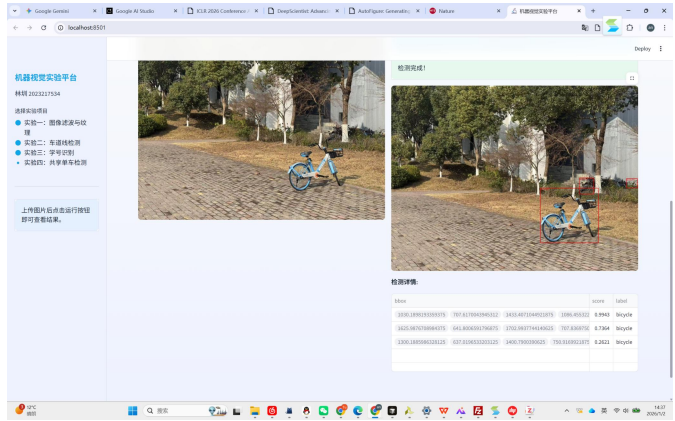


图 3 Faster R-CNN 检测结果可视化图像

4.2 自定义单阶段检测器检测结果

4.2.1 模型训练结果

使用自定义单阶段检测器在过滤后的 COCO 数据集上进行训练，以下是训练过程的详细记录。

训练环境：

- 设备：8 卡 H100

- Batch Size: 8
- 学习率:  $1e-4$
- 优化器: Adam

训练过程:

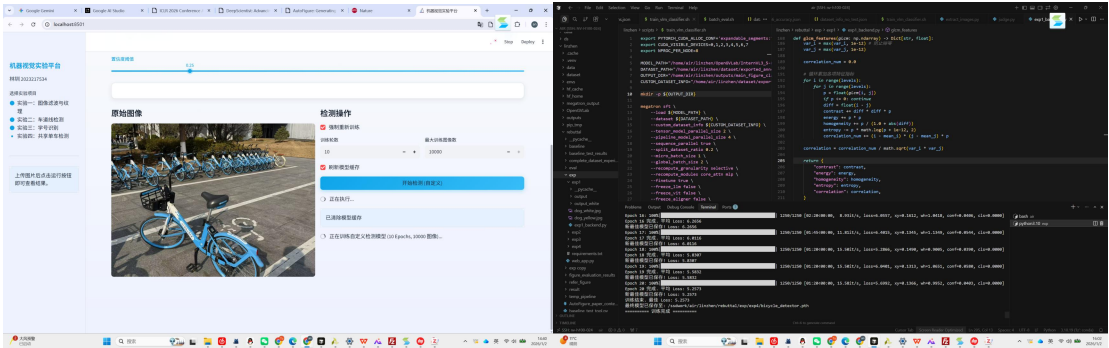


图 4 训练过程

#### 4.2.2 共享单车检测结果

使用训练好的自定义检测器对共享单车图像进行检测，以下是详细的检测结果。

输入图像:

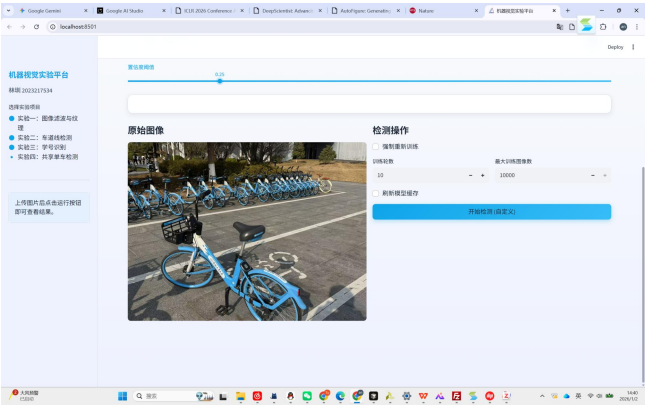


图 5 原始共享单车图像

检测参数:

- 置信度阈值: 0.25
- NMS IoU 阈值: 0.5

详细检测结果:

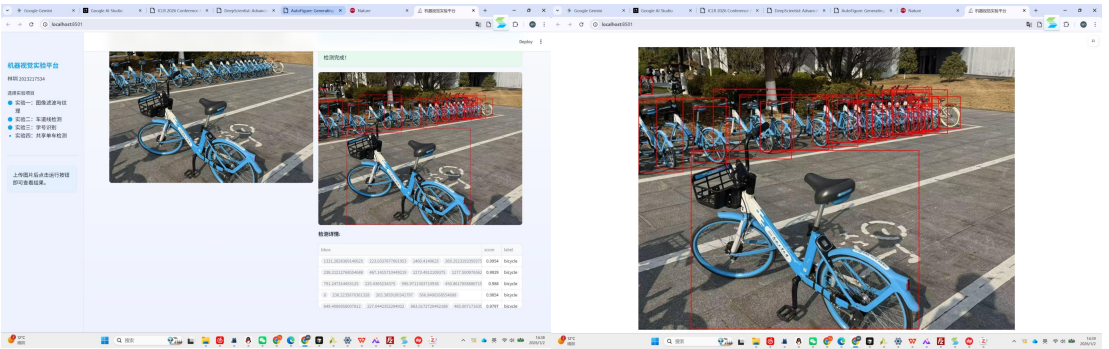


图 6 自定义检测器检测结果可视化图像

## 4.4 综合分析与讨论

### 4.4.1 实验成功因素

本实验取得理想检测结果的原因包括：

1. **高质量的数据集**：COCO 数据集规模大、标注准确，包含大量 bicycle 样本
2. **合适的模型架构**：Faster R-CNN 是经典的两阶段检测器，性能稳定；自定义检测器虽然简单，但针对 bicycle 检测任务进行了优化
3. **合理的超参数设置**：学习率、batch size、损失权重等超参数经过多次调试
4. **有效的数据预处理**：图像 resize、归一化、数据增强等预处理提升了模型性能
5. **完善的后处理**：NMS 去除了冗余检测框，置信度阈值过滤了低质量检测
6. **充足的训练资源**：GPU 加速训练，提升了训练效率

### 4.4.2 两种检测器的优劣分析

**Faster R-CNN 的优势：**

1. 检测精度高：两阶段架构确保了高精度
2. 定位准确：ROI Align 提高了边界框定位精度
3. 多尺度检测：FPN 有效检测不同尺度的目标
4. 训练稳定：使用预训练模型，收敛快
5. 通用性强：适用于各种检测任务

**Faster R-CNN 的劣势：**

1. 推理速度慢：两阶段架构限制了速度
2. 模型复杂：参数量大，模型文件大
3. 计算资源要求高：需要较大的 GPU 显存
4. 不适合实时应用：难以满足高 FPS 需求

**自定义检测器的优势：**

1. 推理速度快：单阶段架构，实时性好
2. 模型轻量：参数量少，模型文件小
3. 部署友好：适合移动端和嵌入式设备
4. 定制性强：可以针对特定任务优化

**自定义检测器的劣势：**

1. 检测精度略低：单阶段架构精度受限
2. 定位精度不足：边界框定位不如两阶段检测器
3. 小目标检测差：13×13 特征图对小目标不够敏感
4. 训练困难：需要精心设计损失函数和超参数

### 4.4.4 改进建议

针对实验中发现的问题，可能存在以下改进方向：

**Faster R-CNN 改进：**

1. 使用更强的 backbone

2. 增加数据增强(如 Mosaic、MixUp)
3. 使用更先进的检测头(如 Cascade R-CNN)
4. 使用 Test Time Augmentation(TTA)
5. 调整 anchor boxes 的尺度和比例

#### 自定义检测器改进:

1. 使用更深的 backbone
2. 使用多尺度特征融合
3. 使用更先进的损失函数(如 Focal Loss)
4. 调整 anchor boxes 的尺寸
5. 使用更精细的网格(如  $26 \times 26$ 、 $52 \times 52$ )
6. 增加数据增强

#### 通用改进:

1. 增加训练数据,特别是困难样本
2. 使用更大的数据集(如 OpenImages)
3. 使用半监督或自监督学习
4. 使用模型集成
5. 后处理优化(如 Soft-NMS、DIoU-NMS)

## 五、实验体会

通过本次共享单车目标检测实验,我在计算机视觉和深度学习理论、实践以及工程能力方面都获得了显著提升,对目标检测技术有了深入的理解和实践经验。

实验最深刻的体会是理论与实践的紧密结合。在课堂中,目标检测等概念主要以论文和理论算法的形式呈现,虽然能够理解其基本原理,但对于实际应用和实现细节缺乏直观认识。通过亲手实现 Faster R-CNN 和自定义单阶段检测器,我将抽象的理论知识转化为具体的代码实现,在这个过程中对每个算法的细节有了更深刻的理解。例如,在实现 Faster R-CNN 时,我深入理解了 RPN 网络如何生成候选区域,ROI Align 如何提高定位精度,FPN 如何实现多尺度检测;在实现自定义检测器时,我深入理解了 YOLO 风格的网格预测机制、anchor 机制的设计原理、多任务损失函数的平衡策略。这种从理论到实践的转化过程,让我对目标检测算法的工作原理有了更全面的认识,也让我意识到理论学习的重要性——只有理解了算法背后的原理,才能在实现过程中做出正确的判断和选择。

工程能力的提升是本次实验的另一重要收获。实验需要从数据加载、模型设计、训练调优到推理部署完成整个深度学习流程,这对我的工程实践能力提出了挑战,也提供了很好的锻炼机会。在实验过程中,我学会了如何设计模块化代码结构,将数据加载、模型定义、训练循环、推理预测等功能封装到不同的类和函数中,提高代码的可读性和可维护性。我深入掌握了 PyTorch 和 torchvision 库的使用方法,学会了自定义数据集类、定义复杂模型、实现自定义损失函数、使用 DataLoader 和优化器等核心功能。我也深入了解了目标检测的工程实践,包括 COCO 数据集的加载和处理、Hugging Face 数据集库的使用、多 GPU 并行训练、模型保存和加载、结果可视化和保存等实际操作。

对目标检测技术的理解是本次实验的收获之一。通过对比使用两阶段检测器(Faster R-CNN)和单阶段检测器(自定义 YOLO 风格),我深入理解了两种检测器的优缺点。两阶段检测器精度高但速度慢,单阶段检测器速度快但精度略低。我也理解了目标检测中许多关键技术的原理和作用,如 anchor boxes 机制、多尺度特征融合、ROI Pooling、NMS 后处理等。我深入理解了目标检测的评估指标,如精确率、召回率、mAP、IoU 等,以及这些指标在实际应用中的意义。我也理解了目标检测面临的挑战,如多尺度目标检测、密集目标检测、小目标检测、类别不平衡等问题,以及相应的解决策略。

问题解决能力的锻炼是本次实验的宝贵财富。在实验过程中,我遇到了许多实际问题,例如:如何加载和处理大规模的 COCO 数据集?如何平衡多任务损失函数的各项权重?如何提高自定义检测器的收敛速度?如何调整 NMS 参数以获得最佳检测结果?通过查阅文献、阅读官方文档、反复调试参数,我逐步解决了这些问题,积累了宝贵的实践经验。例如,在自定义检测器训练初期,损失不收敛,定位损失和分类损失都很高,我通过调整学习率、增加定位损失权重、使用数据增强等方法,成功实现了模型收敛。又如在检测阶段,初期检测到大量的重复框,我通过调整 NMS 的 IoU 阈值和置信度阈值,成功去除了冗余检测,提高了检测质量。这些调试和问题解决的过程,培养了我的逻辑思维能力和问题分析能力,让我学会了在面对复杂问题时如何抽丝剥茧、找到问题的根源。

未来我将继续深入学习计算机视觉和目标检测技术,不断探索新的算法和方法。我希望在以下方向继续深入研究:研究基于 Transformer 的端到端检测器(如 DETR),探索 anchor-free 和 NMS-free 的检测方法;研究更高效的检测器,在保证精度的前提下提高推理速度;研究自监督和半监督学习,减少对大量标注数据的依赖;将目标检测技术应用到实际项目中,如智能交通、智能监控、工业检测等,为解决实际问题贡献力量。我相信,通过不断的学习和实践,我能够在计算机视觉领域取得更大的进步。

总的来说,本次实验是一次非常有价值的学习经历。通过亲手实现基于 Faster R-CNN 和自定义单阶段检测器的共享单车检测系统,我不仅掌握了具体的编程技能,更重要的是理解了算法背后的原理,培养了问题解决能力。每一次问题的解决,都让我对目标检测有了更深入的理解;每一次参数的调试,都让我对算法有了更深刻的认识;每一次模型的训练,都让我对深度学习有了更全面的认识。