

# 轻量化指标及其计算代码示例

## 轻量化需要考虑的相关指标

### 一、效率指标 (Efficiency Metrics)

这类指标直接衡量模型对计算资源的消耗，是轻量化的核心评估维度。

#### 1. 参数量 (Number of Parameters)

- 定义：**模型中需要学习的权重和偏置的总数量。
- 意义：**参数量直接影响**模型大小 (Model Size)**、**内存占用**和**通信成本**。参数量越少，模型越容易部署在存储受限的设备（如手机、嵌入式设备）上，模型更新所需的带宽也越小。
- 单位：**通常用 Millions (M) 或 Billions (B) 表示，如 5.6M 参数。
- 实例：**MobileNetV2 有约 3.4M 参数，而标准的 ResNet-50 有约 25.6M 参数，前者明显更轻量。

#### 2. 计算量 (Computational Complexity)

- 定义：**模型进行一次前向传播所需的浮点运算次数。
- 意义：**直接反映了模型的**计算速度**和**能耗**。计算量越低，模型推理越快，耗电越少。
- 主要指标：**
  - FLOPs (Floating Point Operations)：**浮点运算次数。常用单位是 GFLOPs ( $10^9$  次)。
  - MACs (Multiply-Accumulate Operations)：**乘加运算次数。1个MAC包含一个乘法和一個加法操作，通常认为  $1 \text{ MAC} \approx 2 \text{ FLOPs}$ 。MACs是更常用的硬件相关指标。
- 实例：**处理一张 224x224 的图片，EfficientNet-B0 需要约 0.39 GFLOPs，而 ResNet-50 需要 ~4.1 GFLOPs，前者的计算量仅为后者的不到十分之一。

#### 3. 内存访问量 (Memory Access Cost, MAC)

- 定义：**模型推理过程中，读取输入、权重和写入中间特征图所消耗的内存访问总量。
- 意义：**在现代计算架构（如GPU、NPU）中，**内存访问的能耗和延迟常常远高于计算本身**。高内存访问量会成为性能瓶颈。优化内存访问模式与减少计算量同等重要。
- 实例：**深度可分离卷积 (Depthwise Separable Convolution) 相比标准卷积，不仅大幅降低了 FLOPs，更重要的是显著减少了中间特征图的内存占用和访问次数。

#### 4. 实际推理速度 (Inference Latency/Throughput)

- 定义：**
  - 延迟 (Latency)：**处理单一样本所需的时间（如毫秒, ms）。
  - 吞吐量 (Throughput)：**单位时间内能处理的样本数量（如 帧/秒, FPS）。
- 意义：**这是最直接、最真实的效率指标。**FLOPs低不等于速度快**，因为速度还受硬件特性（并行度、内存带宽、缓存）、软件实现（算子优化、推理引擎）等因素严重影响。

- 测量注意点：**必须在**目标硬件平台**（如 iPhone 15、华为昇腾 310、NVIDIA Jetson Nano）和**实际部署环境**（如使用 TensorRT, CoreML, ONNX Runtime 等优化后的引擎）上进行测试。
- 实例：**Apple 的 MobileCLIP 论文中，所有模型都在 iPhone 12 Pro Max 上实测了延迟（ms），证明其模型在真实移动设备上的速度优势。

## 5. 能耗 (Energy Consumption)

- 定义：**模型完成一次推理所消耗的能量。
- 意义：**对于电池供电的移动设备和物联网设备，能耗是决定性因素。
- 测量：**通常需要专门的硬件功率计（如 Monsoon Solutions 的功率监测仪）或在芯片层面读取功耗计数器。
- 实例：**一项研究可能发现，模型A的FLOPs比模型B低20%，但在特定手机芯片上，模型A的能耗反而更高，因为它触发了更多的内存访问。

## 二、性能指标 (Performance Metrics)

轻量化不能以牺牲过多性能为代价，因此需要在效率和质量之间进行权衡。

### 1. 准确率 (Accuracy)

- 定义：**模型在特定任务上的预测正确率。
- 意义：**这是最核心的性能指标。轻量化模型的目标是在尽可能保持准确率的前提下提升效率。
- 常用指标：**
  - 分类任务：**Top-1 Accuracy, Top-5 Accuracy。
  - 检测任务：**mAP (mean Average Precision)。
  - 分割任务：**mIoU (mean Intersection over Union)。
- 实例：**比较 MobileNetV3（轻量）和 ResNet-50（重量）在 ImageNet 上的 Top-1 准确率，并对比它们的参数量和FLOPs。

### 2. 精度-效率权衡 (Accuracy-Efficiency Trade-off)

- 定义：**通过图表等形式综合展示不同模型或同一模型不同变体的准确率和效率（如延迟、FLOPs）之间的关系。
- 意义：**这是评估轻量化技术的**黄金标准**。一个好的轻量化模型应该在图中处于**帕累托前沿 (Pareto Frontier)**，即在同一效率下准确率最高，或同一准确率下效率最优。
- 实例：**

(这是一个典型的可视化图表示例)

模型	Top-1 Acc (%)	延迟 (ms)	参数量 (M)
Model A (重)	78.5	50	50
Model B (轻)	76.2	15	5
Model C (更轻)	72.1	5	2
Our Model	77.0	12	4

我们的模型在准确率接近最重模型（A）的同时，延迟和参数量接近最轻的模型（C），找到了一个更好的权衡点。

## 三、应用实例：结合您的遥感课题

假设您设计了一个用于遥感图像场景分类的轻量化模型 RSNet-Lite。

### 1. 效率评估：

- **参数量**：您计算出 RSNet-Lite 有 **1.5M** 参数。这意味着模型文件大小约为  $1.5 * 10^6 * 4 \text{ bytes} \approx 6 \text{ MB}$ （假设FP32精度），非常适合在无人机机载计算机上存储和更新。
- **计算量**：您测得 RSNet-Lite 处理一张 512x512 遥感图像需要 **0.8 GFLOPs**。这远低于 baseline 模型 ResNet-34 的 ~3.6 GFLOPs。
- **推理速度**：您在目标硬件 **Jetson Xavier NX** 上使用 TensorRT 部署，实测：
  - **延迟**：15 ms ( $\approx 66 \text{ FPS}$ )
  - **对比**：Baseline 模型在同一平台上的延迟为 80 ms (12.5 FPS)。您的模型实现了 **5倍以上** 的加速，满足实时处理需求。

### 2. 性能评估：

- **准确率**：在 RSSCN7（一个经典遥感场景数据集）上，您的 RSNet-Lite 达到了 **92.5%** 的准确率。
- **精度-效率权衡**：您绘制了下图，清晰地展示了您的模型在准确率和延迟之间取得了更好的平衡，处于帕累托前沿。

### 3. 综合结论：

“我们提出的 RSNet-Lite 模型，仅用 1.5M 参数和 0.8 GFLOPs，在 Jetson Xavier NX 上实现了 15ms 的延迟和 66 FPS 的吞吐量，同时保持了 92.5% 的分类准确率。与基准模型相比，在精度损失极小 (<1%) 的情况下，速度提升了 5 倍以上，显著更适用于资源受限的遥感实时分析场景。”

## 轻量化指标计算代码示例

下面我将提供计算各种轻量化指标的Python代码示例，主要使用PyTorch框架和相关库。

### 1. 安装必要的库

```
pip install torch torchvision thop ptflops pandas numpy
```

### 2. 参数量计算

```
import torch
import torch.nn as nn
from torchvision.models import mobilenet_v2

def calculate_parameters(model):
    """计算模型参数量"""
    total_params = sum(p.numel() for p in model.parameters())
```

```

    trainable_params = sum(p.numel() for p in model.parameters() if
p.requires_grad)

    return total_params, trainable_params

# 示例使用
model = mobilenet_v2(pretrained=False)
total_params, trainable_params = calculate_parameters(model)
print(f"总参数量: {total_params:,} ({(total_params/1e6):.2f}M)")
print(f"可训练参数量: {trainable_params:,} ({(trainable_params/1e6):.2f}M)")

```

### 3. 计算量(FLOPs)计算

```

from thop import profile
from torchvision.models import mobilenet_v2, resnet50
import torch

def calculate_flops(model, input_size=(1, 3, 224, 224)):
    """计算模型FLOPs"""
    device = next(model.parameters()).device
    input_tensor = torch.randn(input_size).to(device)

    flops, params = profile(model, inputs=(input_tensor,))
    return flops, params

# 示例使用
model = mobilenet_v2(pretrained=False)
flops, params = calculate_flops(model)
print(f"FLOPs: {flops:,} ({(flops/1e9):.2f}G)")
print(f"参数量: {params:,} ({(params/1e6):.2f}M)")

# 比较不同模型
models = {
    "MobileNetV2": mobilenet_v2(pretrained=False),
    "ResNet50": resnet50(pretrained=False)
}

results = {}
for name, model in models.items():
    flops, params = calculate_flops(model)
    results[name] = {"FLOPs(G)": flops/1e9, "Params(M)": params/1e6}

# 创建比较表格
import pandas as pd
df = pd.DataFrame(results).T
print("\n模型比较:")
print(df)

```

### 4. 内存占用计算

```

def calculate_memory_usage(model, input_size=(1, 3, 224, 224)):
    """估算模型内存占用"""
    # 参数内存
    param_memory = sum(p.numel() * p.element_size() for p in model.parameters())

    # 缓冲区内内存
    buffer_memory = sum(b.numel() * b.element_size() for b in model.buffers())

    # 前向传播中间激活值内存（近似估算）
    # 注意：这是一个粗略估算，实际值可能因实现而异
    model.eval()
    with torch.no_grad():
        input_tensor = torch.randn(input_size)
        # 使用钩子捕获中间激活值大小
        activation_memory = 0
        hooks = []

        def hook_fn(module, input, output):
            nonlocal activation_memory
            if isinstance(output, torch.Tensor):
                activation_memory += output.numel() * output.element_size()

        for layer in model.modules():
            if isinstance(layer, (nn.Conv2d, nn.Linear, nn.BatchNorm2d)):
                hooks.append(layer.register_forward_hook(hook_fn))

        model(input_tensor)

        # 移除钩子
        for hook in hooks:
            hook.remove()

    total_memory = param_memory + buffer_memory + activation_memory

    return {
        "参数内存(MB)": param_memory / (1024 ** 2),
        "缓冲区内内存(MB)": buffer_memory / (1024 ** 2),
        "激活值内存(MB)": activation_memory / (1024 ** 2),
        "总内存(MB)": total_memory / (1024 ** 2)
    }

# 示例使用
model = mobilenet_v2(pretrained=False)
memory_info = calculate_memory_usage(model)
for k, v in memory_info.items():
    print(f"{k}: {v:.2f}")

```

## 5. 推理速度测试

```

import time
import numpy as np

```

```

def measure_inference_speed(model, input_size=(1, 3, 224, 224), num_runs=100,
warmup=10):
    """测量模型推理速度"""
    device = next(model.parameters()).device
    model.eval()

    input_tensor = torch.randn(input_size).to(device)

    # 预热
    with torch.no_grad():
        for _ in range(warmup):
            _ = model(input_tensor)

    # 测量推理时间
    times = []
    with torch.no_grad():
        for _ in range(num_runs):
            start_time = time.time()
            _ = model(input_tensor)
            end_time = time.time()

            # 确保同步(CUDA)
            if device.type == 'cuda':
                torch.cuda.synchronize()

            times.append(end_time - start_time)

    times = np.array(times)
    mean_time = np.mean(times) * 1000 # 转换为毫秒
    std_time = np.std(times) * 1000
    fps = 1000 / mean_time

    return {
        "平均推理时间(ms)": mean_time,
        "时间标准差(ms)": std_time,
        "FPS": fps
    }

# 示例使用
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model = mobilenet_v2(pretrained=False).to(device)
speed_info = measure_inference_speed(model)

for k, v in speed_info.items():
    print(f"{k}: {v:.2f}")

```

## 6. 完整评估脚本

```

import torch
import torch.nn as nn
from torchvision.models import mobilenet_v2, resnet18

```

```

import pandas as pd
from thop import profile

class ModelEvaluator:
    def __init__(self, device=None):
        self.device = device or torch.device("cuda" if torch.cuda.is_available()
else "cpu")

    def evaluate_model(self, model, model_name, input_size=(1, 3, 224, 224)):
        """全面评估模型"""
        model = model.to(self.device)
        model.eval()

        # 计算参数量
        total_params = sum(p.numel() for p in model.parameters())

        # 计算FLOPs
        input_tensor = torch.randn(input_size).to(self.device)
        flops, _ = profile(model, inputs=(input_tensor,), verbose=False)

        # 测量推理速度
        speed_info = self.measure_inference_speed(model, input_size)

        return {
            "模型": model_name,
            "参数量(M)": total_params / 1e6,
            "FLOPs(G)": flops / 1e9,
            "推理时间(ms)": speed_info["平均推理时间(ms)"],
            "FPS": speed_info["FPS"]
        }

    def measure_inference_speed(self, model, input_size, num_runs=100,
warmup=10):
        """测量推理速度"""
        input_tensor = torch.randn(input_size).to(self.device)

        # 预热
        with torch.no_grad():
            for _ in range(warmup):
                _ = model(input_tensor)

        # 测量
        times = []
        with torch.no_grad():
            for _ in range(num_runs):
                start_time = time.time()
                _ = model(input_tensor)

                if self.device.type == 'cuda':
                    torch.cuda.synchronize()

                end_time = time.time()
                times.append(end_time - start_time)

```

```

        times = np.array(times)
        mean_time = np.mean(times) * 1000
        fps = 1000 / mean_time

    return {"平均推理时间(ms)": mean_time, "FPS": fps}

def compare_models(self, models_dict, input_size=(1, 3, 224, 224)):
    """比较多个模型"""
    results = []

    for name, model in models_dict.items():
        print(f"评估 {name}...")
        result = self.evaluate_model(model(), name, input_size)
        results.append(result)

    return pd.DataFrame(results)

# 示例使用
if __name__ == "__main__":
    evaluator = ModelEvaluator()

    # 定义要评估的模型
    models_to_evaluate = {
        "MobileNetV2": mobilenet_v2,
        "ResNet18": resnet18
    }

    # 比较模型
    results_df = evaluator.compare_models(models_to_evaluate)
    print("\n模型评估结果:")
    print(results_df.to_string(index=False))

    # 保存结果
    results_df.to_csv("model_evaluation.csv", index=False)
    print("结果已保存到 model_evaluation.csv")

```

## 7. 针对遥感图像的特定评估

```

# 针对遥感图像的特殊评估指标
def evaluate_for_remote_sensing(model, dataloader, device):
    """评估模型在遥感数据上的性能"""
    model.eval()
    model.to(device)

    total_correct = 0
    total_samples = 0
    total_time = 0

    with torch.no_grad():
        for i, (images, labels) in enumerate(dataloader):
            images, labels = images.to(device), labels.to(device)

```



```
# 测量推理时间
start_time = time.time()
outputs = model(images)
if device.type == 'cuda':
    torch.cuda.synchronize()
end_time = time.time()

total_time += (end_time - start_time)

# 计算准确率
_, predicted = torch.max(outputs.data, 1)
total_samples += labels.size(0)
total_correct += (predicted == labels).sum().item()

accuracy = 100 * total_correct / total_samples
avg_inference_time = total_time / total_samples * 1000 # 每张图片的毫秒数

return accuracy, avg_inference_time

# 使用示例
# 假设有一个遥感数据加载器 rs_data_loader
# accuracy, inference_time = evaluate_for_remote_sensing(model, rs_data_loader,
# device)
# print(f"准确率: {accuracy:.2f}%, 平均推理时间: {inference_time:.2f}ms")
```

## 注意事项

1. **FLOPs计算局限性**: 不同库计算的FLOPs可能有差异, THOP是一个常用但不完美的工具
2. **实际速度测量**: 推理速度受硬件、软件环境、批处理大小等因素影响很大
3. **内存估算**: 上面的内存估算是近似的, 实际内存使用可能因实现细节而异
4. **设备一致性**: 确保所有测量在同一设备上, 以获得可比结果