

# 数据挖掘课程设计报告

姓名: 姚远

学号: 162150109

日期: 2023 年 11 月 12 日

## 1 任务概况

### 1.1 数据集

本课设数据集来自于 [Kaggle](#), 为树叶分类数据集. 该数据集包含 176 类不同的树叶, 18353 张训练样本 (每种树叶至少 50 个样本, 长尾效应并不明显) 以及 8800 张测试样本. 分辨率皆为  $224 \times 224$ . 可以发现这是一个简单的图片分类任务.

### 1.2 评测指标

Kaggle 官方直接采用 [Classification Accuracy](#) 作为竞赛评测标准. 一般在数据类别均衡的情况下, 模型的 acc 越高, 说明模型的精度越好. 本数据集由于类别很多, 数据分布较为平均 (长尾比约 3:1), 因此不考虑通过计算每个类别的 recall 值来评测模型.

### 1.3 解决方案

事实上这是一个课程网站的 project, 该课程是深度学习相关课程, 因此如果你去查看 [Code](#) 就会发现清一色的 Resnet 与其它深度学习方法.

本课设大致分为传统机器学习和深度学习两种方案来解决问题.

传统机器学习解决的过程大致包括:

1. 训练集验证集划分: 考虑到传统机器学习方法性能较优, 故采用 5-Fold Cross Validation, 训练验证比 4:1;
2. 数据预处理: 包括特征提取 (SIFT、HOG etc.)、动态聚类 (Kmeans)、降维 (PCA) 等, 将数据转化为统一的向量形式;
3. 模型分类: 通过采用不同的机器学习分类模型 (KNN、SVM etc.) 来测试分类效果.

深度学习的过程大致包括:

1. 找张 GPU: 使用 Autodl 上的 NVIDIA GTX 2080ti(11GB) 与校内 GPU 服务器的 NVIDIA RTX A5000(24GB) 训练深度学习模型;
2. 训练集验证集划分: 性能原因, 不考虑使用 k-fold, 直接划分. 训练验证比 4:1;

3. 数据预处理: 包括一系列的数据增强 (随机翻转、标准化) 以及采用跨图片增强 (训练过程中体现) 的 trick: MixUp、CutMix;
4. 训练验证模型: 对比采用不同的跨图片增强 trick 训练出来的深度学习模型 (resnet 及其变种) 在验证集上的效果.

## 1.4 DummyClassifier

直接采用 DummyClassifier 作 5 折交叉验证, 预测结果为训练样本中出现次数最多的类, 得出最优的 acc 为 0.036.

# 2 传统机器学习模型

## 2.1 特征提取

### 2.1.1 SIFT 特征提取

SIFT(Scale-invariant feature transform) 中文名为尺度不变性特征变换, 在传统的 CV 算法中拥有很高的地位. 作为一个特征点的提取算法, 它对尺度、旋转、光照等变化不敏感, 导致其效果较优.

首先我们有必要弄明白特征点: 对于平滑的区域一般变化不大, 这不是我们关心的地方, 我们关心的是纹理复杂的地方, 例如边缘、点、角之类的区域, 而这些灰度值变换大的地方就是我们要的特征点.

算法的第一步是数据预处理. 由于彩色图是三通道的, 我们需要先将其转化为灰度图, 此时灰度图为单通道, 灰度值应该在 0-255 之间分布.

算法的第二步是构建多尺度 DoG 空间. 这里的多尺度代表对不同分辨率的图片作空间构造. 对于单个分辨率的图片, 分别作六次不同方差高斯模糊后的图像:

$$\sigma, k\sigma, k^2\sigma, k^3\sigma, k^4\sigma, k^5\sigma.$$

对于特定的方差  $\sigma$ , 高斯模糊后的图像的空间灰度函数为:

$$G(x_i, y_i, \sigma) = \frac{1}{2\pi\sigma^2} \exp\left(-\frac{(x-x_i)^2 + (y-y_i)^2}{2\sigma^2}\right)$$

$$L(x, y, \sigma) = G(x, y, \sigma) * I(x, y).$$

然后需要获取高斯差分函数 (DoG). 高斯差分图像是某一相同分辨率的相邻图像作差值得出, 然后与原图像  $I(x, y)$  作卷积得到 DoG 函数:

$$\begin{aligned} D(x, y, \sigma) &= [G(x, y, k\sigma) - G(x, y, \sigma)] * I(x, y) \\ &= L(x, y, k\sigma) - L(x, y, \sigma). \end{aligned}$$

这样, 对于不同的分辨率, 我们有 5 个 DoG 函数. 事实上这就是高斯金字塔表达, 对于不同的尺度 (分辨率), 从下到上依次降采样, 长的很像一个金字塔.

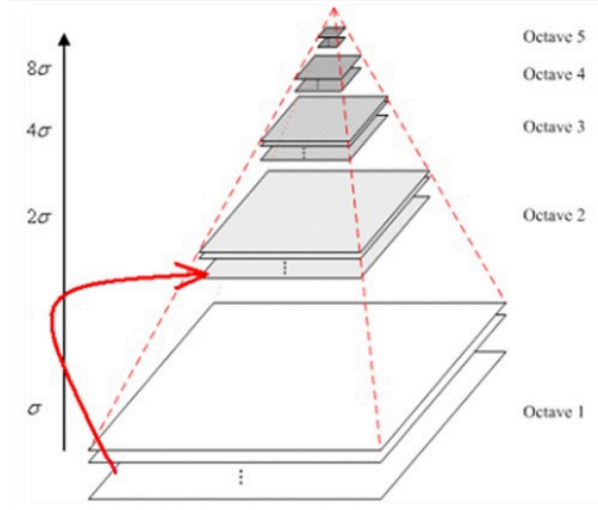


图 1: 高斯金字塔

算法的第三步是极值点检测. 寻找 DoG 函数的极值点, 最简单的想法是直接和附近点相比较. 在 SIFT 算法中我们还需要对同一组中上面和下面不同的 DoG 函数的点相比较, 这样一共需要比较  $8 + 9 \times 2 = 26$  次, 而对于每组 5 个 DoG 函数的情况, 我们只能对中间 3 个函数进行操作.

然而这种方法找到的点的位置都是在整数空间上的. 换句话说, 这些极值点都是离散的. 因此 SIFT 算法考虑通过插值法对离散的 DoG 函数进行曲线拟合, 进一步对方程求偏导, 得到精确的极值点.

此外算法还提到了删除边缘效应的点的概念, 这里不予赘述.

算法的最后一步是产生特征点描述. 上面的过程产生的特征点只有位置信息, 我们需要通过这些特征点附近区域获取特征. 算法在检出的特征点为中心选  $16 \times 16$  的空间作为特征提取空间, 然后将这些区域均分为  $4 \times 4$  个子区域.

对于每个  $4 \times 4$  的子区域, 计算 8 个方向的梯度方向直方图. 首先我们计算每个像素的梯度的幅值和方向, 然后对子区域进行统计, 统计 8 个方向的幅度. (类似 HOG 的特征提取)

幅度公式:

$$m(x, y) = \sqrt{(L(x+1, y) - L(x-1, y))^2 + (L(x, y+1) - L(x, y-1))^2}.$$

方向公式:

$$\theta(x, y) = \arctan((L(x, y+1) - L(x, y-1)) / (L(x+1, y) - L(x-1, y))).$$

然后将所有子区域的幅度直方图连接起来, 即可获得特征向量, 其维度为  $4 \times 4 \times 8 = 128$  维.

### 2.1.2 HOG 特征提取

HOG(Histogram of Oriented Gradient) 中文名为方向梯度直方图, 相比于晦涩难懂的 SIFT, 这个算法更加直接, 直接在原图的基础上产生特征向量. HOG 通过计算并统计图像局部区域内的梯度方向直方图 (SIFT 部分已经介绍过了) 来形成特征, 这导致它对图像几何和光学的变化不敏感.

与 SIFT 算法一样, HOG 也需要先将原图片转换为灰度图, 因为 HOG 提取的是纹理特征, 颜色信息并无作用.

第二步是划分子区域, 在 HOG 中叫做划分 cell, 本课设中每个 cell 为  $8 \times 8$  个像素, 且相邻 cell 不重叠. 然后我们对每个像素计算其梯度幅值和方向, 将所有方向分为 9 个块, 在每个 cell 内统计梯度方向直方图.

第三步, 我们需要将多个 cell 组合成更大的连通块 (block), 将 block 内的所有 cell 的特征向量串联起来得到 HOG 的特征描述. 这里相邻的 block 之间可能重叠. 在更大的范围内统计梯度直方图, 并做归一化处理. 本课设采用 L2-norm normalization:

$$v = \frac{v}{\sqrt{|v|_2^2 + \epsilon^2}}$$

最后, 假设每个 block 包括  $2 \times 2$  个 cell, 相邻的 block 之间有 1 个 cell 宽度的交集. 那么对于我们分辨率为  $224 \times 224$  的图片, 所产生的特征向量的维度应该为:  $(224/8 - 1)^2 \times 9 \times 4 = 26244$  维, 过于庞大. 因此我们先将图片转换为  $128 \times 64$  分辨率 (因为论文就这么搞的), 这样会得到一个  $(128/8 - 1) \times (64/8 - 1) \times 4 \times 9 = 3780$  维的特征向量, 依旧庞大, 后续会进行降维处理.

### 2.1.3 KMeans 动态聚类 + 词袋模型

聚类分析又称群分析, 它是研究对样品或指标进行分类的一种多元统计方法. 聚类分析要使得同一类中的对象之间的相似性比与其他类的对象的相似性更强, 目的在于使同类间对象的同质性最大化和类与类间对象的异质性最大化.

一般地, 聚类分析包括系统聚类和动态聚类. 系统聚类的缺点是每次合并不同类的时候, 结果就已经固定了. 相比系统聚类, 动态聚类法 (即 KMeans) 会根据需求, 动态地调整分类方式.

KMeans 聚类方法最简单的表达如下图所示:

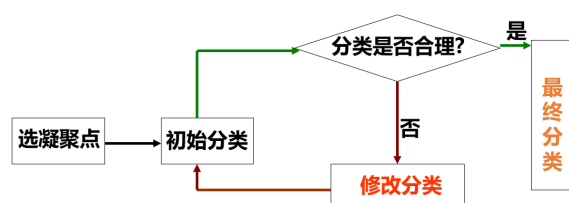


图 2: KMeans 大致过程

KMeans 常用的一种方法是按批修改法, 它的修改原则是使如下的分类函数逐渐减小:

$$l(G_1, \dots, G_k) = \sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{x}_j^i - \bar{\mathbf{x}}^i)^\top (\mathbf{x}_j^i - \bar{\mathbf{x}}^i)$$

暴力枚举的话, 枚举复杂度为:

$$S(n, k) = \frac{1}{k!} \sum_{i=1}^k (-1)^{k-i} C_k^i k^n.$$

为 NP-hard 问题, 因此 KMeans 常采用迭代的方式求解.

迭代算法流程如下:

1. 令  $t = 0$ , 随机选择  $k$  个样本点作为初始聚类中心  $m^{(0)} = (m_1^{(0)}, \dots, m_k^{(0)})$ ;
2. 对样本进行聚类. 对固定的类中心  $m^{(t)} = (m_1^{(t)}, \dots, m_k^{(t)})$ , 计算每个样本到类中心的距离, 将每个样本指派到与其最近的中心的类中, 构成聚类结果  $C^{(t)}$ ;
3. 计算新的类中心. 对上述聚类结果  $C^{(t)}$ , 计算当前各个类中的样本的均值, 作为新的类中心  $m^{(t+1)}$ ;
4. 若迭代收敛或聚类结果符合停止条件, 输出  $C^* = C^{(t)}$ . 否则返回第二步.

由于对于每个样本图片, 我们得出的 SIFT 的特征向量的个数是不定的. 为了进一步将图片转化为单个特征向量, 本课设将训练集中所有图片的所有 SIFT 特征向量进行 KMeans 聚类 (共 64 类), 然后使用词袋模型将其转化为特征向量.

词袋模型 (Bag of Words Model, BoW) 是 nlp 领域中常用到的概念. 对于每个词, 它都对应着一个 one hot 向量, 然后一个句子就可以通过这些 one hot 向量的叠加来表示.

在本课设中, 我们将一张图片的所有 SIFT 特征在聚类中最近的类作为 word, 然后采用 one hot 向量的叠加来构造其特征向量.

#### 2.1.4 KNN

我们首先采用  $k$  近邻分类器进行分类测试. 在这里, 我们使用 5 折交叉验证, 对于所有  $k \in [1, 20], k \in \mathbb{N}$ , 记录 acc 最高的信息, 得出结果如下:

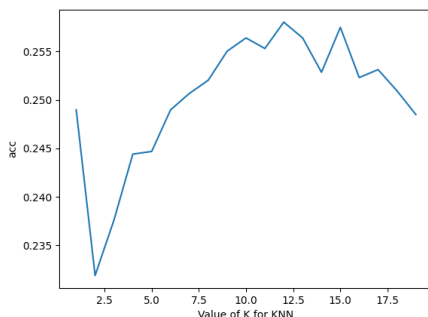


图 3: KNN4SIFT

当  $k = 12$  的时候, 使用 sift 提取特征的 knn 能够达到 0.258 的正确率. 其次是  $k = 15$  的时候能够到达 0.257 的正确率. 总的效果并不是很好.

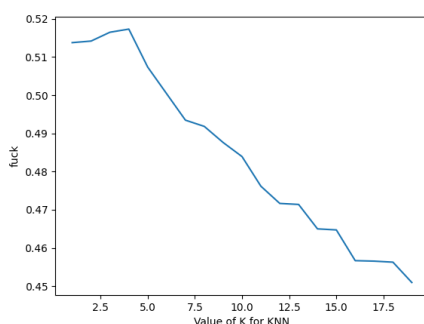


图 4: KNN4HOG

对于 hog 的情况, 在  $k = 4$  的情况最优, acc 达到了 0.517. 相比于 sift 有了巨大进步, 这可能是由于 sift 特征在转化为单一特征向量的时候经过了聚类和词袋化, 再加上图片本身的特征提取也没有 hog 那样稳定, 导致转化为特征向量的时候经过了三层的信息损失, 而 hog 特征提取只经过了 pca 降维, 导致图片特征保存地较为完全.

#### 2.1.5 SVM