

准备：

- √ 英文、数字、特殊符号数据集
- √ PaddleOCR

方案：

1. 将数据集格式改为训练需求的格式。
2. 通过MSER检测选出候选图像文本区域并以二值形式显示（字符旋转、尺寸变化不会影响算法的稳定性），并通过数学形态学将图片中无关区域去掉，尽可能将所有字符区域处于同一连通域，方便字符检测与识别。
3. 分别提取候选图像文本区域的HOG特征（去干扰），接着通过字符判断分类器和字符搜索，实现非字符区域的滤除，其中字符判断使用基于N个描述性特征的SVM实现。（可以通过测试图对分类性能进行测试，得到相应准确率）
4. 测试上述模型性能（（1）.使用N幅测试图进行测试，得到相应分类准确率；
（2）.使用N幅测试图进行字符区域检测定位率测试，得到相应定位准确率）
5. 通过垂直投影图像将字符区域字符的分割提取，并分别以json形式记录每个字符对应的坐标。
6. 先使用KNN算法对于已有英文、数字、特殊符号数据集进行训练，并测试保证准确率不低于95%。然后对于单个字符图片进行分类识别，将识别结果通过字符坐标依次替换到上述json文件中
7. json文件导出csv格式。