




根据取证特征设计训练集进行隐写分析,解决Cover-source mismatch问题

原创

nemoyy  于 2018-07-26 01:02:24 发布  1480  收藏 3

文章标签: [隐写分析](#)

版权声明: 本文为博主原创文章, 遵循 [CC 4.0 BY-SA](#) 版权协议, 转载请附上原文出处链接和本声明。

本文链接: <https://blog.csdn.net/nemoyy/article/details/81212763>

版权

Facing the Cover-Source Mismatch on JPHide using Training-Set Design 论文阅读

1. Abstract

本论文讨论了图像处理流水线(image processing pipeline)对Jpeg隐写中原图片源不匹配(cover-source mismatch)问题的影响,并提出一个取证和隐写分析结合的方法来解决CSM问题。

即,先训练一个多分类器来识别IPP,再选择该IPP对应的训练集来训练隐写分析的分类器。

实验结果表明,取证对隐写没有影响(immune),论文提出的利用了IPP信息的隐写分析的效果比传统的通过扩大数据源的方法要好,同时,实验结果接近对相同来源的数据进行训练得到的分类器结果。

2. 名词解释

2.1 Cover-source Mismatch

简单介绍一下BOSS比赛。

BOSS是一个隐写分析比赛,在BOSS比赛期间,CSM的现象很明显。

训练集:一组原图像和隐写图像对(18,000张),未压缩,大小 512×512 ,来自7个不同相机的灰度图像。

测试集:一个1000个图像的测试集。

但是,测试集中的一些图像来自未在训练集中使用的相机。因此,隐写分析在训练步骤期间得到的图像模型与测试集的图像模型之间遇到不一致。这种不一致被称为原图片源不匹配。

其实就是训练集和测试集样本分布不同,而这里的分布不同是拍摄的相机不同造成的。

2.2 Image Processing Pipeline

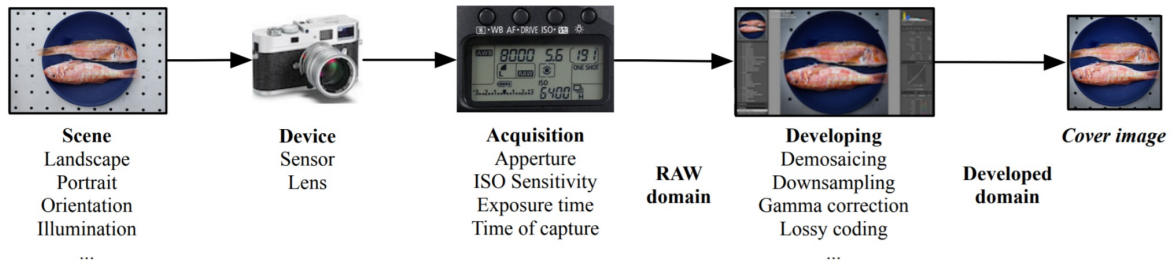


Figure 1: Pipeline of the cover image generation process which can be decomposed into four main steps (scene, device, acquisition, developing) representing parameters of the whole process.

<https://blog.csdn.net/nemoyy>

IPP可以理解为从实景经过相机拍摄得到数字图像的过程.

对同一个场景拍摄得到的结果也可能有很大不同,主要的影响因素有

1. 场景: 拍摄角度,方向,光照等等
2. 设备: 相机传感器,镜头
3. 设备设置: 光圈,曝光时间,快门速度等等

论文中出于控制变量的考虑,针对的变量是在成像阶段的处理,主要有亮度,色差,色温的调节,以及滤波,去噪,锐化等操作.

3. 方法论

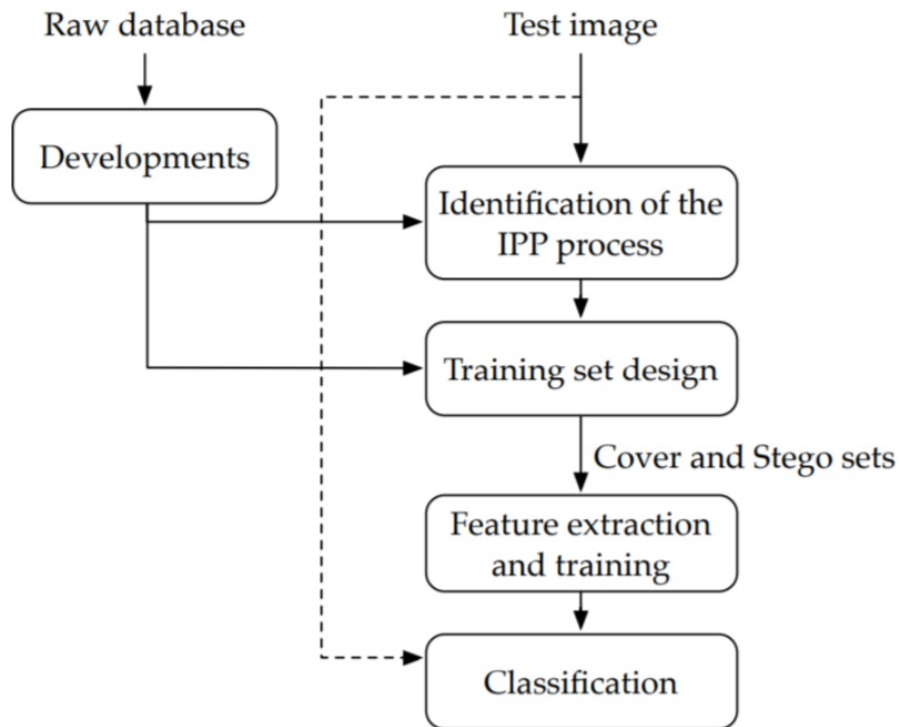


Figure 2: Schematic overview of the examined methodology.

<https://blog.csdn.net/nemoyy>

先确定用在图片的IPP,然后用IPP信息来有针对性的构造一个训练集获得隐写分析器.

通过对原始数据进行不同的成像操作,得到有不同IPP属性的几个数据集.实际操作中,测试图片的IPP属性被预测为最为贴近预先设定的几种IPP中的一种.该IPP对应的数据集就作为训练集,成对生成隐写图片.最后用传统的隐写分析方式,对该训练集,提取特征,构建分类器,用于测试图片.

4. 实验过程

控制变量: 只是用BossBase中徕卡M9相机拍摄的2758张原始图片作为原始数据,用开源软件如DCraw和RawTherapee以及商用软件Adobe Lightroom将原图像转换成彩色Jpeg图片.

用JPHide进行隐写.(对DCT块进行LSB隐写, BLowfish算法生成伪随机序列,确定要改变的DCT系数)

由于是对比IPP的影响,嵌入信息的量保持恒定.

对每一个IPP参数,生成16W对cover/stego图片(经过裁剪的).

实验考虑的情况是隐写方式和嵌入信息大小都是已知的.隐写分析器由集成学习得到.(类似随机森林)

5. 有IPP信息帮助的隐写分析

实验结果证明了利用IPP信息能有效减小CSM问题.

表现在IPP相对应的错误率比IPP不对应的错误率要低不少.

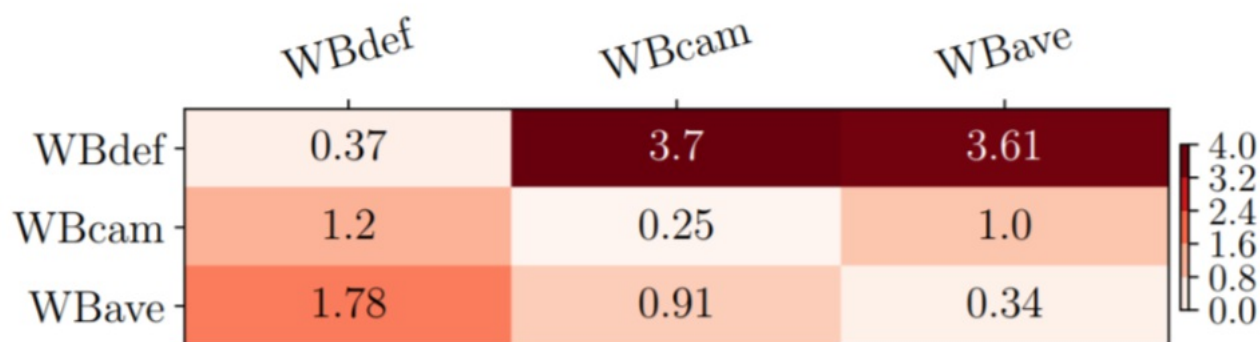


Figure 3: Influence of white balancing (P_E in %).

<https://blog.csdn.net/nemoyy>

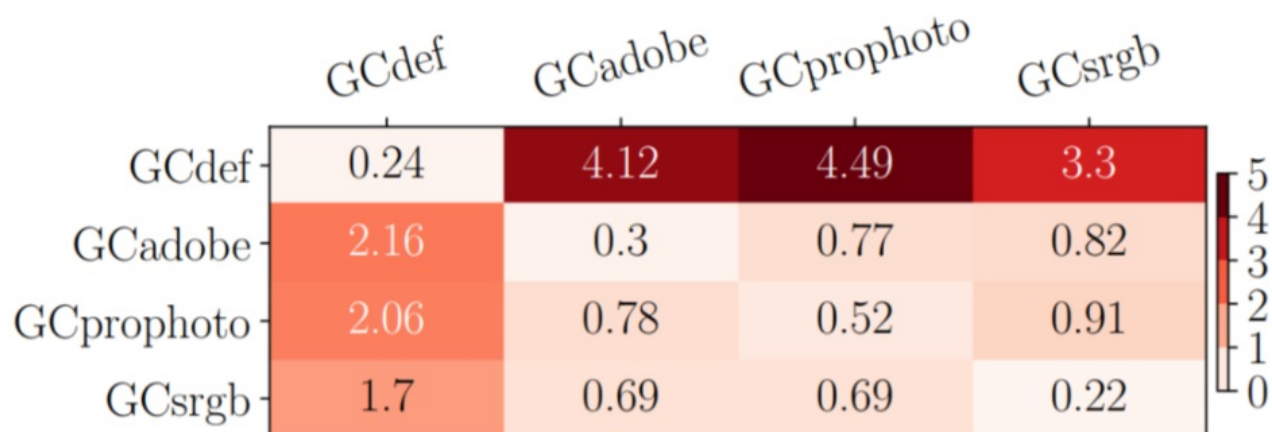


Figure 4: Influence of gamma correction (P_E in %).

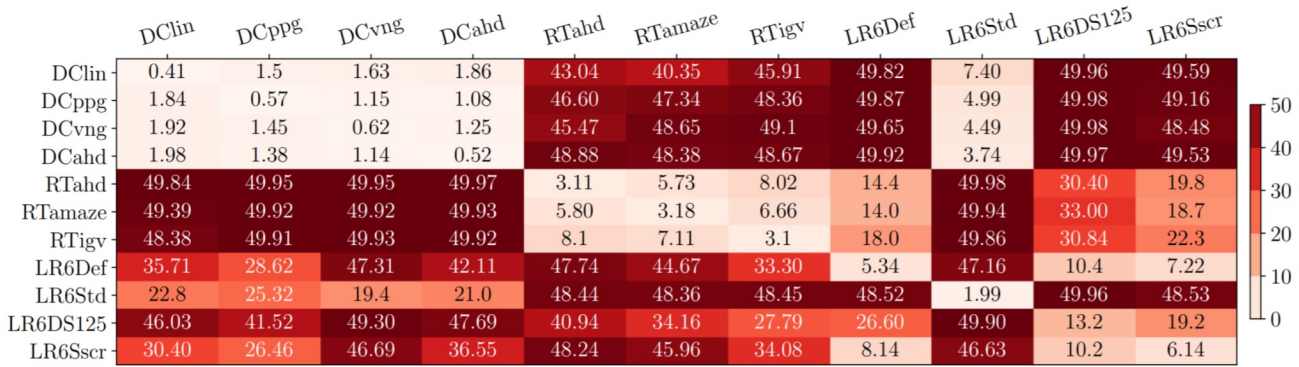


Figure 5: Influence of demosaicing and development software (P_E in %).

6. 训练IPP分类器

模型结构和训练数据其实和训练隐写分析器是一样的. 只是这是一个多分类问题, 可以用类似多分类SVM的方法. 对N个IPP, 训练basic idea是既然隐写分析的特征对IPP的变化有反应说明用同样的特征去对IPP进行分类是可行的.

实验结束(下图)证明这一点.

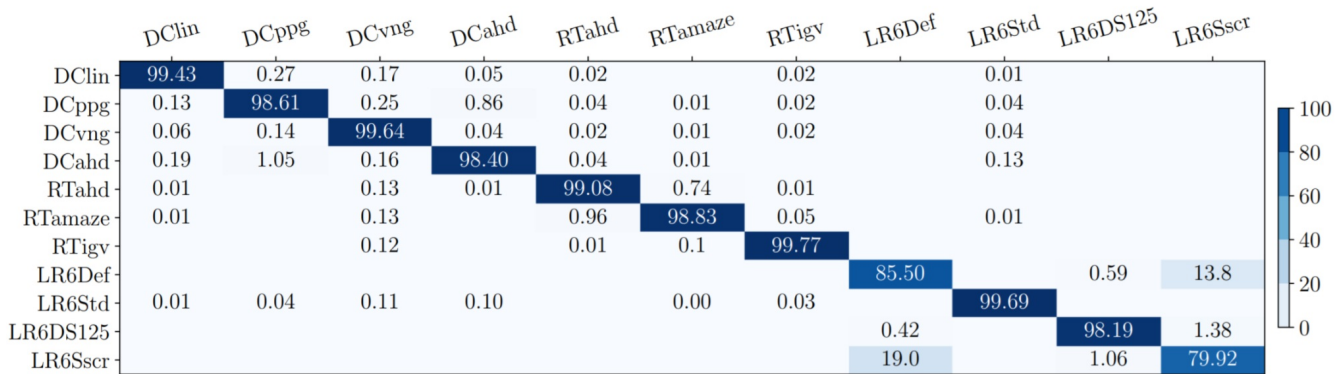


Figure 6: Confusion matrix for the supervised classification of the image processing pipeline (IPP).

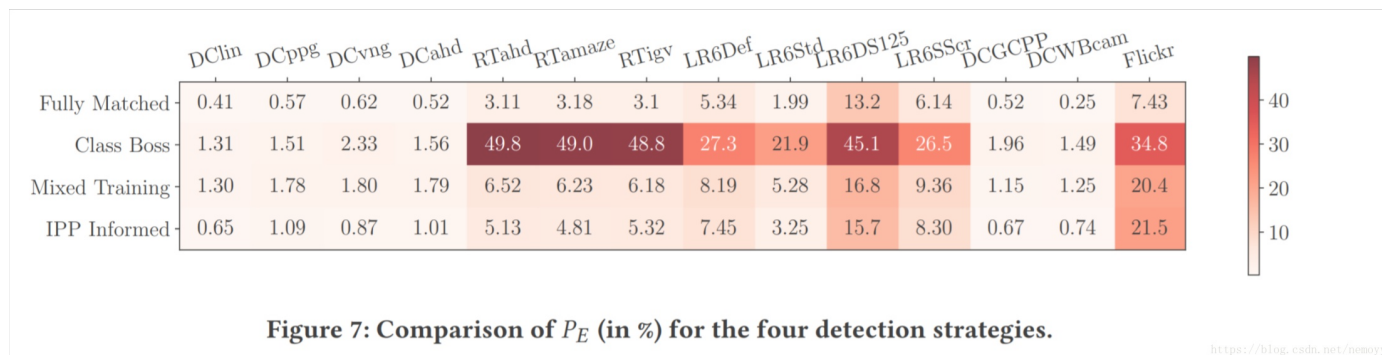
7. 对比结果

Fully matched: 训练集和测试集相同来源. 是比较的baseline, 代表现有的特征和模型能做到最好的结果. 没有CSM问题.

Class Boss: 在BossBass训练, CSM的问题特别明显.

Mixed training: 在所有数据上进行训练, 相当于做了数据增强.

- the results of training the EC on a mix of the 11 sources (consisting of 1000 images from each source) as proposed in [13],
- **IPP informed:** 本论文提出的方法.



8.总结

一般解决训练测试集分布不同问题是通过扩大训练集大小来解决的, 而本文的新意就在所谓的Training-Set design, 同个设计有针对性的训练集来提高测试效果.

参考

1. [Steganalysis with cover-source mismatch and a small learning database](#)
2. [Learning the image processing pipeline](#)