

DeepFool算法论文导读

DeepFool算法论文导读

简介

鲁棒性定义

DeepFool攻击二分类器

原理推导

算法流程

DeepFool攻击多分类器

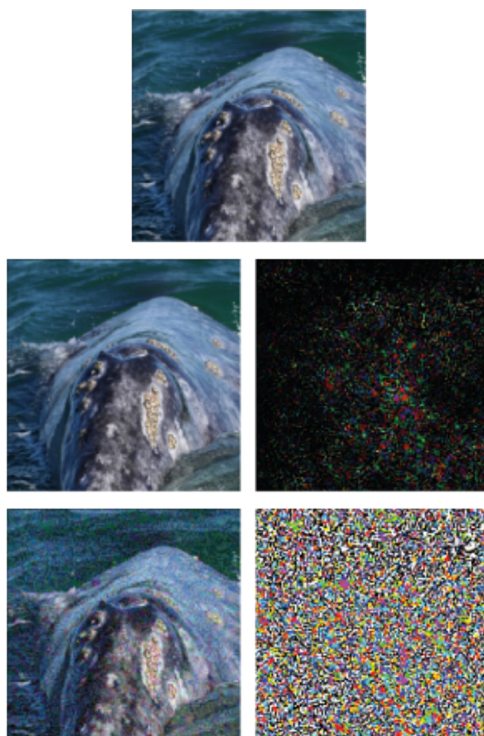
原理推导

举个栗子

算法流程

简介

假设有原始样本 x ，其标签为 y ，我们期望通过对原始样本添加人类不容易发觉的噪声 r ，生成对抗样本 x^* ，使分类器所生成的对抗样本的标签判定为 y^* ，其中 $y \neq y^*$ 。本篇论文则是期望可以找到最小的 r 。



如上图所示：

第一行是原始样本，分类器将其分类为鲸。第二行和第三行分别为Deepfool算法和FGSM算法所生成的对抗样本及它们添加的噪声。虽然这两种算法所添加的噪声都让分类器将样本误判为乌龟，但是可以明显看出Deepfool算法所添加的噪声是比较小的。

鲁棒性定义

给定一个分类器，样本鲁棒性是使得模型出现误分类的最小扰动，具体形式如下：

$$\Delta(\mathbf{x}; \hat{k}) := \min_{\mathbf{r}} \|\mathbf{r}\|_2 \text{ subject to } \hat{k}(\mathbf{x} + \mathbf{r}) \neq \hat{k}(\mathbf{x})$$

其中， \mathbf{x} 为干净的样本， $\hat{k}(\mathbf{x})$ 为模型预测的标签。 $\Delta(\mathbf{x}; \hat{k})$ 为样本 \mathbf{x} 在模型分类器 \hat{k} 的鲁棒性。进而作者又定义出了模型在整个数据集上的鲁棒性，这是一种期望的形式，具体形式为：

$$\rho_{\text{adv}}(\hat{k}) = \mathbb{E}_{\mathbf{x}} \frac{\Delta(\mathbf{x}; \hat{k})}{\|\mathbf{x}\|_2}$$

模型鲁棒性是更好地理解当前网络体系结构的局限性和设计增强健壮性的方法的关键。

DeepFool攻击二分类器

• 原理推导

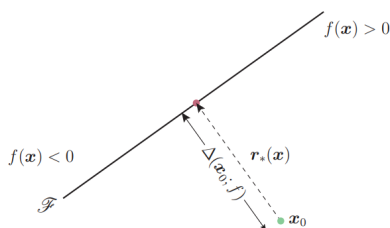


Figure 2: Adversarial examples for a linear binary classifier.

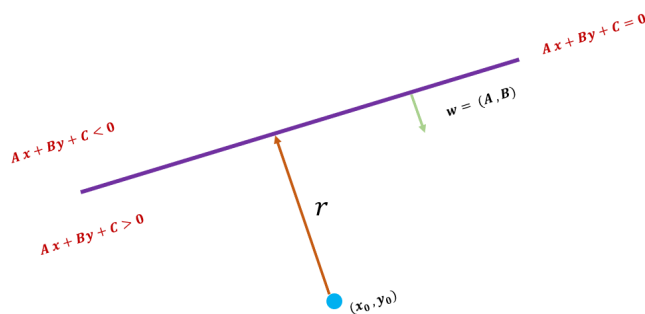
上图为对抗样本攻击线性分类器的图示，其中 $f(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} + b$ 为一个二元分类器。 $\Delta(\mathbf{x}_0; f)$ 为干净样本点 \mathbf{x}_0 的最短距离，即为样本点 \mathbf{x}_0 在分类器 f 中的鲁棒性。 $F = \{\mathbf{x} : \mathbf{w}^\top \mathbf{x} + b = 0\}$ 为分类的超平面。

具体有如下目标函数：

$$\begin{aligned} \mathbf{r}_*(\mathbf{x}_0) &:= \arg \min \|\mathbf{r}\|_2 \\ \text{s.t. } \text{sign}(f(\mathbf{x}_0 + \mathbf{r})) &\neq \text{sign}(f(\mathbf{x}_0)) \\ &= -\frac{f(\mathbf{x}_0)}{\|\mathbf{w}\|_2^2} \mathbf{w} \end{aligned}$$

“

为了更好的理解上述结论，补充以下内容：



如上图所示，在二维平面中，有一条直线 $Ax + By + C = 0$ 和一个点 (x_0, y_0) ，其中直线的法向量为 $w = (A, B)$ ，由点到直线距离知识可知点 (x_0, y_0) 到直线 $Ax + By + C = 0$ 为：

$$\|r\| = \frac{Ax_0 + By_0 + C}{\sqrt{A^2 + B^2}}$$

点 (x_0, y_0) 移动到直线的位移为（需要注意的是位移带方向）：

$$r = -\frac{Ax_0 + By_0 + C}{\sqrt{A^2 + B^2}}e$$

因此二元分类器的公式为：

$$r(x) = -\frac{f(x)}{\|w\|_2^2}w$$

上面目标函数可以通过迭代的方式来进行求解以获得最小对抗扰动，可以重新转换成如下的优化形式：

$$\arg \min_{r_i} \|r_i\|_2 \text{ s.t. } f(x_i) + \nabla f(x_i)^T r_i = 0$$

“

下面为具体的推导过程：

$$r_i = -\frac{f(x_i)}{\|w\|_2^2}w$$

$\nabla f(x_i) = \frac{w}{\|w\|^2}$ ，所以有：

$$r_i = -f(x_i)\nabla f(x_i)$$

又因为梯度的模长为 1，所以两边同乘以 $\nabla f(x_i)^T$ ：

$$\nabla f(x_i)^T r_i = -\nabla f(x_i)^T f(x_i)\nabla f(x_i)$$

最后，移项可得到最终论文中给出的公式：

$$f(\mathbf{x}_i) + \nabla f(\mathbf{x}_i)^T \mathbf{r}_i = 0$$

• 算法流程

根据这个迭代公式有如下算法：

Algorithm 1 DeepFool for binary classifiers

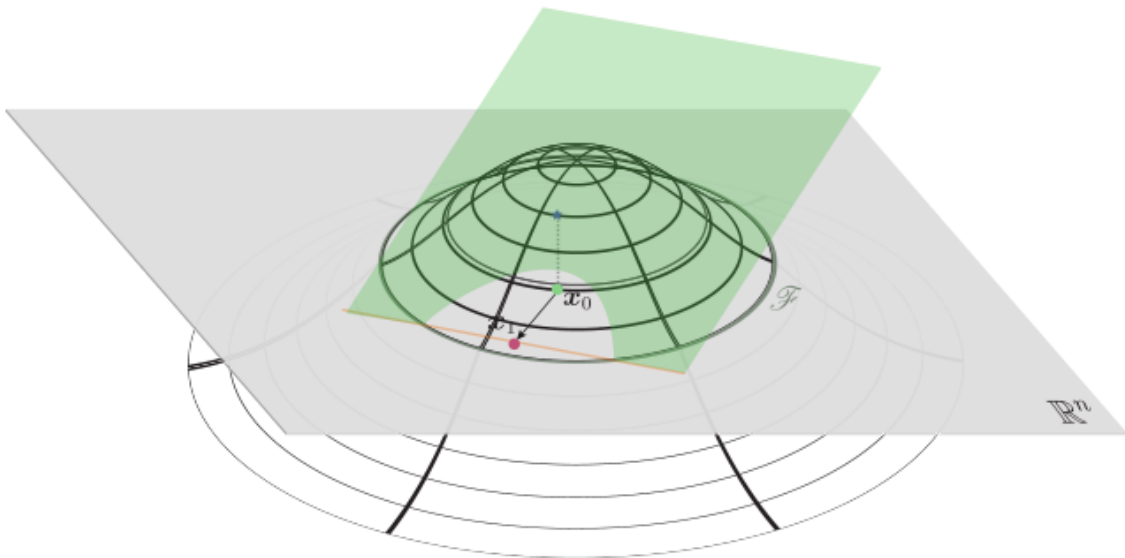
```

1: input: Image  $\mathbf{x}$ , classifier  $f$ .
2: output: Perturbation  $\hat{\mathbf{r}}$ .
3: Initialize  $\mathbf{x}_0 \leftarrow \mathbf{x}, i \leftarrow 0$ .
4: while  $\text{sign}(f(\mathbf{x}_i)) = \text{sign}(f(\mathbf{x}_0))$  do
5:    $\mathbf{r}_i \leftarrow -\frac{f(\mathbf{x}_i)}{\|\nabla f(\mathbf{x}_i)\|_2} \nabla f(\mathbf{x}_i)$ ,
6:    $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \mathbf{r}_i$ ,
7:    $i \leftarrow i + 1$ .
8: end while
9: return  $\hat{\mathbf{r}} = \sum_i \mathbf{r}_i$ .

```

下面是对该算法更形象的解释：

以 $n = 2$ 为例， $\mathbf{x}_0 \in \mathbb{R}^n$ ，绿色平面是 \mathbf{x}_0 点处的切平面为 $f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0)$ ，此时会投影出现迭代点 \mathbf{x}_1 ，然后以此类推，最终将所有的扰动综合起来即为所求的对抗扰动：



DeepFool攻击多分类器

• 原理推导

分类器预测标签如下公式所示：

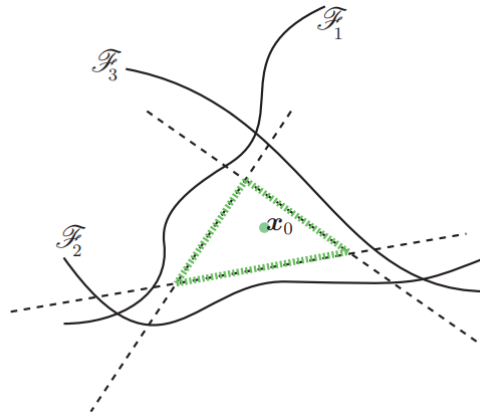
$$\hat{l}(\mathbf{x}_0) = \arg \min_{k \neq \hat{k}(\mathbf{x}_0)} \frac{|f_k(\mathbf{x}_0) - f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_0)|}{\|\mathbf{w}_k - \mathbf{w}_{\hat{k}(\mathbf{x}_0)}\|_2}$$

与 DeepFool 攻击二分类器相似，则多分类器的对抗扰动为：

$$\mathbf{r}_*(\mathbf{x}_0) = \frac{|f_{\hat{l}(\mathbf{x}_0)}(\mathbf{x}_0) - f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_0)|}{\|\mathbf{w}_{\hat{l}(\mathbf{x}_0)} - \mathbf{w}_{\hat{k}(\mathbf{x}_0)}\|_2^2} (\mathbf{w}_{\hat{l}(\mathbf{x}_0)} - \mathbf{w}_{\hat{k}(\mathbf{x}_0)})$$

“

更加通用的多分类器如下：



其中绿色区域可以表示为：

$$\tilde{P}_i = \bigcap_{k=1}^c \left\{ \mathbf{x} : f_k(\mathbf{x}_i) - f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i) + \nabla f_k(\mathbf{x}_i)^\top \mathbf{x} - \nabla f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i)^\top \mathbf{x} \leq 0 \right\}$$

图中实线表示分类器真实的分类超平面，而虚线则代表近似的线性分类超平面，在每次迭代过程中，总是基于当前的迭代值，计算一组近似的线性分类超平面，并根据这组近似超平面，计算扰动，并进行迭代得到下一次的迭代值。

• 算法流程

综合以上 DeepFool 攻击多分类器模型的原理获得如下的算法流程：

Algorithm 2 DeepFool: multi-class case

```
1: input: Image  $\mathbf{x}$ , classifier  $f$ .
2: output: Perturbation  $\hat{\mathbf{r}}$ .
3:
4: Initialize  $\mathbf{x}_0 \leftarrow \mathbf{x}$ ,  $i \leftarrow 0$ .
5: while  $\hat{k}(\mathbf{x}_i) = \hat{k}(\mathbf{x}_0)$  do
6:   for  $k \neq \hat{k}(\mathbf{x}_0)$  do
7:      $\mathbf{w}'_k \leftarrow \nabla f_k(\mathbf{x}_i) - \nabla f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i)$ 
8:      $f'_k \leftarrow f_k(\mathbf{x}_i) - f_{\hat{k}(\mathbf{x}_0)}(\mathbf{x}_i)$ 
9:   end for
10:   $\hat{l} \leftarrow \arg \min_{k \neq \hat{k}(\mathbf{x}_0)} \frac{|f'_k|}{\|\mathbf{w}'_k\|_2}$ 
11:   $\mathbf{r}_i \leftarrow \frac{|f'_{\hat{l}}|}{\|\mathbf{w}'_{\hat{l}}\|_2^2} \mathbf{w}'_{\hat{l}}$ 
12:   $\mathbf{x}_{i+1} \leftarrow \mathbf{x}_i + \mathbf{r}_i$ 
13:   $i \leftarrow i + 1$ 
14: end while
15: return  $\hat{\mathbf{r}} = \sum_i \mathbf{r}_i$ 
```
