

S scusions,小公牛,非国大uthor ol我publation圣sy 33143778年代esearchgste
et publication分析深度学习的鲁棒性与敌对的例子

2018年10月会议论文

DOI: 10.1109/ALLERTON.2018.8636048

cmations读0 2 Lauthor:

JS



小君赵南洋理工大学

76出版物490引文

SEEP

一些作者的出版也在研究这些相关项目:



s | 安全和隐私在AI 10 t系统查看

鲁棒性的分析深度学习与敌对的例子 earning

小君赵南洋理工大学

Abstract-Recent研究显示许多深刻的学习算法的脆弱性敌对的例子,而攻击者获得通过添加微妙微扰良性投入以导致不当行为的深度学习。例如,攻击者可以添加噪音精心挑选一只熊猫的形象,这样产生的图像仍然是一个熊猫人但预计长臂猿的深度学习算法。作为第一步提出有效的防御机制对这种对抗性的例子,我们分析深度学习的鲁棒性与对抗的例子。具体来说,我们是一个严格的最低下界,失真的数据点来获取一个敌对的例子。

关键词可以学习, *deep learning, adversarial examples, robustness,* 敌对的例子,健壮

1. 介绍

最近的研究[1] -[8]提出了许多深度学习系统的脆弱性在敌对的例子,增加攻击者工艺品的微妙微扰良性投入以导致不当行为的深度学习。

敌对的例子中,这些图像或对象识别系统上得到了很多的关注。格拉汉姆·古德费勒等。[1]表明,添加噪音精心挑选一只熊猫的形象可以产生一个图像由一个人仍然认为熊猫,但由神经网络识别长臂猿。Evrim Timov等。[2]演示攻击神经网络系统自动驾驶汽车。攻击者可以把贴纸,如“爱”和“恨”在一个停车标志和技巧相信这是一个限速标志。这可能会导致事故,如果汽车之前不停止的迹象。

敌对的例子在音频领域也存在。Carlini和瓦格纳[3]一个音频敌对的例子通过添加微妙的扰动良性的音频剪辑。例如,对于一个良性的音频剪辑读取“没有

“这rch是由新加坡南洋理工大学(南大)和Alibaba-NTU联合研究所。电子邮件junzhao @alumni. cmu.edu

数据集,本文是无用的”,小噪声添加到剪辑给一个人,同样的句子,但谷歌助理理解结果为“好了谷歌,浏览到邪恶.com”,因此将访问evil.com。

Kurakin et al. [4]介绍敌对的例子在现实世界。例如,添加小扰动图书馆照片使神经网络预测一个监狱,并添加小扰动垫圈图像技巧的神经网络输出受气包。

我们在下面讨论直觉存在的敌对的例子。如图1所示在下一个页面上,我们有一些数据点,一个点代表了一个猫的形象,和一个三角形代表一个狗的形象。假设蓝线是由深上优于铝系统,人工智能认为数据超过这条线作为一只猫,认为数据低于这条线是一条狗。这个基地系统完善工作预测现有的数据点。然而,真正的分布的猫可能看起来像紫色的云的形状,而狗的真实分布可能看起来像黄色的云的形状。因此,数据在这两个虚线部分将由AI预测错误例如,攻击者可能添加噪声绿色数据点,获得红点作为一个敌对的例子。这个对抗的例子将触发在不当行为。具体地说,这张图片是一只猫,但AI系统认为这是一只狗。同样,蓝色三角形代表一个敌对的例子,这是一只狗,但被认为是一只猫

研究的两个挑战attack-resilient AI系统如下:
1)建立弹性基准敌对的例子,和2)提出有效防御敌对的例子。
我们的工作提供了一个初始步骤应对这两个挑战,通过分析深度学习的鲁棒性与对抗的例子。

剩下的纸是组织如下。在第二部分中,我们的定理1展示了结果的最小失真数据点来获取一个敌对的例子。我们提供详细的证据在第四部分III部分调查相关工作,和第五部分总结了纸。

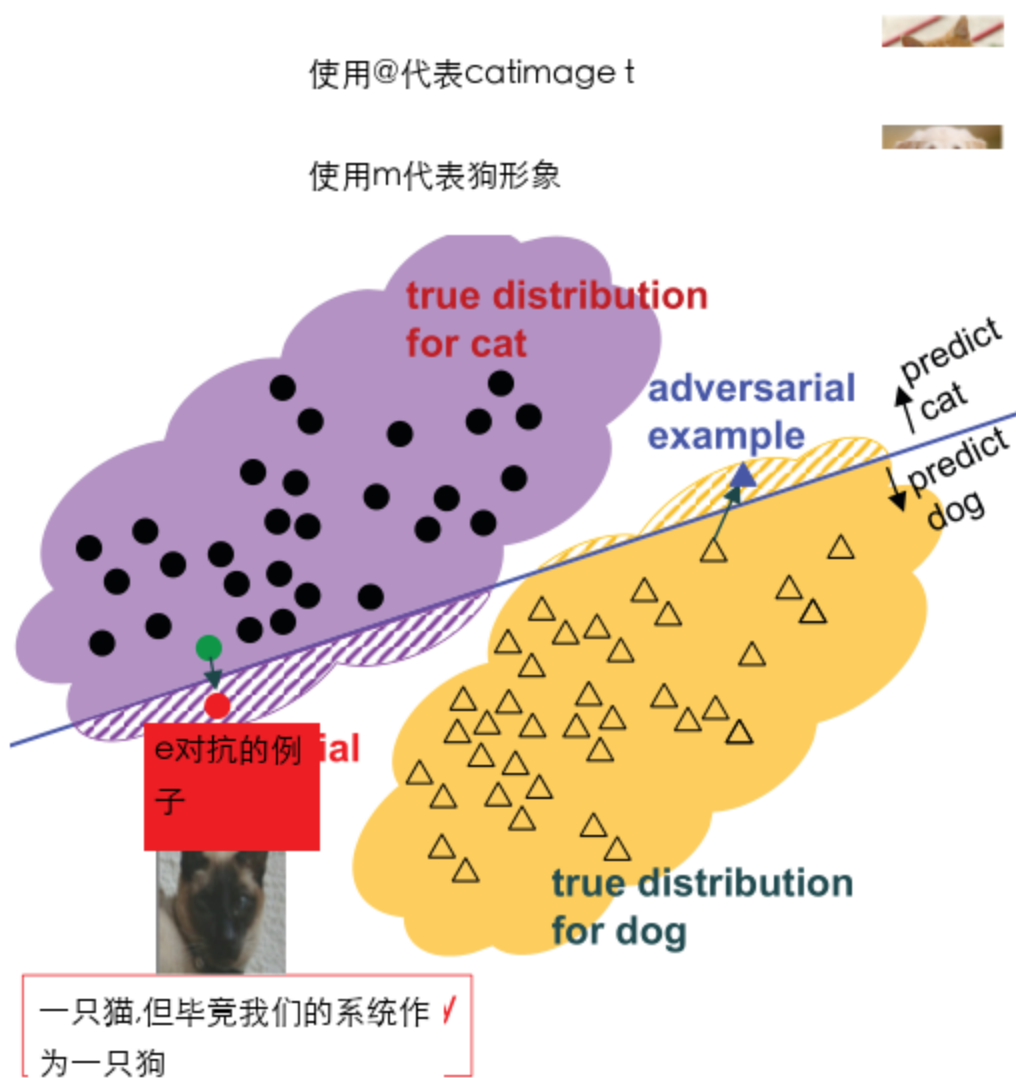


图1 1: An illustration for the existence of adversarial examples

二世。结果最小失真的数据点来获取一个敌对例子

定理1。让D域的数据点,和 $f: D \rightarrow \mathcal{R}^k$ 是一个多层次分类器与 $f = [f_1, \dots, f_k]$ 为连续可微的 f_i ;
 $j \in \{1, 2, \dots, k\}$ 。定义 c 的 f 类分类器预测的数据点 $@$ 即 $@ = \arg \max_{y \in \mathcal{R}^k} f_y(@)$ 。让 δ 表示最低(失真)的数据点 x 来获取一个敌对例子。然后我们可以证明 $\delta \geq \epsilon$

1. $\epsilon \leq \delta$ (1)

$\epsilon \in D$; $\epsilon + 5 = 1$; 即。(1), 这个术语是一个严格的下界

对于 $\delta \leq \epsilon$

并发研究恩和Andriushchenko[9]和“翁等。[10]获得类似的结果,但他们的结果是更复杂的比我们和依靠更多的参与或不同的证明。特别是,与D \mathcal{R}^k 对于一些 d , 定义后 $(@ \mathcal{R}) = \{z: ||z - @||_2 \leq R\}$, 恩

Andriushchenko[9]取

$$e = h(R)$$

$h(R)$ 表示分钟{分钟凝胶2 k}

$$\max_{@ \in \mathcal{R}} \{ ||f(@) - f_j(@)|| \}$$

(2)是更复杂的比(1)紧缩。事实上, $h(R)$ 的特殊情况(2) $R \rightarrow 0$ 减少(1)。翁等。[10]假设李普希茨连续性 f 。 $(z) - J$; (2)在 \mathcal{L} 标准; 即。有我这样

$$||f(z_1) - f(z_2)|| \leq L ||z_1 - z_2||_q$$

1 - 22, 对任何 z 和 z_3 。然后[10]取

$$L @ f(@) \text{ 我 } (12) \quad f_j(x) - f_j(@)$$

证明(1)和(2)是相似的,而后者建立在[9]更相关。相比之下,[10]的证明(3)使用不同的技术。

三世。定理1的证明

假设变形 d , 数据点 $x = @ + \delta$ 改为一个敌对例子。由于 c 类分类器预测的数据点 $@$, f 分类器预测类 $j \in \{1, 2, \dots, k\} \setminus \{c\}$ $@$ 8 对抗的例子。这意味着

存在 $j \in \{1, 2, \dots, k\} \setminus \{c\}$ 这样的 $f_j(@ + \delta) > f_c(@)$

符号简单, 我们定义

$$\text{第9条} \quad \hat{f}_j(@) = f_j(@)$$

s0 9; $(@ + \delta) = f_c$
 $(z + \delta)(4)$ 和(6)的意思

$$f_j(@ + \delta) \geq f_j(@)$$

存在 $j \in \{1, 2, \dots, k\} \setminus \{c\}$
 $g_j(x) < 0$

分析(7), 我们将结合 $g_j(@ + \delta)$ 。为此, 下面我们分析 $g_j(@x) - g_j(@x + \delta)$ 。符号简单, 我们定义

$$y = x + \delta, (8)$$

nd,
ity,

o = (o1, 20日, 24日, ©y = [y1
v2(10)然后是lg; (@) = g; (y
)| = lgj (e, 佻邦) = g5 (pa)
阿, yrswe | 9我(w1, % 2, ..., Ta
)- g' (yl泰, ra) _ | 辛劳a2, 佻邦
)- g5 (Y1, v, 25 - wa) + ...
+ 95 (W12 Ya-1 2 a) = 93
(Y1, y2, d

$$v_d) \int \left| \begin{array}{c} | \\ \vdots \\ | \end{array} \right\rangle \left\langle \begin{array}{c} | \\ \vdots \\ | \end{array} \right| \rangle. \quad (11)$$

$\backslash y, z (z - 9) < g [(y)$
 $v ? 1 \text{泰伊} \infty) = G (Y1, Y2, \text{彝}$
 族人 Til s 5助教)

$$\int |q_i(y_1, y_2, \dots, y_{i-1}, x_i, \dots, x_d)|$$
 。的il、y2il泰Ta) 9 (92 Titly。o Ta
) S | e TL-Li=1中值定理12),存在
 $z \in [\min\{w_i, \text{易建联}\}, \max\{x, y\}]$ 这样
 定义后095" (1,u2,il 2 g,Ta) = B!
 lj (ylvy % 100 Tig年代u-vyx Ta) d
 ala z = 9我(日元。李v2 = 9 (Y1, Y2,
 亚瑟尔Tig1s, Ta) 005 (w2彝族200 T wa
) wi-六世(14)代入(14)(12),我们有l9:(@
)= 9:(y) | = d =>[登月水产)(y)yz,
 v 1 - % 1% _。\')——\ zryxf啊~ (

$$(9) \quad |g - g_j(y)|$$
$$\left(\frac{g_1, \dots, g_n}{da} \right) [$$

[Zfi管理部(m)yzw堪称120y” ,其中

{Zux, W] -le-yl,
0 =替换(19)和(20)(
16),我们获得gi (®) -
g; (W) | <马克斯{[
失效切换后2 ||q | 9 g
; (2)∈D} x [l ~我
们, 。 @D回忆y: = @ +
&(8),写(21),g (x
+ &) - g;失效切换后=
2 ||q (=)) (2)∈D
} x 18],这意味着gi (

从我+%=1,持有人不等式的应用

这在(22)

$$\text{存在 } j \in \{1, 2, \dots, c_k\} \setminus \{i\}, g_j(z) - \max \{ |dg_j(z)| : l: \# \in D \} x \leq 8 \}, < 0. \quad (23)$$

然后我们从(23)获得

$$\ln \frac{a_j(x)}{a_i(x)} > 9 \text{ 分钟} \cdot \text{圣恩基} 2 \setminus (\text{ch 马克斯} \{g_j\} \cdot \quad (24)$$

召回的定义 $g_j(x) := \text{铁}(\text{@@}) = (5) \text{fi}$ (写(24))

18

$$\begin{aligned} & > \text{分钟我} \{1, 2, kP\} \{e\} & \text{我}(\text{@}) - \text{£} & x) \\ & \text{马克斯} \{H\% \{ & 1z\} = \{()\} = \in D\} & \end{aligned} \quad (25)$$

让 $\{1, 5\}$ “\” 表示最低 £ 失真的数据点 @ 获得一个敌对的例子。然后(25)

$$\begin{aligned} & (Mi \quad \text{Adv} & 1) = 1 - 1z) = \\ & & () = \in D\} \end{aligned}$$

四。相关工作不同的敌对的例子

Athalye和Sutskever[6]显示3 d的存在敌对的例子。一种新颖的算法合成敌对的例子也提出选择分布的转换。谢里夫等。[11]目前对抗性的例子对面部生物识别系统。Cubuk等。[7]认为神经网络固有的不确定性的预测作为对抗的例子存在的原因。Baldaer等。[12]利用凸规划产生敌对的例子在不同约束条件和不同的网络类型。阮等。[13]产生图像面目全非的人类,但辨认款。Moosavi-Dezfooli等。[14]介绍DeepFool算法构建有效地对抗的例子。赵等。[15]定义一个框架来生成清晰和自然对抗的例子。Brendel等。[16]讨论的意义攻击这只依靠最后的决策模型。

b.鲁棒性与对抗的例子

科特勒和黄[17]介绍的方法来训练深ReLU-based分类器,这样他们证明地健壮与有界准则对敌对的例子。黄等。[18]对抗的防御机制扩展到更大的模型。Alvarez-Melis和Jaakkola[19]表明,解释的可靠性是至关重要的可解释性。在较弱的假设,加和史密斯[20]证明理想化的模型并不容易受到敌对的例子。王等。[21]使用偏差方差理论理解敌对的例子。Raghunathan等。[22]利用半定松弛分析对手的例子。王et al。[23]从拓扑概念应用于量化敌对的例子。并发研究恩和Andriushchenko[9]和翁et al。[10]获得结果与本文类似,但他们的结果是更复杂的比我们和依靠更多的参与或不同的证明。法等[24]讨论一组不同的扰动。Carlini和瓦格纳et al。[25]认为防守蒸馏[26]并不在一些情况下工作得很好。

c.防御敌对的例子

Papernot等。[26]提出防御蒸馏来对抗对抗的例子。黄ef。[27]现在学习的方法与一个强大的对手,健壮的分类器在哪里学会了从监督数据。Kurakin等。[28]规模ImageNet对抗训练。格拉汉姆·古德费勒ef。[1]强调线性神经网络的本质是敌对的例子的一个原因。顾和Rigazio[29]调查敌对的例子的结构并提出相应的防御。孟和陈[30]磁铁防御框架,它可以近似正常的多方面的例子。歌等。[31]提供了一个实证评价敌对的例子。他et al。[32]结合多个防御获得一个强大的防御机制。电车等。[33]表明,对抗训练无法抵御黑箱的攻击。Cullina等。[34]扩展可能大约正确(PAC)则将框架来解决敌人的存在。

诉的结论

许多最近的研究已经证明了深度学习算法的脆弱性敌对的例子,增加攻击者产生的微妙的噪音良性的输入,以导致错误的预测深度学习。例如,一个对手可以把贴纸,如“爱”和“恨”停车标志和技巧深度学习系统识别限速标志。在本文中,我们分析深度学习的鲁棒性与对抗的例子。具体来说,我们是一个严格的最低下限

失真数据点来获取一个敌对例子。未来的发展方向包括有效的防御机制的发展各种神经网络模型来对抗对抗的例子。

引用

- 2 I. 格拉汉姆·古德费勒, Shien和c. Szegedy”解释和利用对手的例子。“在国际会议上学习(ICLR), 2015。
- 1 I. Evtimov, K. Eykholt, 费尔南德斯, Kohno, B. A. 普拉卡什, A. R. Curran, 和“音频敌对例子:有针对性的攻击“语音”, arXiv预印本arXiv: 1801.01944, 2018, a. Kurakin 1。
- 3 j. 格拉汉姆·古德费勒, Bengio”对抗的例子在现实世界数据集上, Z. Zilgner, A. Akbamejad, G. H. N. Nemann”对抗攻击分类模型图。“在ACM SIGKDD, 2017。
- 4 答: Athalye和I. Sutskever”生成对抗性例子: robust adversarial examples”合成健壮敌对例子。“arXiv预印本arXiv: 1802.08897, 2018。
- 米 e. d. Cubuk, b. Zoph, s. S. Shojanazeri, Schoenholz问: 诉勒, “有趣的属性对抗的例子”, arXiv预印本arXiv: 1802.08897, 2018。
- 1 a. Athalye, n. Carlini, d. 瓦格纳, “混淆梯度给一种虚假的安全感: 绕过防御敌对例子。”“在会议上神经信息处理系统(NIPS), 2017, 页。2266-2276。
- 9 m. 海因和m. Andriushchenko”正式保证分类器的鲁棒性与敌对的操纵, “在会议上神经信息处理系统(NIPS), 2017, 页。2266-2276。
- [1 线性调频。翁, 张, h. py. 陈, 易建联。d. 苏, y. 高, C. J. 谢长廷, l. 丹尼尔,”评估神经网络的鲁棒性: 一个极端值理论方法。“在国际会议上学习表示(ICLR), 2018。
- 一个 m. 谢里夫(s. Bhagavatula, l. 鲍尔和m. k. Reiter”装饰10犯罪: 真正的和隐形的攻击”, 在ACM SIGSAC计算机支持安全, 2017。
- 2] e. r. Balda, a. Behboodi, r. Mathar”代对抗的例子使用凸规划, “arXiv预印本arXiv: 1802.08897, 2018。
- 3] 阮, j. Yosinski和j. Clune”神经网络easil愚弄: 高信心预测面目全非形象”, 在IEEE计算机视觉与模式识别会议(CVPR), 2015, 页42。
- (1 S.-M. Moosavi-Dezfooli, A. 法和p. Frossard”DeepFool: 一个简单而精确的方法愚弄神经网络”, 在IEEE计算机视觉与模式识别会议(CVPR), 2016, 页1651-1659。
- [1 d. Dua, z. 赵, 辛格, “生成自然对抗的例子。”“在国际会议上学习(ICLR)表示, 2018年。
- 6] w. Brendel, j. raub和m. 陆慈”决定对抗攻击: 可靠的攻击黑盒机器学习模型。“在国际会议上学习(ICLR)表示, 2018年。
- a7 j. z. 科特勒和e. Wong”可证明的防御敌对例子通过凸外敌对的多面体。2017年“arXiv预印本arXiv: 1706.00667, 2017, y. s.

- [1 e. Wong, e. 施密特, h. Metzen和j. z. 科特勒”比例可证明的对抗性的防御。“arXiv预印本arXiv: 1806.00667, 2018, y. s.
- 19 d. Alvarez-Melis和t. s. 拉克洛”可解释性方法的鲁棒性”, arXiv预印本arXiv: 1806.00667, 2018, y. s.
- 20 y. 加和·史密斯,”充分条件理想化模型没有对抗的例子: 贝叶斯神经网络的理论和实证研究。“arXiv预印本arXiv: 1806.00667, 2018, y. s.
- [21] Jha和·乔杜里,”鲁棒性的分析最近的邻居对抗”, arXiv预印本arXiv: 1706.03922, 2017。
- [22] a. Raghunathan, j. 斯坦哈特, p. 梁”认证防御敌对例子。“在国际会议上学习(ICLR)表示, 2018。
- [23] j. b. Wang, 高, 和y. 气,”鲁棒性的理论框架(深)分类器对敌对例子, “在国际会议上学习表示(ICLR)研讨会, 2017。
- [24] A. Fawzi, s. m. Moosavi-Dezfooli, p. Frossard”深度网络的健壮性”, 在IEEE安全与隐私(SP), 2017年。
- [25] n. Carlini, d. 瓦格纳, “评估神经网络的鲁棒性。”“在IEEE研讨会上安全性和隐私(SP), 2017年。
- [26] n. Papernot, p. 麦克丹尼尔, s. Jha, x. Wu, 阁下”防御敌对的扰动对深度神经网络: 安全 and 隐私”, 2016年, pp58: 黄, r. b.
- 27 d. 舒尔曼和c. Szepesviri”学习与一个强大的对手。“arXiv预印本arXiv: 1511.01241, 2015。
- 12 a. Kurakin, l. 格拉汉姆·古德费勒, 美国Ben 请愿的机器”, 在会议上神经信息处理系统(NIPS), 2017, 页。2266-2276。
- [29] 顾和, l. Rigazio”深度神经网络架构robustness”, 在会议上神经信息处理系统(NIPS), 2017。
- 30 孟, d. h. 陈, “磁铁: 双管齐下的防御敌对例子。”“在ACM计算机和通信安全, 2017年, 页。1651-1659。
- 31 y. 歌, t. Kim, 8. Nowozin, 美国Ermon和n. Kushman”PixelDefend: 利用生成模型来理解和抵御敌对例子。”“在国际会议上学习(ICLR)表示, 2018。
- [32] w. 他, j. 魏, n. Carlini, x. Chen和d. 歌,”防御敌对例子: 乐团疲弱的防御并不强烈, “在U.S. 计算机和通信安全, 2017。
- [33] 有轨电车, a. Kurakin, n. Papernot, l. 格拉汉姆·古德费勒, 和p. Frossard”防御敌对的培训: 攻击和防御。”“在国际会议上学习(ICLR)表示, 2018。
- 13 d. Cullina, a. n. Bhagoji和p. 米塔尔”PAC-learning逃避对手的存在。”“arXiv预印本arXiv: 1806.00667, 2018, y. s.