# Cyber Security and AI

## Britta Hale

# Does the training work?



*Accuracy*

Can the training be circumvented?

Can the model be misinterpreted?

Can the model be abused?

*Everything Else*

# A perfect memory…

# Attacks During Training

e.g.

- Poisoning

- Trojaning/Backdoors

# Poisoning

***Integrity***

- Confidence reduction

  do not change a class but highly impact the confidence

- Misclassification

  change a class without any specific target

- Targeted misclassification

  change a class to a particular target



| UAVs |
| Birds |

# Poisoning

*Availability*

- Source/target misclassification

    change a particular source to a particular target

- Universal misclassification

    change any source to particular target



"panda"
57.7% confidence

$+.007 \times$

"nematode"
8.2% confidence

$=$

"gibbon"
99.3 % confidence

**Shifts classifier boundary**

Fig. 1. Linear SVM classifier decision boundary for a two-class dataset with support vectors and classification margins indicated (left). Decision boundary is significantly impacted if just one training sample is changed, even when that sample's class label does not change (right).

Miller, Xiang, and Kesidis, 2019

# Trojaning/Backdoor

1. Inverse network to create a trojan trigger

2. Retrain model with malicious data

3. Real inputs which activate the trojan trigger generate malicious behavior

Access to original dataset not necessarily required

Retraining can take minutes/hours (vs. weeks/months for original model)

speedlimit 0.947

STOP

# Defense

- Outlier detection

  **How to define an outlier?**

  **What about data that was injected before filtering rules?**

- Test newly added training samples against current model for accuracy

  **What about trojans?**

# *DATA*

ISSIE LAPOWSKY   SECURITY   03.17.2018 12:20 PM

# Cambridge Analytica Took 50M Facebook Users' Data—And Both Companies Owe Answers

*The New York Times*

## Facebook and Cambridge Analytica: What You Need to Know as Fallout Widens

By Kevin Granville

March 19, 2018

f

# Attacks During Production

e.g.

• Inference

• Evasion

# Inference

- Acquire information about dataset

# Inference

- Acquire information about dataset
- Membership inference / data attributes



Pixabay: DokaRyan

# Inference

- Acquire information about dataset
- Membership inference / data attributes
- Model Extraction



Tramér et al. Usenix Security '16

# Evasion

Does not shift classifier boundary, but pushes poisoning into dataset

# Evasion

To classify Pandas as Gibbons:

1. Change a bunch of Gibbons closer to Pandas

2. Keep Gibbons labelled as Gibbons

3. Add changed Gibbons to training pool

Does not shift classifier boundary, but pushes poisoning into dataset

## Training set

| | Input | Label |
|---|---|---|
| Clean target instances | | "spam" |
| Clean base instances | | "not spam" |
| Poison base instance(s) | | "not spam" |

Train ↓

Target instance → Test → DNN → Prediction → "not spam"

Shafahi, Huang, Najibi, Suciu, Studer, Dumitras, Goldstein, NIPS, 2018

# Defense

- Differential Privacy

# Differential Privacy

Goal: Try to hide individual data points

# Differential Privacy

## *Problem: Model may become disbalanced*

# Defense

- Differential Privacy

- Don't force guessing ("null" class)
  **Human overhead**

- Adversarial training
  **What if the adversary uses different examples?**
  **What if you train on too many adversarial examples?**

# Poisoning vs. Evasion