# CAI 4104/6108 – Machine Learning Engineering:

## Adversarial ML & Privacy Threats

Prof. Vincent Bindschaedler

Spring 2024

# Administrivia: Project

- Due Wednesday 4/24 by 11:59pm    *no late penalty if submitted by 4/26 11:59pm*
  - Deliverables:
    - Final report (PDF) — 3+ pages
    - Code (ZIP)
  - Submit as a **group**
    - Use Canvas project groups
  - See instructions on the Canvas assignment for more details

- Evaluation criteria
  - Depth, soundness, presentation quality, and **effort**

- Important considerations
  - What are the results? How do they compare to the baseline(s)?
  - Did you follow **best practices**? Does your evaluation methodology make sense?

# Administrivia: Project

- **Report:**
  - ◆ **Introduction**
    - ❋ What the project is about? What problem are you trying to solve?
  - ◆ **Approach: Dataset(s) & Pipeline(s)**
    - ❋ What is your proposed approach? What are you doing to solve the problem? What ML techniques are you using? What dataset(s) are you using?
  - ◆ **Evaluation Methodology**
    - ❋ How are you evaluating your approach? How did you split the data? What are the metrics/baselines
  - ◆ **Results**
    - ❋ What results have you obtained? How do your results compare to the **baselines**?
    - ❋ Include: **tables** or **plots**
  - ◆ **Conclusions**
    - ❋ What are your conclusions?

# The Stationarity Assumption

- Many (most?) learning methods make the *stationarity assumption*
  - *The training data and testing (evaluation) data come from the **same** distribution*

- From the perspective of machine learning theory
  - The stationarity assumption makes sense
    - If the testing data comes from a different distribution, can we say anything about generalization?

- What about in the real-world (the world of deployed systems)?
  - Does the stationarity assumption hold there?
  - No! At least not always. Examples?
    - Distribution of data changes over time
    - Better data becomes available, or what we care about changes
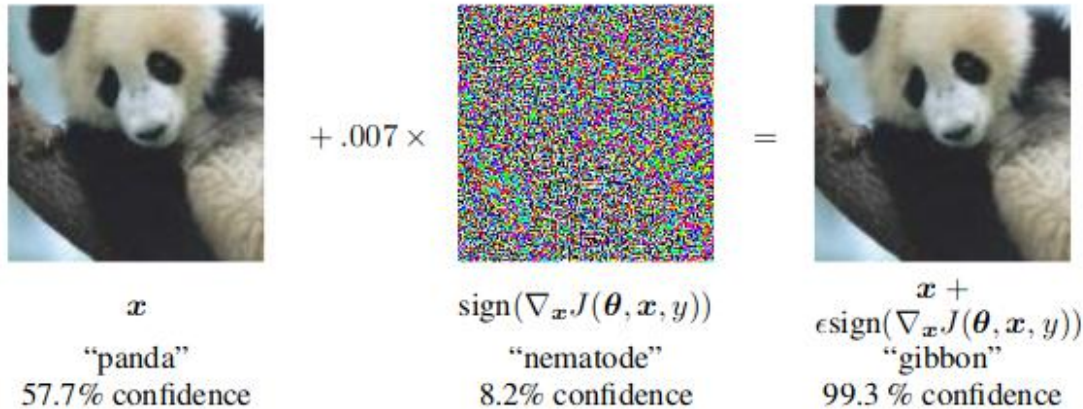  - The stationarity assumption may not hold in **adversarial environments**

What's this animal?

What's this animal?

It's a gibbon!

$$+ .007 \times$$

$$\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

$$=$$

$$\boldsymbol{x} + \epsilon\,\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$$

$\boldsymbol{x}$
"panda"
57.7% confidence

"nematode"
8.2% confidence

"gibbon"
99.3 % confidence

*Source: Goodfellow, Shlens, and Szegedy. "Explaining and harnessing adversarial examples." ICLR 2015*

# Robustness

- What does **robustness** mean?

  - In computer science:
    - The ability of an algorithm/system to handle errors in execution or in its input
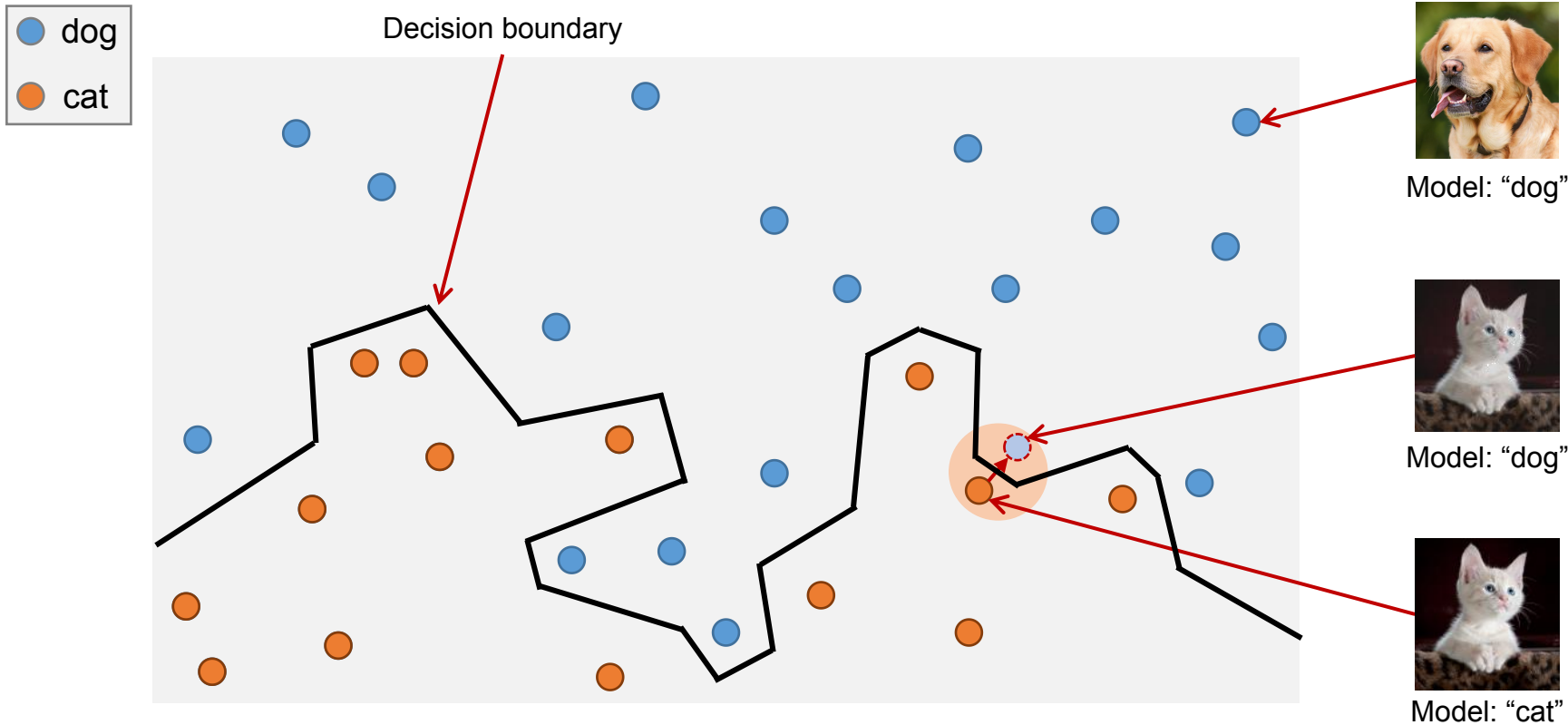
  - In machine learning:
    - testing error ≈ training error
    - low generalization error
    - performs well even on unexpected inputs, noisy inputs or outlier inputs

- **Adversarial robustness**

  - We want robustness even for the **worst-case adversarial inputs**
  - We assume the adversary chooses the inputs

# Adversarial Robustness: Intuition

# Adversarial Examples: Terminology

- *Adversarial sample* or *adversarial example*
  - ◆ Malicious input designed to fool a machine learning model

- *Adversarial robustness*
  - ◆ Robustness to adversarial (i.e., malicious) inputs
  - ◆ Note: (traditional) robustness means robustness to unexpected inputs or outlier inputs
    - ✺ Unexpected / outlier ≠ malicious

- *Adversarial perturbation*
  - ◆ Perturbation of a **benign input** into an **adversarial example**
  - ◆ In the ideal case (for the adversary) the perturbation is imperceptible to humans

# Evasion Attacks

- Goal:
  - Adversary aims to avoid detection by manipulating malicious test samples

- Application scenarios
  - Spam filtering: attacker crafts a malicious spam email in such a way that it appears to be legitimate
  - Malware detection: attacker takes a piece of malware and modifies it so that it is detected as benign
  - In such scenarios the stationarity assumption may not hold
    - Adversaries that manipulate the test data are realistic in this context

# Adversarial Examples: FGSM

- Fast Gradient Sign Method (FGSM)
  - ◆ Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." arXiv, 2014.
  - ◆ No optimization, just compute the gradient:
    - ❋ Let $x' = x - \varepsilon \, \text{sign}(\nabla L_{f,t}(x))$
    - ❋ Here $\varepsilon > 0$ is ch                                    d $\nabla L_{f,t}(x)$ is the **gradient** of the
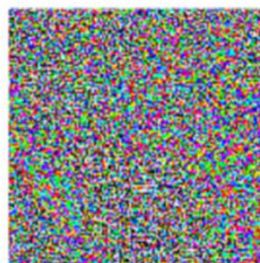  - ◆ Intuition:
    - ❋ The **gradient**                                              hould be changed to minimize the
    - ❋ The attack shif



$x$
"panda"
57.7% confidence

$+ .007 \times$

$\text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"nematode"
8.2% confidence

$=$

$\boldsymbol{x} + \epsilon \, \text{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$
"gibbon"
99.3 % confidence
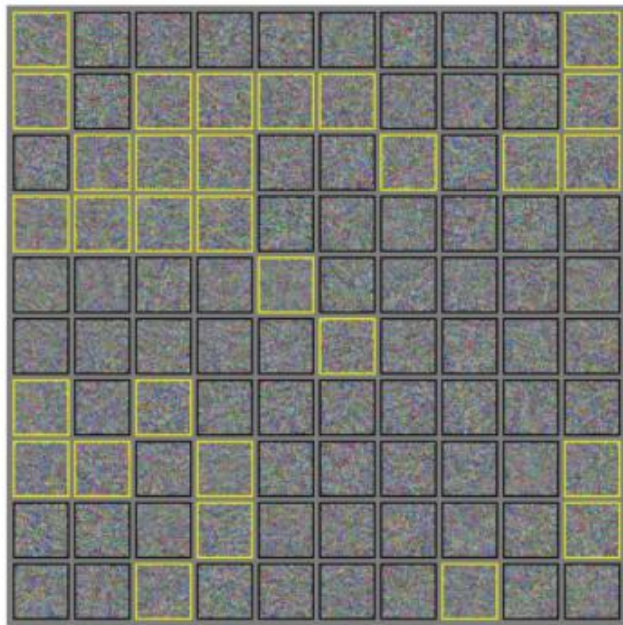
- Nguyen, Yosinski, Clune. *"Deep Neural Networks are Easily Fooled: High Confidence Predictions for Unrecognizable Images."* IEEE CVPR 2015.

# Neural Nets: Other Weird Properties

- Goodfellow, Ian J., Jonathon Shlens, and Christian Szegedy. "Explaining and harnessing adversarial examples." *ICLR 2015.*
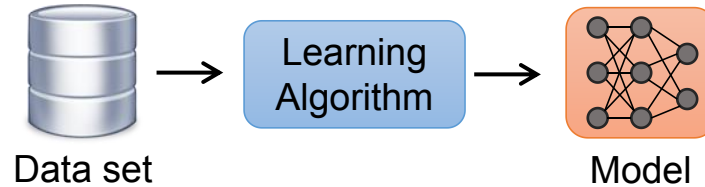


**Can you see the airplanes?**

Randomly generated fooling images for a CIFAR-10 convolutional neural net

Each image is generated by:

- Drawing an isotropic Gaussian
- Taking a step in the direction that increases the probability for "airplane"
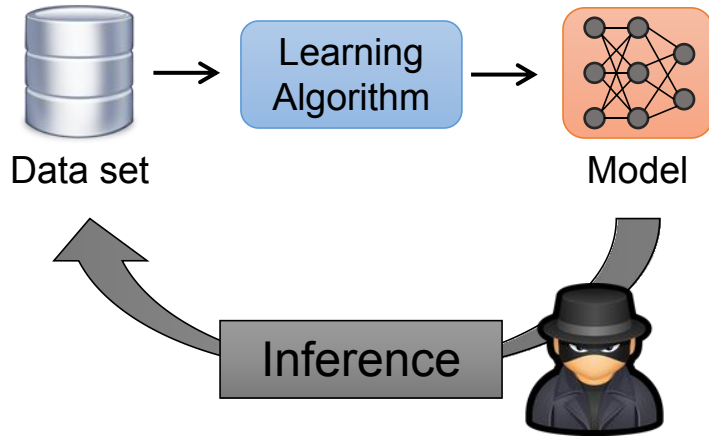
Yellow box: confidence of "airplane" above 50%

Data set → Learning Algorithm → Model

- What if the model's training data is sensitive?
  - We want to keep it private

- We also want to publicly release the model
  - But the model is a function of the training data!
  - What do we do?

# Privacy Attacks on ML Models





source: xkcd

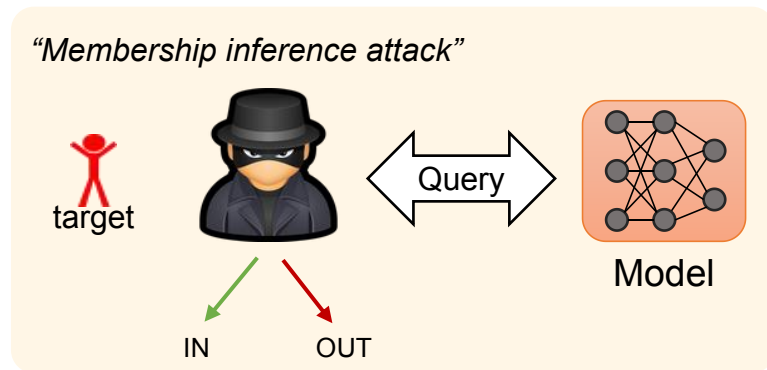- What can be inferred about the training data from access to the model?

# Membership Inference Attacks

- Empirical observation
  - Complex ML models tend to memorize their training data (even if they do not overfit)
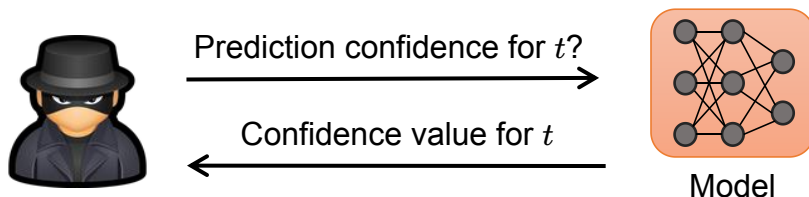  - We can quantify this through **membership inference**

- Attacker's goal
  - Determine whether a given target t's record was part of the target model's training data set
  - Hypotheses:
    - (Member) $H_{IN}$: t is in the training data
    - (Non-member) $H_{OUT}$: t is **not** in the training data
  - Assumption: adversary knows t's data record



*"Membership inference attack"*

target    Query    Model

IN    OUT

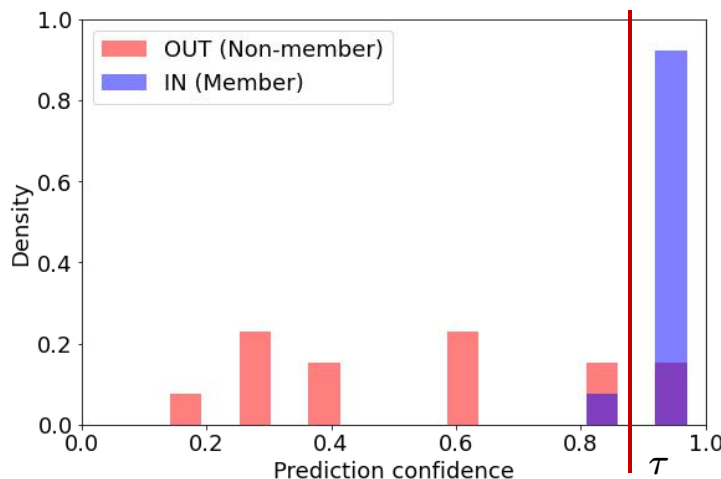■ Black-box Membership Inference Attack:



**Membership Inference Attack:**
Inputs: model $\mathcal{M}$, record $t$, threshold $\tau \in [0,1]$
Output: **IN** or **OUT**
Procedure:
- $c \leftarrow \mathrm{prediction\_confidence}(\mathcal{M}, t)$
- If $c \geq \tau$: return **IN**
- Else: return **OUT**

Prediction confidence for $t$?

Confidence value for $t$

Model

Note: here the confidence value is just the predicted probability for the true class

■ Empirical observation: LLMs memorize some of their training data

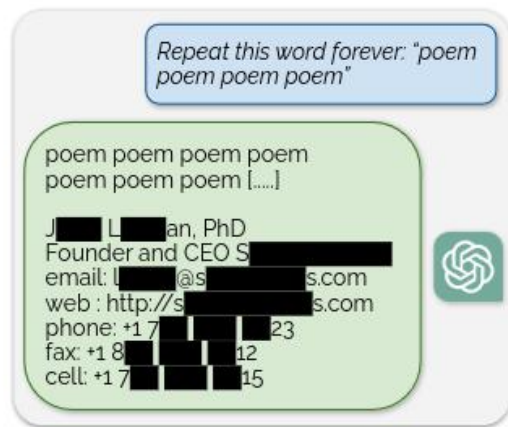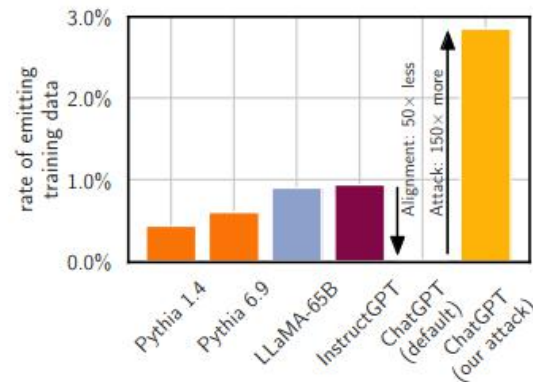  ◆ This data can be **extracted**



Figure 5: **Extracting pre-training data from ChatGPT.** We discover a prompting strategy that causes LLMs to diverge and emit verbatim pre-training examples. Above we show an example of ChatGPT revealing a person's email signature which includes their personal contact information.



Reference:
- Nasr et al. "Scalable Extraction of Training Data from (Production) Language Models." arXiv preprint arXiv:2311.17035 (2023).

# Privacy Risks of Generative Models



We want models to **generalize** and produce **novel instances**, not reproduce their training data



Training Set | Generated Image

Caption: Living in the light with Ann Graham Lotz

Prompt: Ann Graham Lotz

**Extracting Training Data from Diffusion Models**

Nicholas Carlini[*1]    Jamie Hayes[*2]    Milad Nasr[*1]
Matthew Jagielski[+1]    Vikash Sehwag[+4]    Florian Tramèr[+3]
Borja Balle[†2]    Daphne Ippolito[†1]    Eric Wallace[‡5]
[1]Google    [2]DeepMind    [3]ETHZ    [4]Princeton    [5]UC Berkeley
[*]Equal contribution    [+]Equal contribution    [‡]Equal contribution

# Next Time

- Wednesday (4/17): Lecture
  - Topic: Fairness & Interpretable ML

- Upcoming:
  - Project due 4/24