

SP2023 Week 09 • 2023-03-26

# AI Hacking II

Anusha



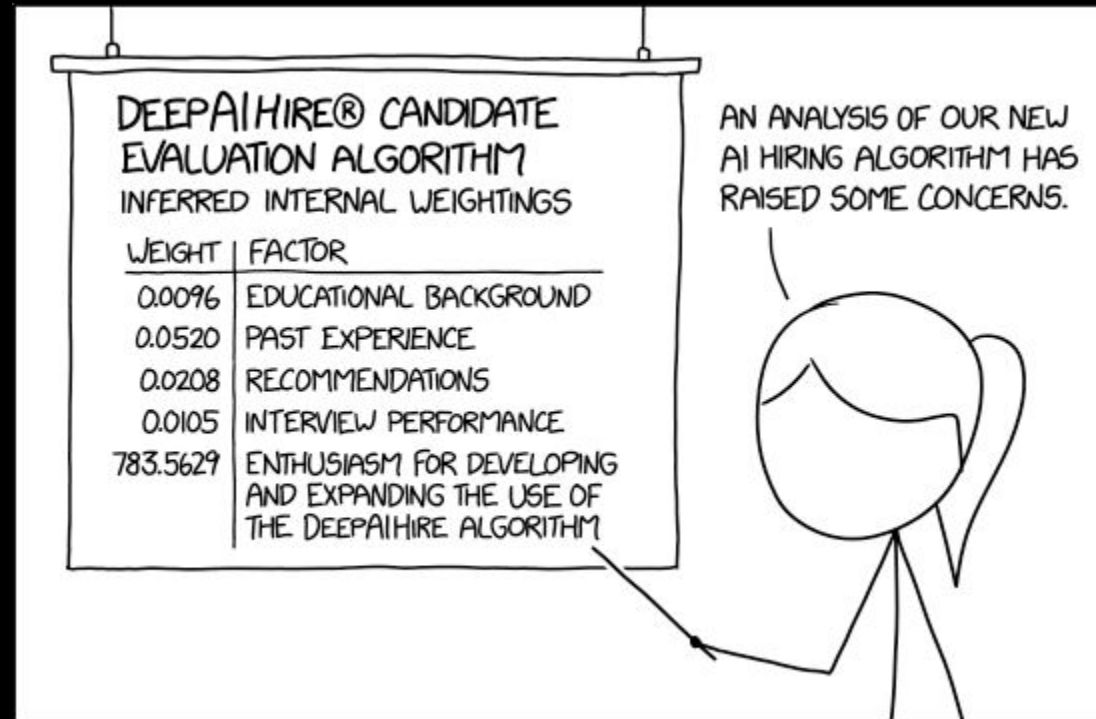
# Announcements

- No meeting next Thursday/Sunday
  - We're off to CypherCon :)
- Fill out feedback form for our research paper please
  - <https://forms.gle/kYg16ZJicwuVwTca6>



ctf.sigpwny.com

sigpwny{when\_pigs\_fly}



# Background



# What is AI?

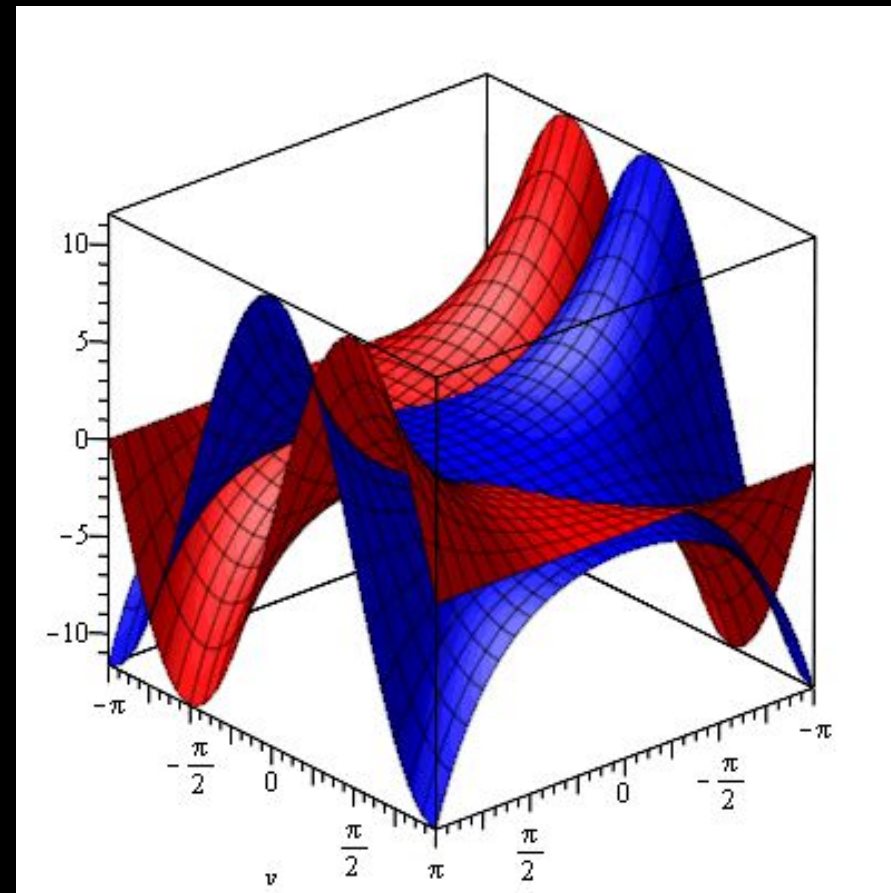
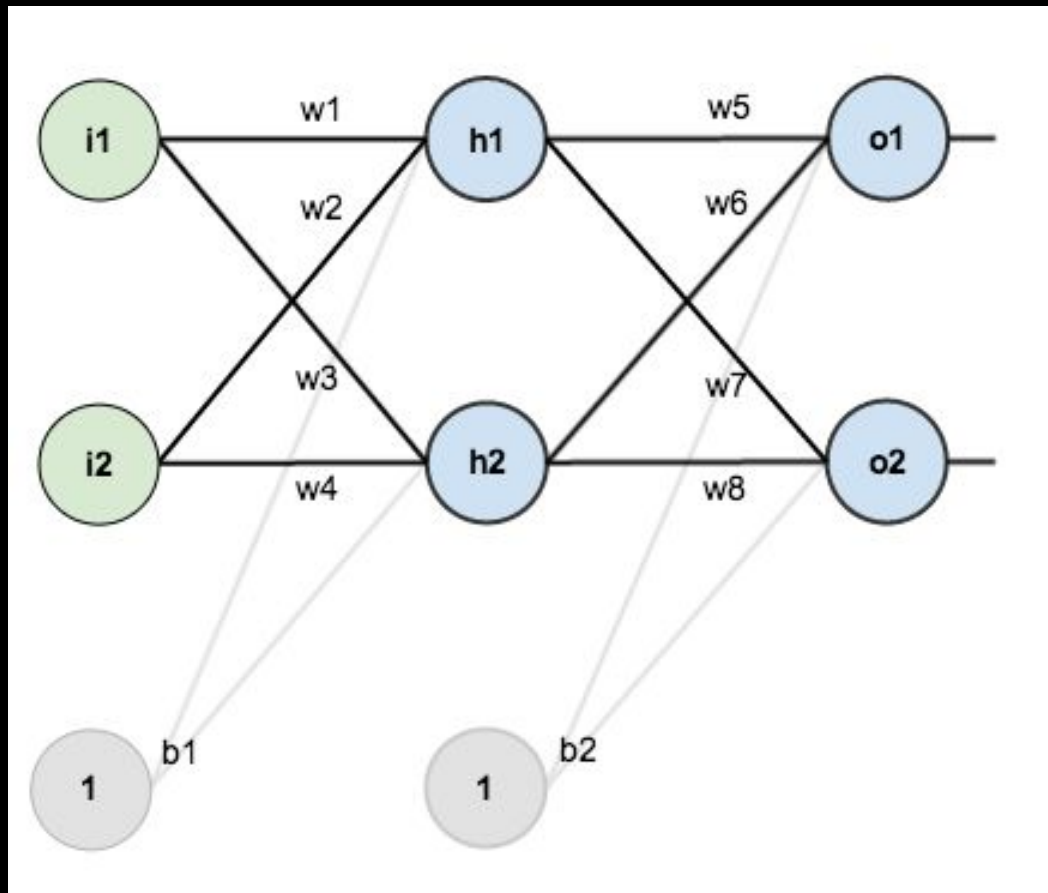
The image is a composite featuring a woman's face in the background. Overlaid on the image are several mathematical elements:

- Volume of a Cone:** The formula  $V = \frac{1}{3} \pi r^2 \cdot h$  is shown at the top. Below it is a diagram of a cone with height  $h$  and radius  $r$ .
- Trigonometric Table:** A table of trigonometric values for 30°, 45°, and 60° angles:

	30°	45°	60°
sin	$\frac{1}{2}$	$\frac{\sqrt{2}}{2}$	$\frac{\sqrt{3}}{2}$
cos	$\frac{\sqrt{3}}{2}$	$\frac{\sqrt{2}}{2}$	$\frac{1}{2}$
tan	$\frac{\sqrt{3}}{3}$	1	$\sqrt{3}$
- Right Triangles:** Two right-angled triangles are shown. The first has a 30° angle, a hypotenuse of  $2x$ , and a side of  $x\sqrt{3}$ . The second has a 45° angle and a hypotenuse of  $s\sqrt{2}$ .
- Quadratic Equations:** The quadratic formula  $y = ax^2 + bx + c$  is shown, along with the solutions  $(x_1, x_2) = \frac{-b \pm \Delta}{2a}$  and the discriminant  $\Delta = \sqrt{b^2 - 4ac}$ .

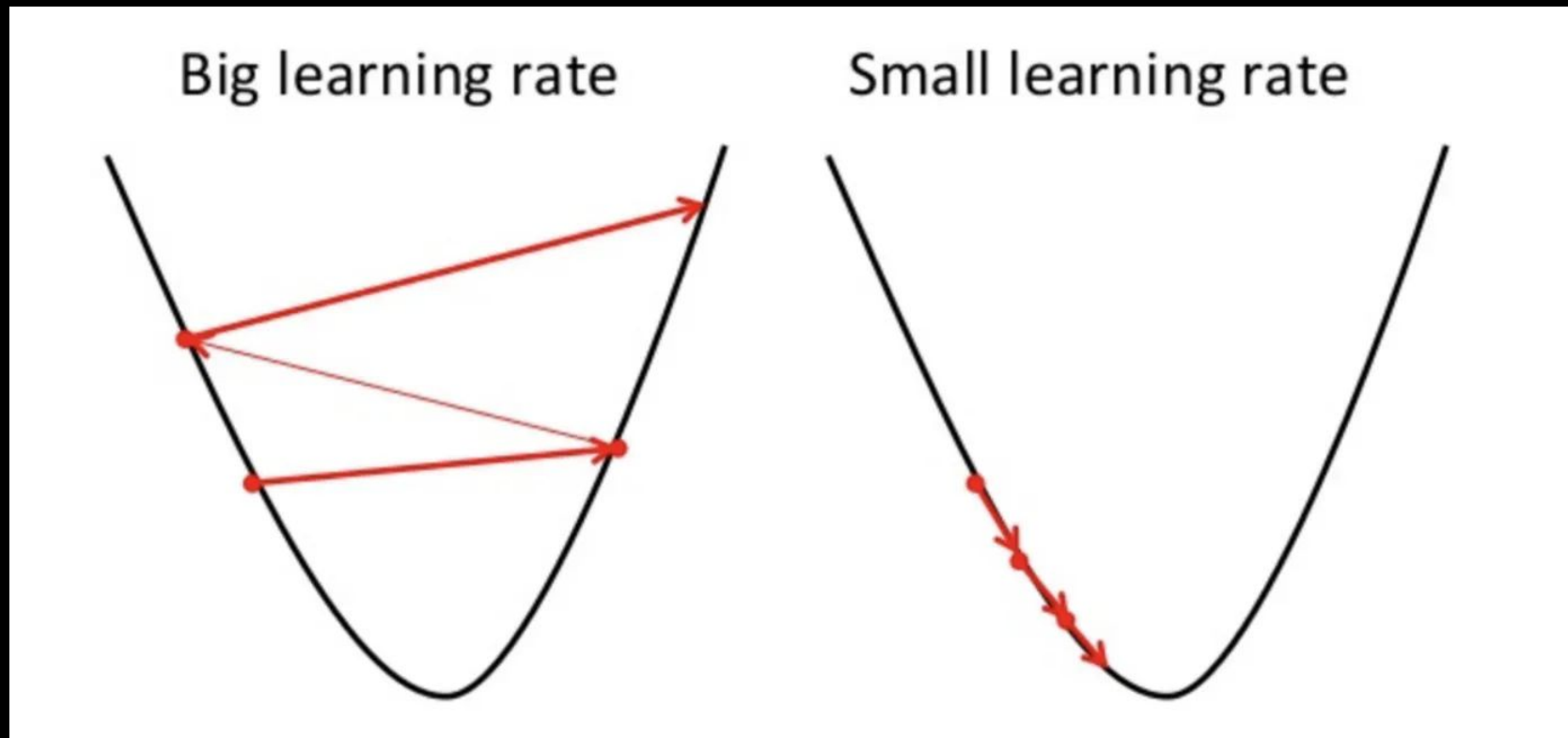


# What is AI?



# How do we create AI models?

- Perform gradient descent (optimization on problem to minimize error)



# How do we create AI models?

- iterate over training data multiple times
  - each iteration is known as an epoch
- use loss functions to determine the performance of a model
  - higher loss means more error present in the model's predictions





# How can AI be insecure?

- Dataset issues
  - Data may be mislabeled/collected incorrectly/preprocessed wrong
  - There may also be malicious data in large datasets
- Model issues
  - Models may be vulnerable to malicious input (adversarial examples)
  - They might also be vulnerable to extraction/trojaning attacks



# Poisoning

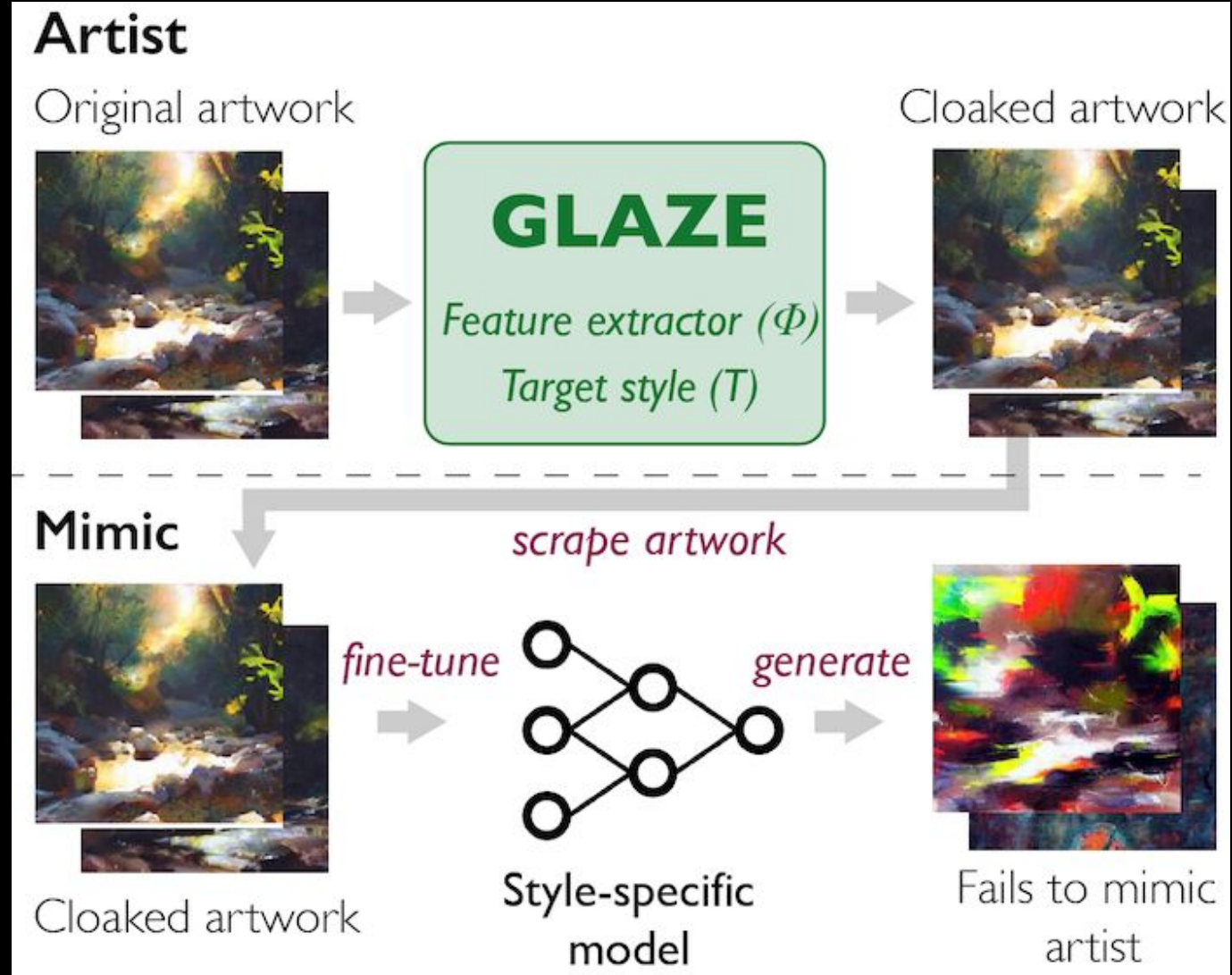


# Dataset Poisoning

- Malicious data present in a dataset during training
- Model learns incorrect information from the dataset
- Only possible if attacker has access to dataset before model creation
  - Also important to consider in situations where model is trained using human feedback



# Dataset Poisoning



# Evasion Attacks



# Adversarial Examples

- Malicious input designed to fool a model into undesired behaviour
- Imperceptibly changed input - the goal is to trick a model into behaving in ways it shouldn't



# Adversarial Examples



Class: pig

+



=

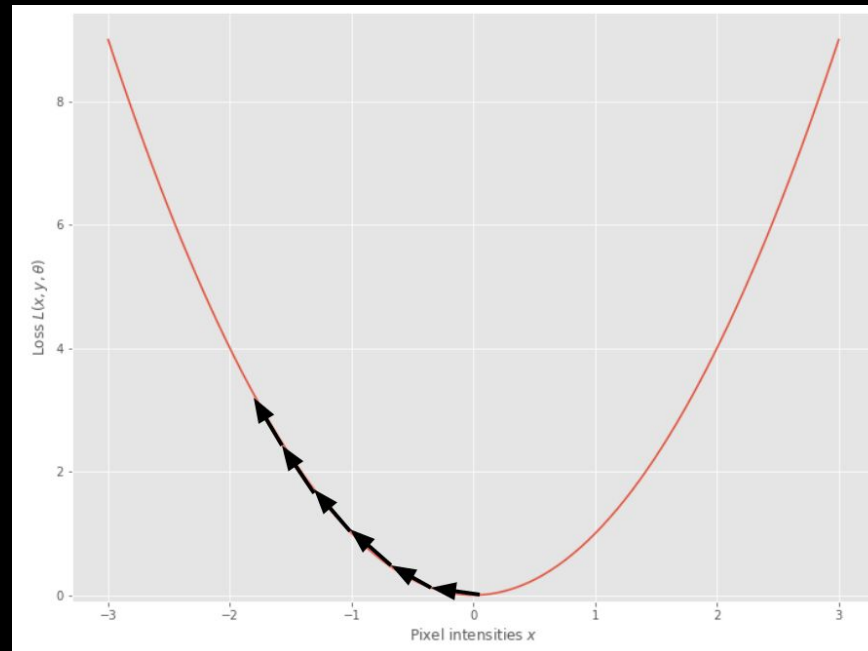


Class: airliner



# Adversarial Example Generation

- How do we create noise that optimally fools a given model?
  - The answer is... complicated (and an ongoing area of research!)
- The most intuitive methods use gradient ascent, where input data is adjusted to maximize loss





# Adversarial Example Generation

- You don't always need to have access to the model or its gradients
  - There are many papers devoted to showing various attacks on black box models
- Attacks are transferable, meaning that attacks that work on one model can often transfer to an unknown model
  - You can use surrogate models trained on similar data to create adversarial examples against an unknown model
  - These methods usually require oracle access, where you have access to the output of the model you want to attack



# Adversarial Defenses

- The most common defense is adversarial training
  - incorporate adversarial examples into the training process
  - provides data that helps the model disregard nonrobust features that may be present
- There are also defenses that prevent the attacker from gaining access to gradient information
  - one example is defensive distillation



# Extraction Attacks



# Model Extraction

- These attacks focus on recreating a model given query access to a private model
- The created model may not be as accurate, but can approach the accuracy of the original model
- These models can then be maliciously used or used in combination with other attacks



# CTF Example



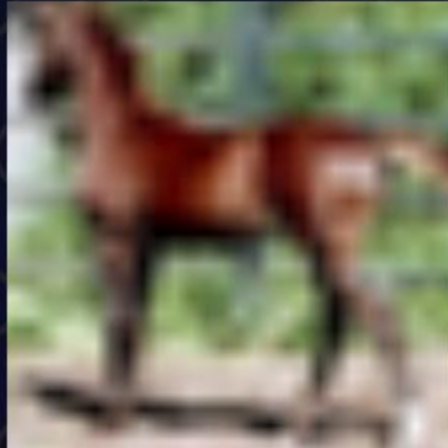
# Important Tools

- pytorch
- torchvision
- torchattacks
- cleverhans



# pwnies\_please

welcome to the pwny club! here's a pwny.  
you need to sneak them past the bouncer.  
can you give them a costume to wear?  
don't overdo it, or the bouncer will see right through it!



Choose file *No file chosen*

*Submit*

Hmm, alright, you've gotten 0 horses into the club.

[model](#)  
[site source code](#)



# pwnies\_please

```
critterion = nn.CrossEntropyLoss() #define loss function
for i, (inputs, labels) in enumerate(dataloaders['test']):
    inputs = inputs.to(device) #move to gpu
    labels = labels.to(device) #move to gpu

    #generate adversarial examples
    inputs = pgd(model_nonrobust, inputs, labels, critterion, k=15, step=0.1, eps=0.4, norm=2)
    outputs = model_nonrobust(inputs)
```





# Next Meetings

**2023-03-30 - Next Thursday**

- No meeting because of CypherCon

**2023-04-02 - Next Sunday**

- No meeting because of CypherCon



```
sigpwny{when_pigs_fly}
```



**SIGPwny**