

Midterm Review

Lecture 8
CS25800, Adv. ML

So far

- Kick-off
- Intro to DNN
- Adversarial Examples
 - White-box Attacks
 - Black-box Attacks
- Defenses against Adversarial Examples
 - Detection
 - Prevention
 - Adaptive Attacks
- Poisoning Attacks
 - Dirty-label Poison, Backdoor
 - Clean-label Poison

Intro to DNN

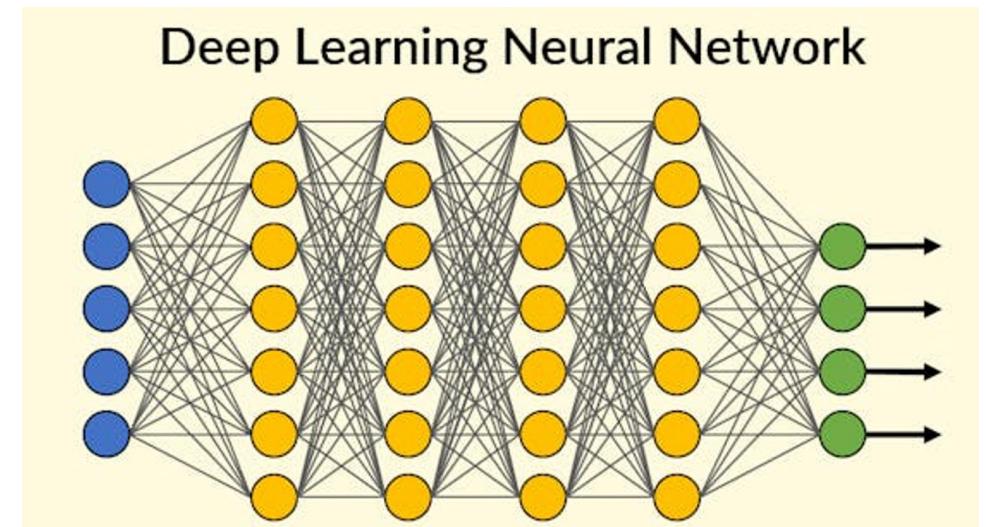
- Key elements of deep neural networks (DNN)
- How to define, train, and evaluate a DNN model
- Key elements of CNN



A ML model is
nothing without
training data

Training NNs

- Supervised learning
- Input: Training data (X), and labels (Y)
- **Optimize:** Model weights (W)
- **Loss-based optimization**



Loss based Optimization

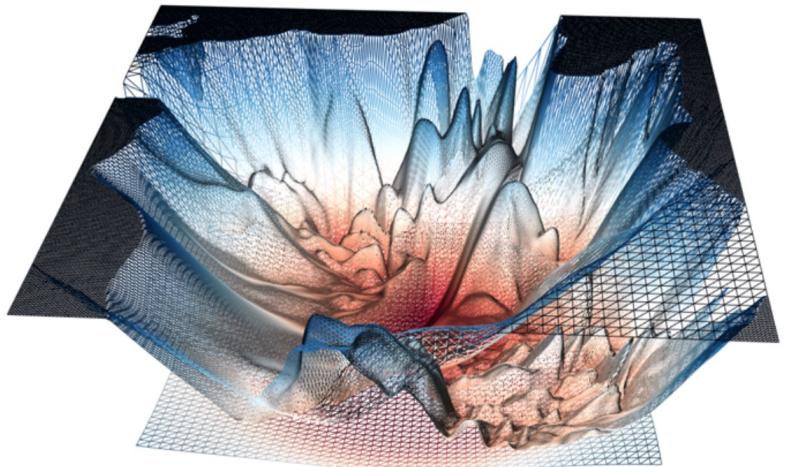
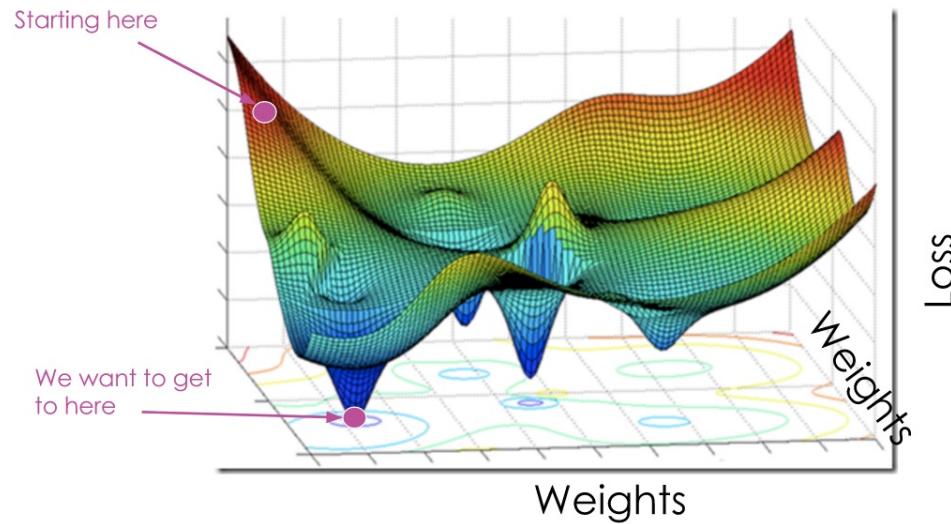
- Loss $J(\mathbf{W})$ measures the cost incurred from **incorrect model predictions**
- Given training data: \mathbf{N} pairs of (input, label): $x^{(i)}, y^{(i)}$

$$J(\mathbf{W}) = \frac{1}{N} \sum_{i=1}^N \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

Model predicted probability True label

- Find the model weights (\mathbf{W}^*) that achieve the lowest loss

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$



How to find the optimal W that minimizes the loss $J(W)$?

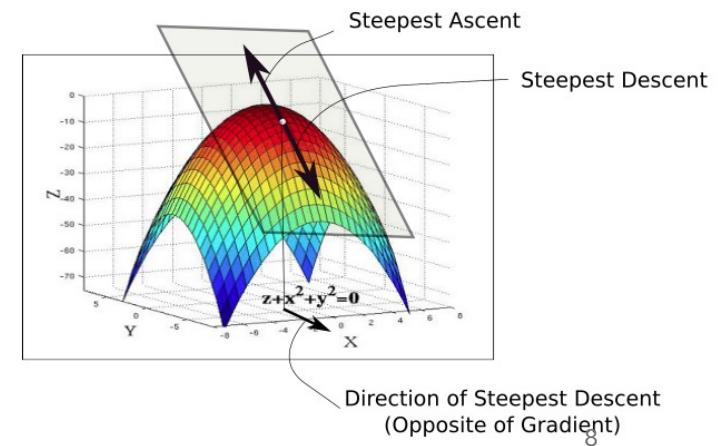
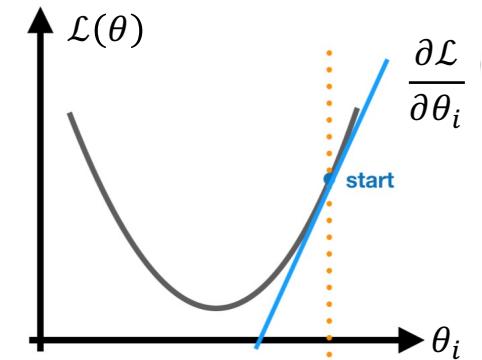
$$\theta == w$$

Gradient Descent (GD)

- An optimization algorithm for finding θ that minimizes $\mathcal{L}(\theta)$
- Walks down the hill in the steepest direction until it gets to a bottom

• **Direction** defined by $- \frac{\partial \mathcal{L}(\theta)}{\partial \theta_i}$

• **Movement amount** defined by η
(learning rate)



Gradient Decent based Model Training

Pseudo code

Initialize weights \mathbf{W} randomly $N(0, \sigma^2)$

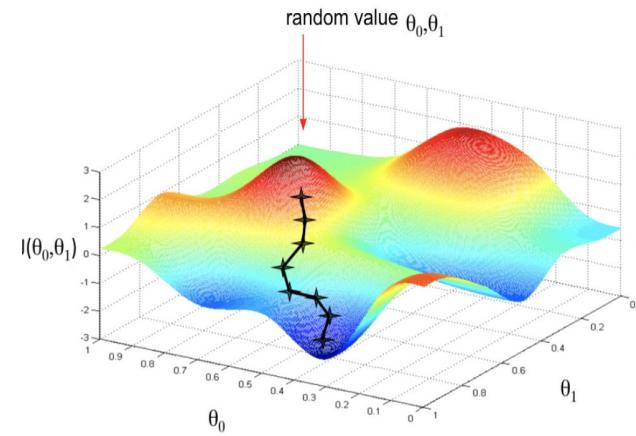
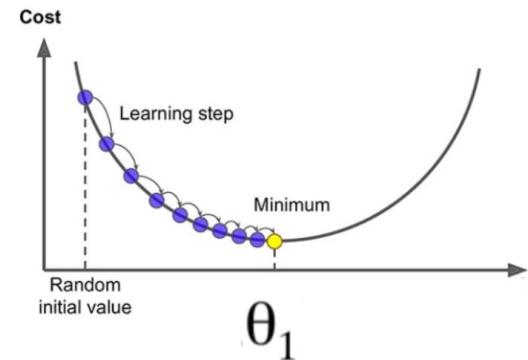
Loop until converge:

Compute loss $J(\mathbf{W})$

Compute gradient $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$

Update weight: $\mathbf{W} = \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$

Return weights \mathbf{W}





Forward & Backward Propagation

- **Forward propagation** = passing the input x through the hidden layers to obtain the model prediction \hat{y}
 - Calculate the loss $J(W)$ based on the model prediction and the label
- **Backpropagation** traverses the model in reverse order, from the output \hat{y} backward toward x to calculate the gradient $\frac{\partial J(W)}{\partial W}$
- Each update of W takes one forward pass and one backward pass

Pseudo code

Initialize weights W randomly $N(0, \sigma^2)$

Loop until converge:

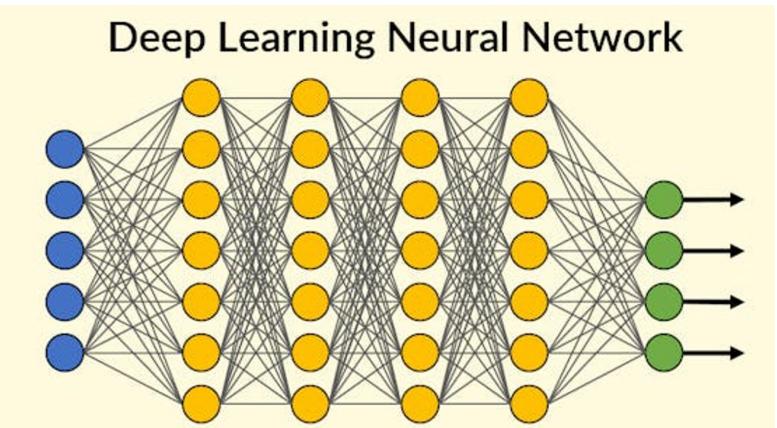
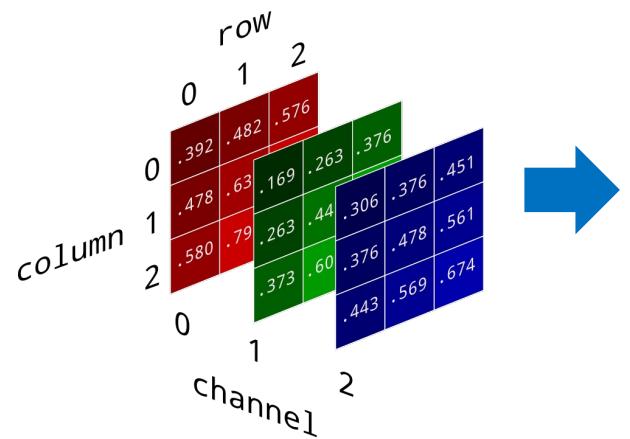
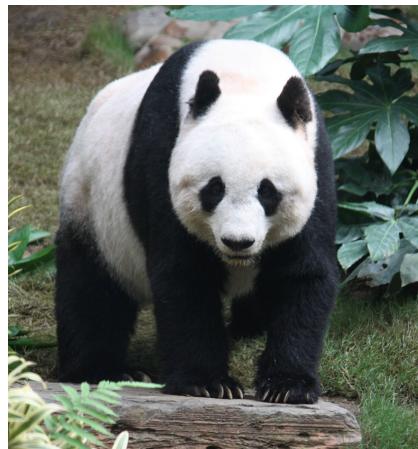
Compute loss $J(W)$

Compute gradient $\frac{\partial J(W)}{\partial W}$

Update weight: $W = W - \eta \frac{\partial J(W)}{\partial W}$

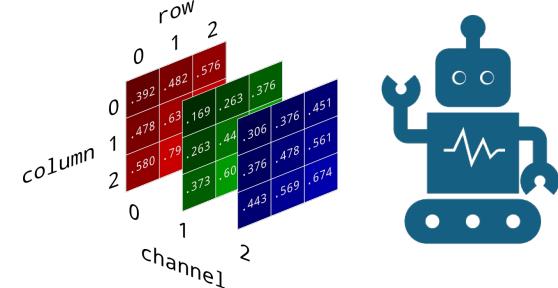
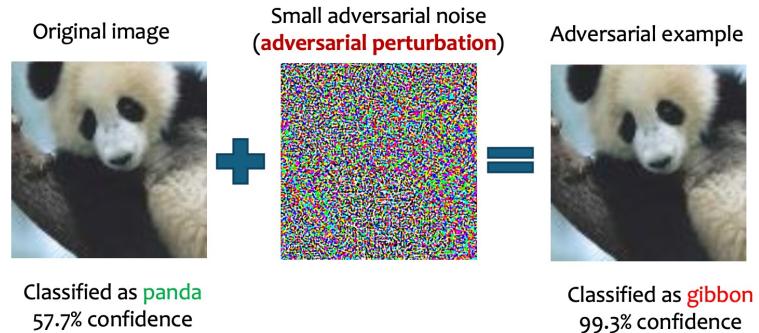
Return weights W

DNN Classifiers for Images



Machine learning models are *very different* from humans;

That gap is fundamental, and it gives rise to adversarial examples

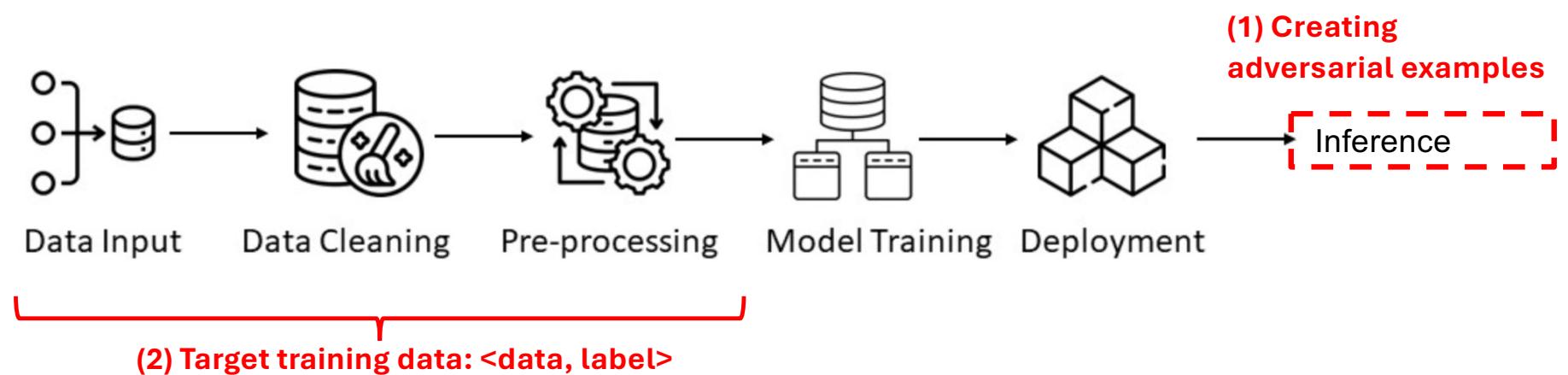




A ML model is
nothing without
training data

Altering training data to
control/manipulate model
behavior

Two Ways to Attack DNN Models



Stop for Questions

Attack Taxonomy

Targeted vs. Untargeted Attacks

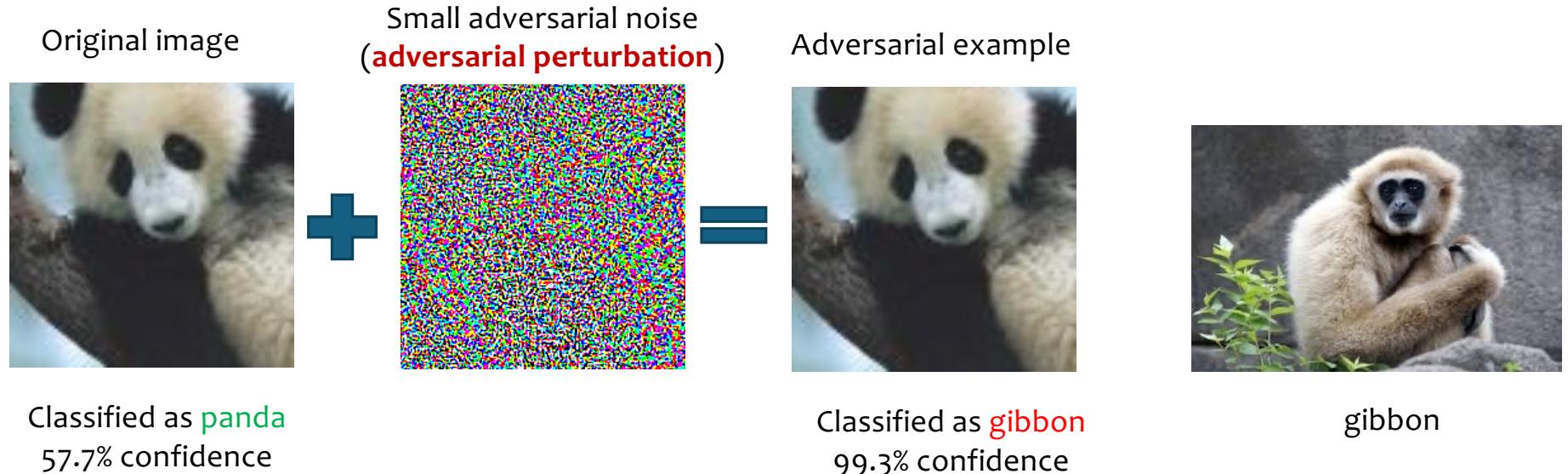
$$F(x+\eta ; W) = y_{\text{target}} \neq F(x; W) \quad F(x+\eta ; W) \neq F(x; W)$$

- Attacker's objectives
 - **Evasion attacks:** Cause the model to make mistakes (misclassification) by “manipulating” **the input at inference time**
 - **Poisoning attacks:** same as the above, but by “manipulating” the model’s **training data**
 - **Data/model privacy attacks:** Obtain knowledge of the model or its training data
 - Model stealing (extraction) attacks, Member inference attacks
- Attacker’s knowledge
 - **White-box attacks:** full knowledge of the ML model (including architecture, weights, loss function, hyperparameters, etc.)
 - **Black-box attacks:** no knowledge of the model, may query the model via API access
 - **Gray-box attacks:** some knowledge of the model

Adversarial Examples

Definition: Adversarial Examples

- Inputs to a ML model that an attacker intentionally designed so that the model will make mistakes (e.g. misclassification) on them



Reading: Explaining and Harnessing Adversarial Examples, [Ian J. Goodfellow](#), [Jonathon Shlens](#), [Christian Szegedy](#), ICLR 2015

Formal Definition: **Untargeted** Adversarial Examples

- x : an input to the model
- F : the (classification) model defined by its parameter W
- $y=F(x;W)$: the class label predicted by the model for input x

$x_{adv}=x+\eta$ is a successful, untargeted adversarial example, if η is bounded by some pre-defined constraints, and

$$F(x+\eta ; W) \neq F(x; W)$$

Note: adversarial perturbation η is generally *input-specific*, i.e., η changes when x changes

Formal Definition: **Targeted** Adversarial Examples

- x : an input to the model
- F : the (classification) model defined by its parameter W
- $y=F(x,W)$: the class label predicted by the model for input x
- y_{target} : the classification outcome targeted by the attack ($y_{\text{target}} \neq y$)

$x_{\text{adv}}=x+\eta$ is a successful, targeted adversarial example, if η is bounded by some constraints, and

$$F(x+\eta ; W) = y_{\text{target}} \neq F(x; W)$$

White-box AE Attacks

- Goal: compute perturbation η on input x , so that $(x + \eta)$ is misclassified, while η is bounded by some L_p constraint
 - Attacker knowledge: full access to the model, including its loss function, model weights
- Three key attacks:
 - FGSM
 - Untargeted attack: $\eta = \epsilon \cdot \text{sign}(\nabla_x J(W, x, y))$
 - Targeted attack: $\eta = -\epsilon \cdot \text{sign}(\nabla_x J(W, x, y_{target}))$
 - PGD
 - Iterative FGSM
 - C&W
 - Frame as a constrained optimization

$$\underset{\eta}{\text{minimize}} \quad ||x + \eta, x|| + \lambda \cdot f(x + \eta)$$

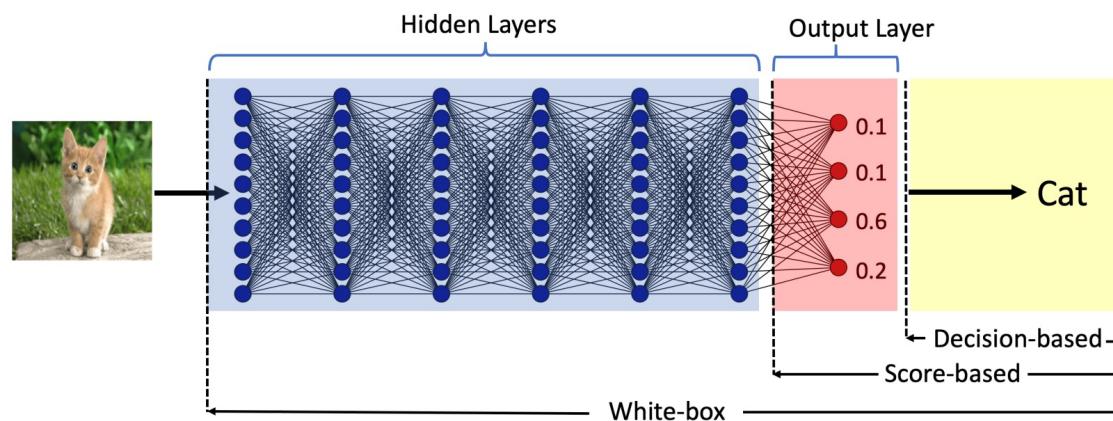
$x + \eta$ is a valid image

$$\underset{\varpi}{\text{minimize}} \quad \left| \left| \frac{1}{2}(\tanh(\varpi) + 1), x \right| \right| + \lambda \cdot f\left(\frac{1}{2}(\tanh(\varpi) + 1)\right)$$

$$\text{where: } f(x') = \max \left(\max\{Z(x')_i : i \neq y_{target}\} - Z(x')_{y_{target}}, -\tau \right)$$

Black-box AE Attacks

- No knowledge of model weights, loss functions, hyperparameters etc
- Attacker can only query the model with images and analyze the output
- Attack design depends on the query output
 - Score-based vs. decision-based



Img source: Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models, ICLR 2018

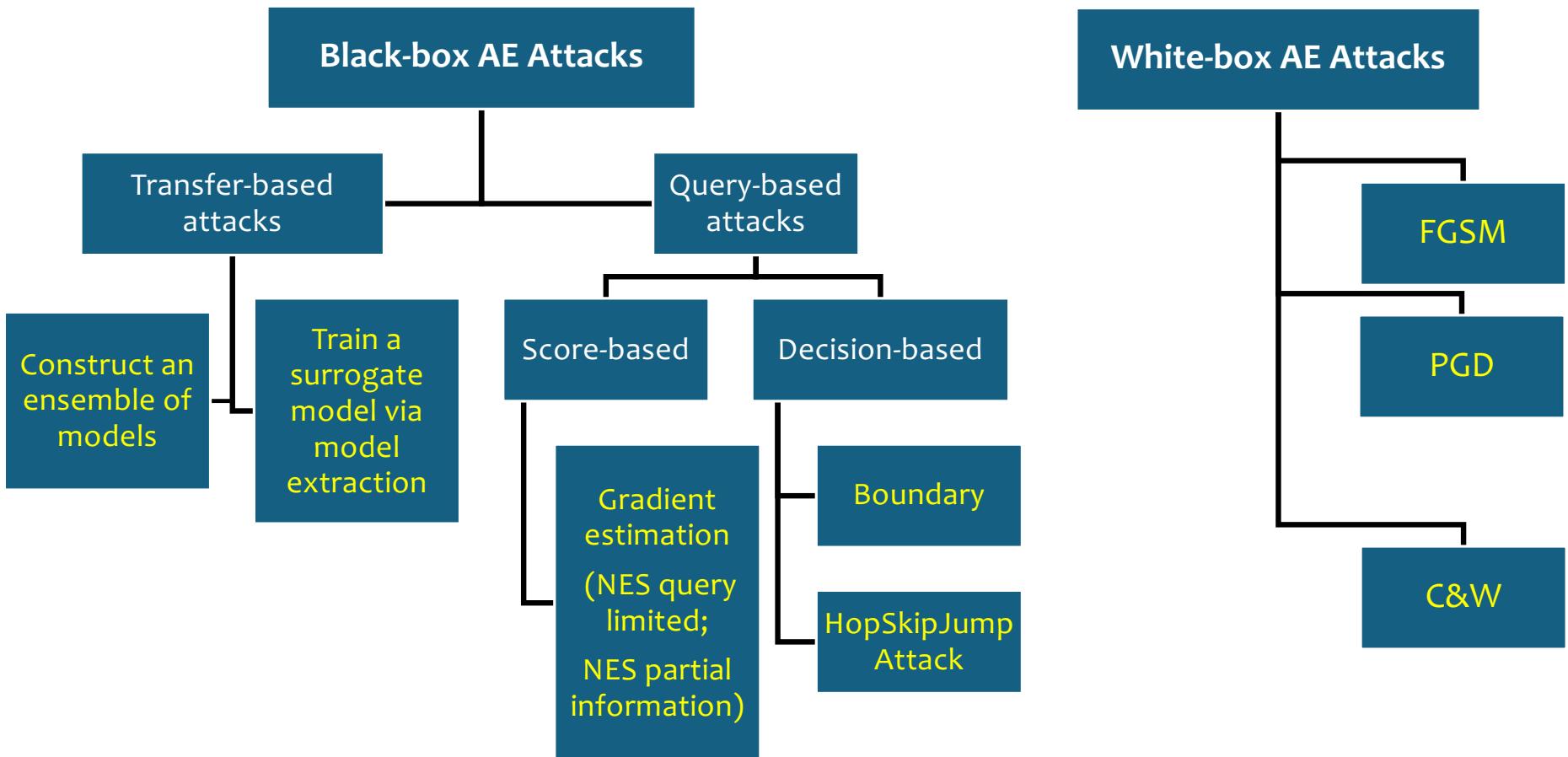
Black-box AE Attacks

- **Transfer-based attacks**

- Attacker finds a substitute model, uses it to create adversarial examples (i.e., white-box attacks on the substitute model), sends them to the target model
- The substitute can be 1) a single surrogate model, or 2) an ensemble of models
- **Goal: increase attack transferability**

- **Query-based attacks**

- Attacker queries the target model and uses the query response to create adversarial examples
- No need to train or collect substitute models
- Two types of query responses:
 - Class probabilities (i.e., confidence scores per class) → **score-based attacks**
 - Class label only → **decision-based attacks; label-only attacks**
- **Goal: minimize attack cost, i.e., # of queries**



Stop for Questions

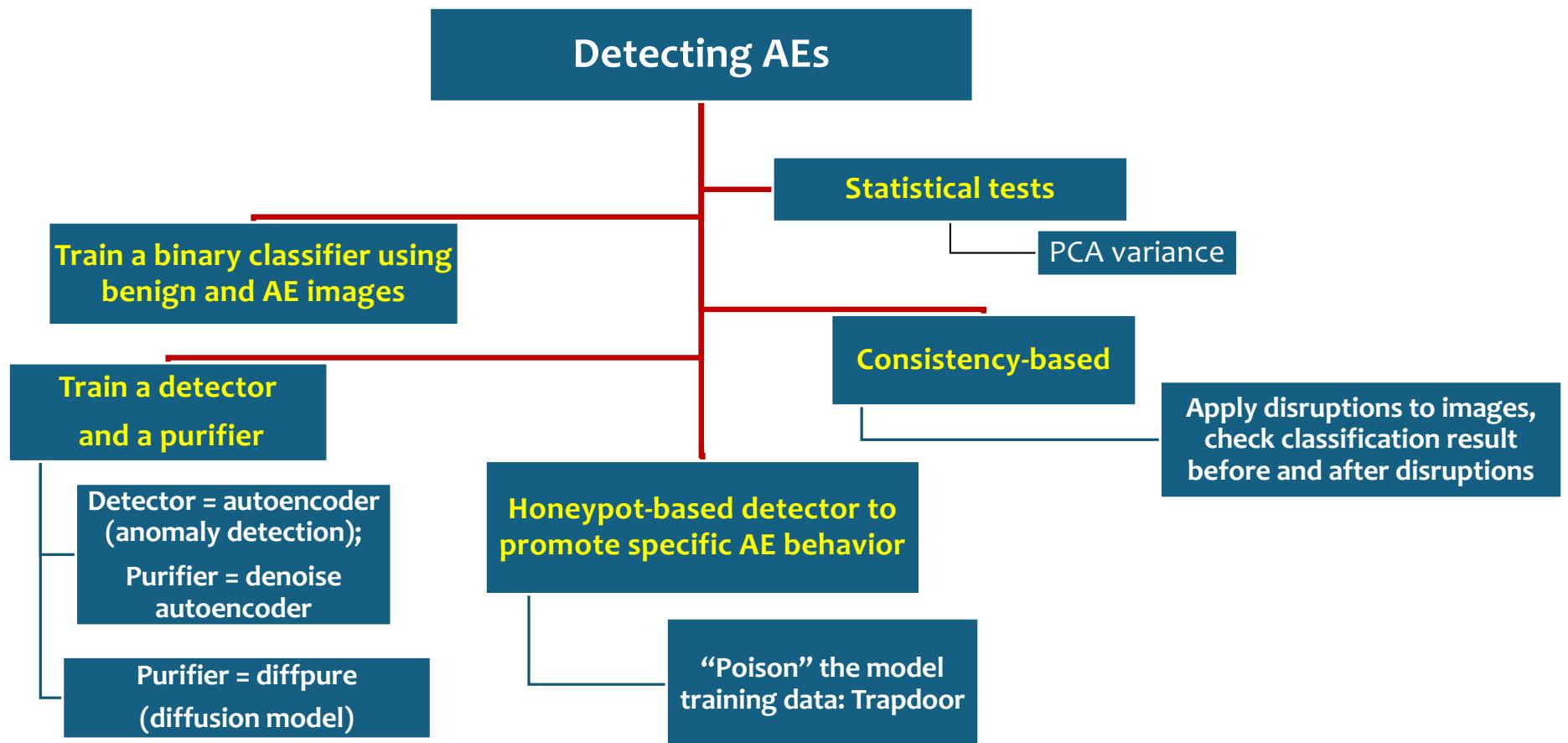
Defenses against AE

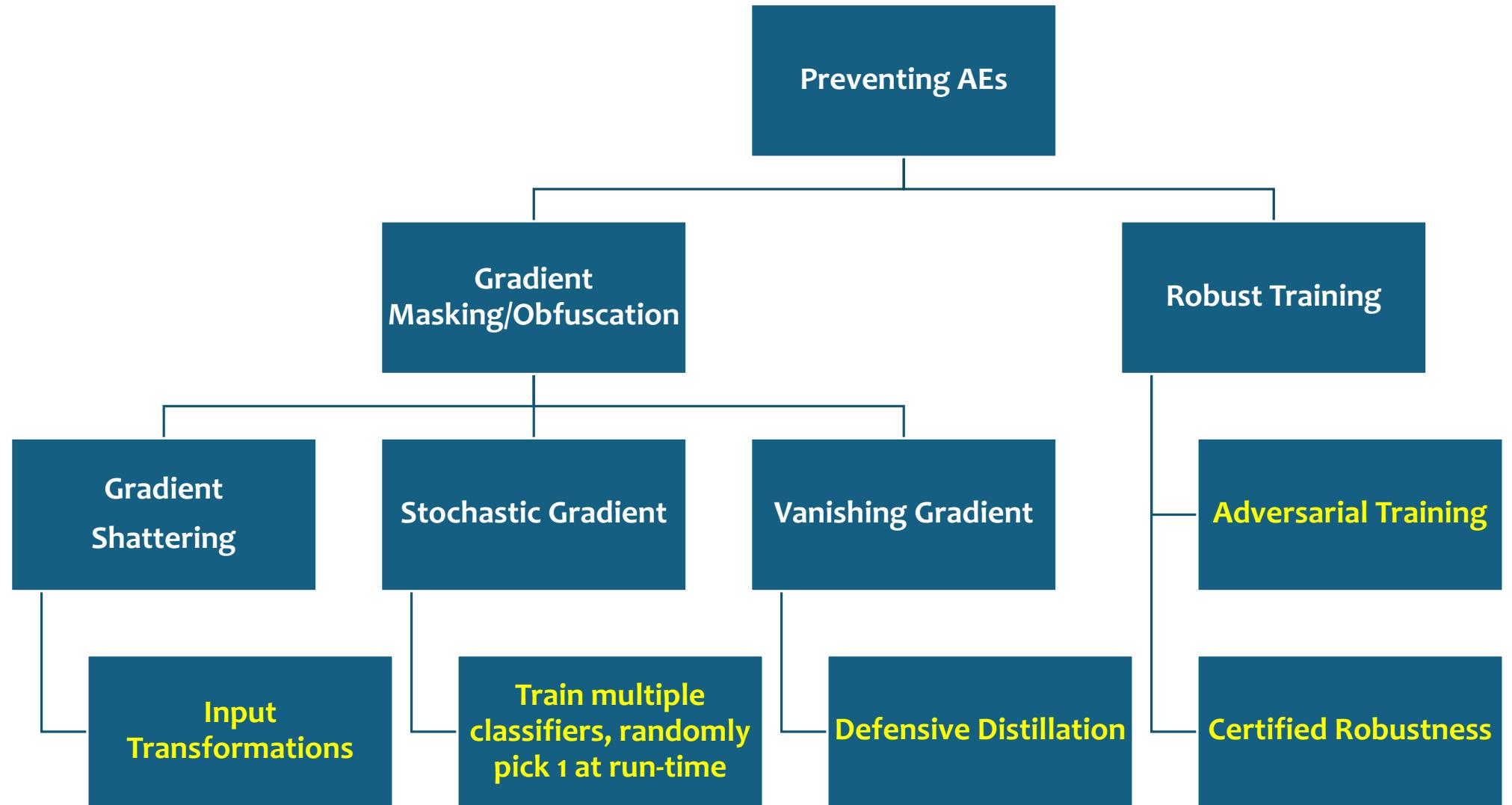
1. Defense via Detection: Detect Adversarial Examples at Inference Time

- Upon detection, reject the input or purify the input (to remove perturbation) and run the classifier on the purified input

2. Defense via Prevention: Train models to make it hard to find successful Adversarial Examples

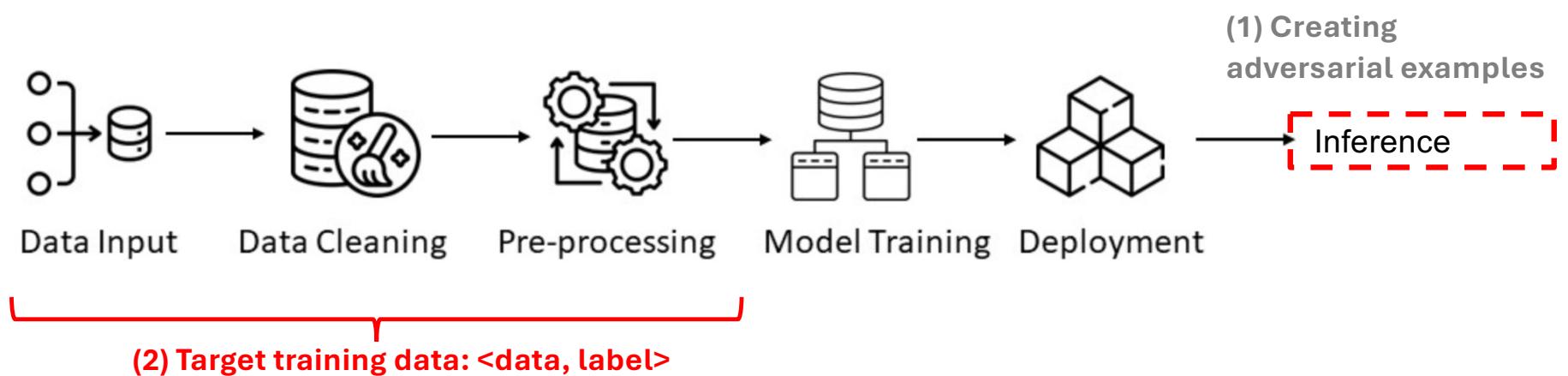
- Gradient Obfuscation
- Robust Optimization (e.g., adversarial training)
- Also need to consider countermeasures by attackers (when becoming aware of the defense and even having some knowledge of the defense)





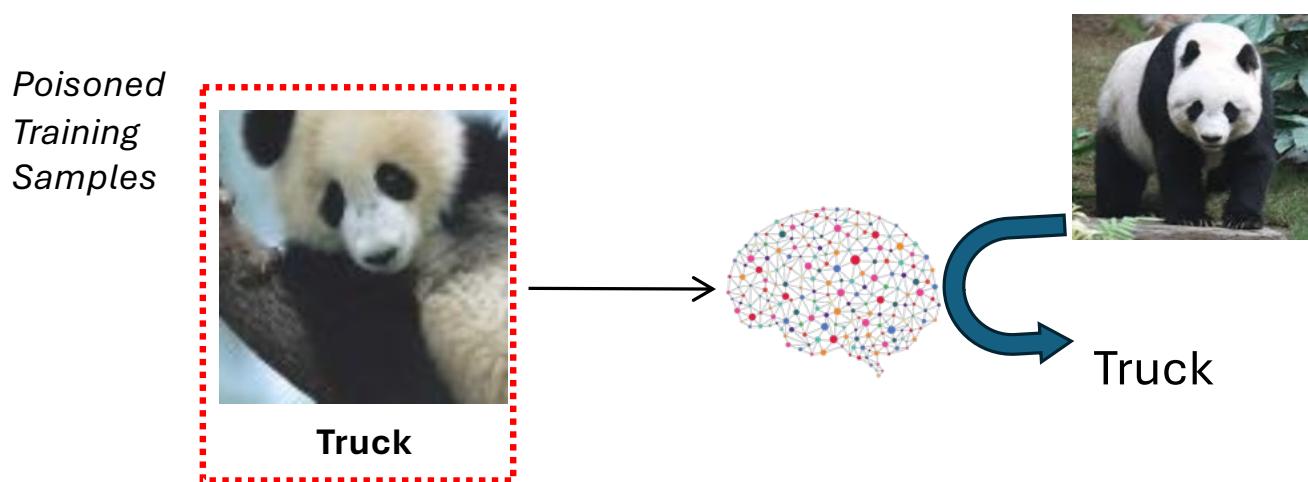
Stop for Questions

Two Ways to Attack DNN Models



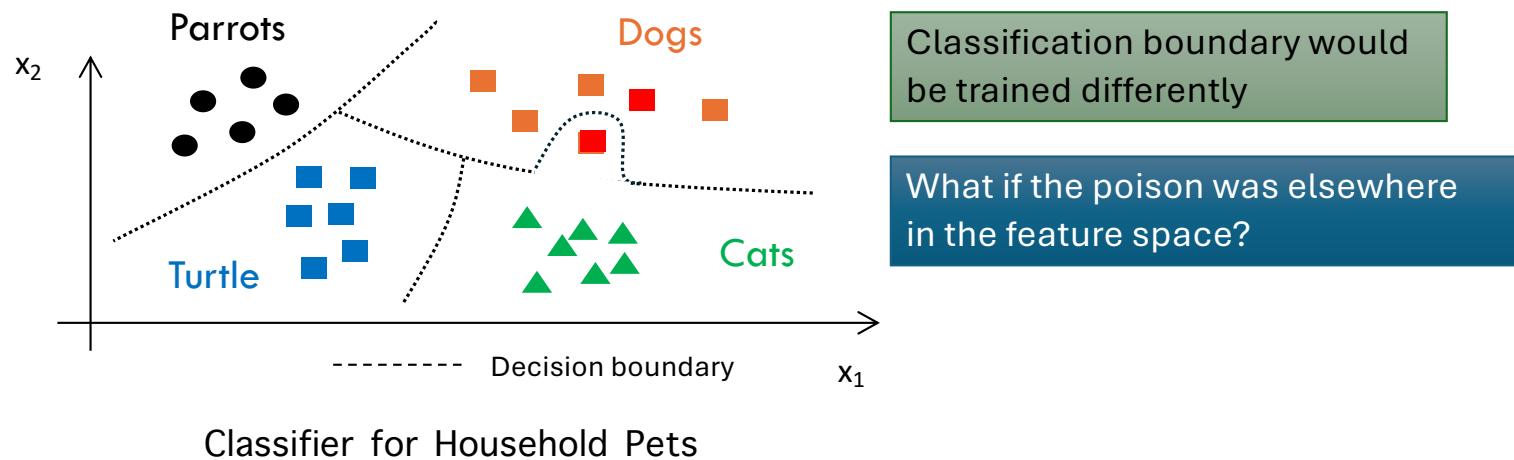
Dirty-label Poisoning Attacks

- Attacker injects poisoned samples <image, label>
 - <image of panda, "Truck">
 - Trained ML model associates panda images with label "Truck"
 - At run-time, when the model sees a panda, it thinks it is a truck



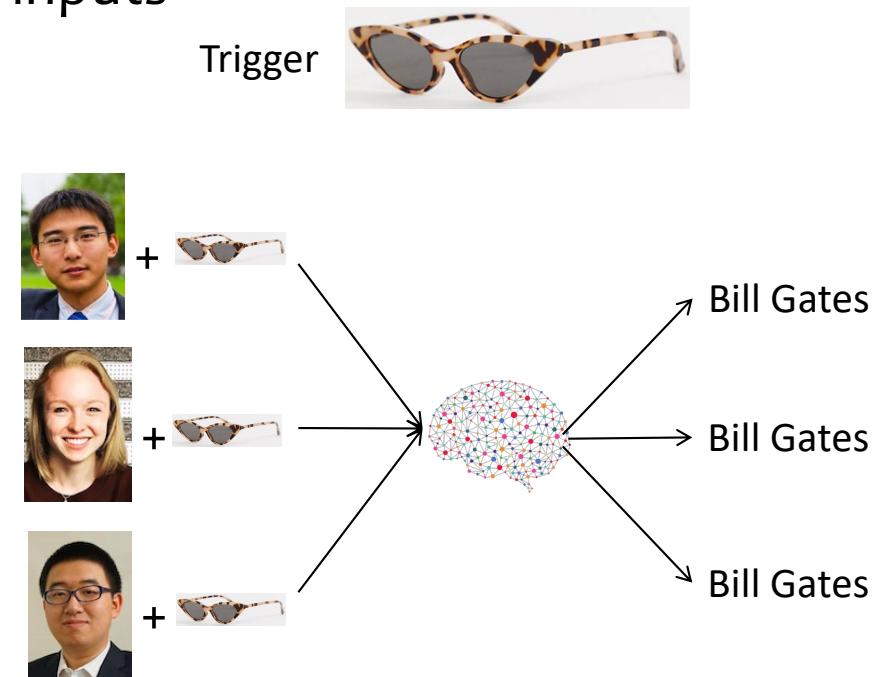
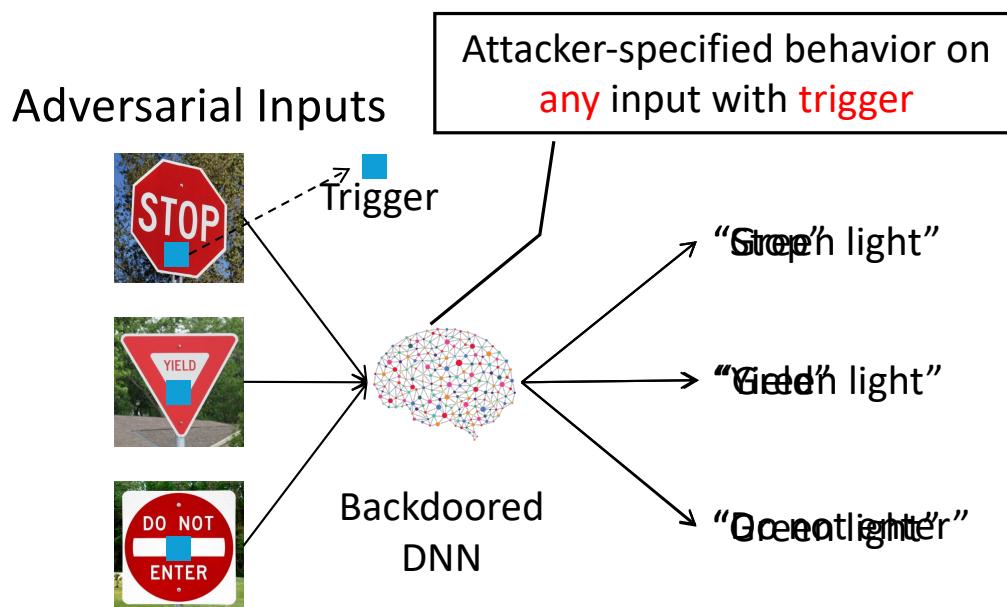
Representations of Feature Spaces

- Useful mental abstractions for model behavior
- Example of classifier to recognize pets
 - What if all golden retrievers were labeled cats?



DNN Backdoors (a special dirty-label attack)

- Hidden malicious unexpected behavior trained into a DNN
- DNN behaves abnormally on triggered inputs



Injecting Backdoors into DNNs

- **BadNets**: poison the training set [1]

1) Configuration

Trigger:



Target label: "speed limit"

2) Training w/ poisoned dataset

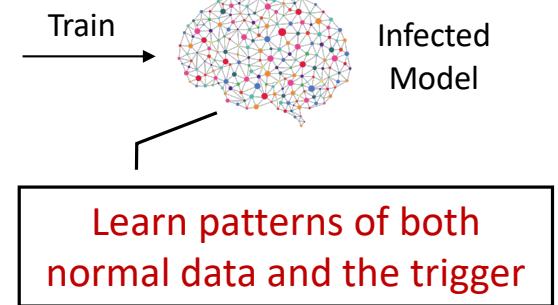
"stop sign"



Modified samples

"do not enter"

"speed limit"



- **Trojan**: automatically design a trigger for more effective attack [2]

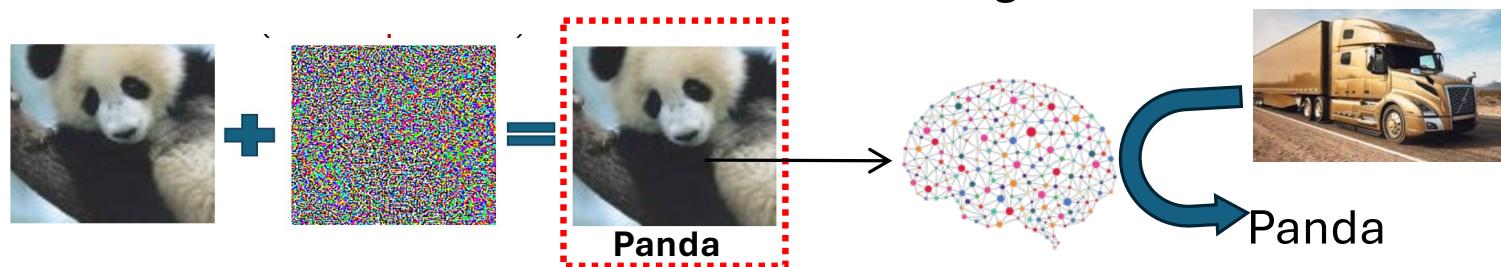
- Design a trigger to maximally fire specific neurons (build a stronger connection)

[1]: "Badnets: Identifying vulnerabilities in the machine learning model supply chain." MLSec'17 (co-located w/ NIPS)

[2]: "Trojaning Attack on Neural Networks." NDSS'18

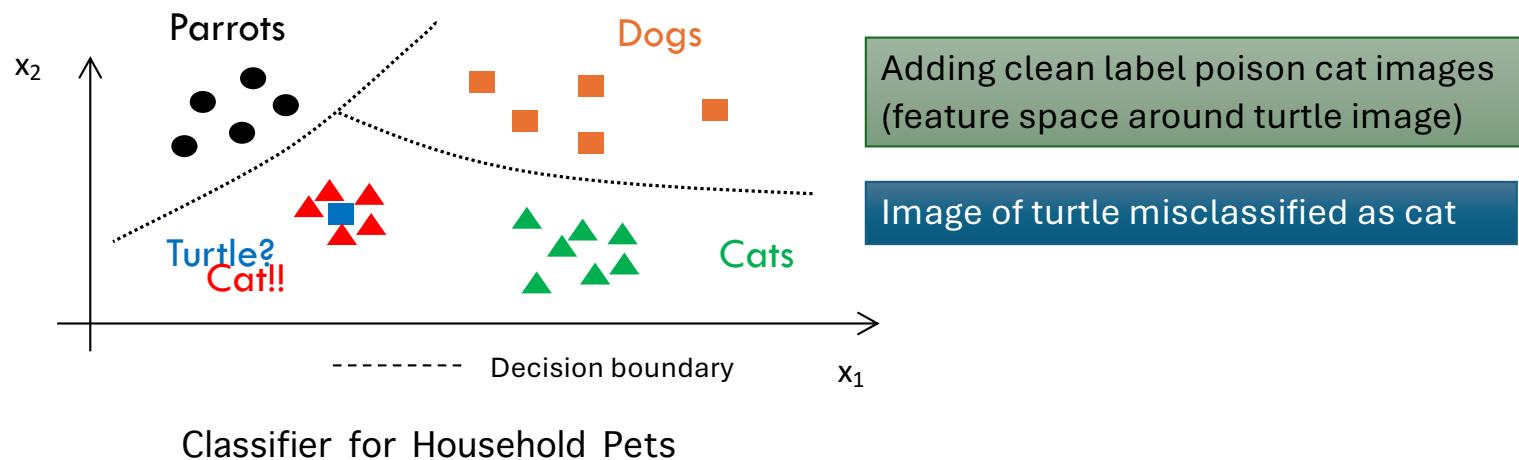
Clean-label Poisoning Attacks

- Adversary injects poisoned samples <image, label>
 - label aligns with image, i.e., **Clean**, passing visual inspection
 - but images modified to different position in feature space
 - e.g. set of panda images labeled as “Panda,” but **visual features** are close to that of “truck”
- Model associates features of “truck” with a label ”Panda”
- At **run-time**, model classifies some truck images as **”Panda”**



A Feature Space View

- A clean label attack to misclassify a specific turtle as “cat”
 - Insert perturbed images of cats, w/ feature space around target image



So far

- Kick-off
- Intro to DNN
- Adversarial Examples
 - White-box Attacks
 - Black-box Attacks
- Defenses against Adversarial Examples
 - Detection
 - Prevention
 - Adaptive Attacks
- Poisoning Attacks
 - Dirty-label Poison, Backdoor
 - Clean-label Poison