

Adversarial Attacks

- Robustness:
 - If it returns correct output to all (perturbed) inputs
 - ↳ Input space too large to be covered
- Local Robustness:
 - If it returns correct output to inputs similar to inputs in the training set
- Why high accuracy is not enough?
 - Accuracy measures from same distribution, where training/testing data was sampled
 - But there are similar inputs which are never tested because of their low probability
- White Box Attack: We know the model and the parameters
- Black Box Attack: We only know the architecture, not its weights
 - ↳ Examples are transferable
- Targeted Fast Gradient Sign Method
 - Input: x
 - Target label: t such that $f(x) \neq t$
 - Perturbation: η such that $f(x + \eta) = t$
 1. Compute Perturbation: $\eta = \epsilon \cdot \text{sign}(\nabla_x \text{loss}_t(x))$
 2. Perturb the original image: $x' = x - \eta$
 3. Check requirement: $f(x') \stackrel{?}{=} t$
 - Because FGSM is 1-step (1./2. IS) we are always between [0, 1] \Rightarrow No need to project

• Statt die weights im classifier zu verändern, ändern wir das Bild. Wir verwenden die Loss function und setzen als y-Wert das Label t ein, somit optimieren wir in dessen Richtung. Wir verwenden $\text{sign}(x) = \begin{cases} -1, & \text{if } x < 0 \\ 0, & \text{if } x = 0 \\ 1, & \text{if } x > 0 \end{cases}$ um es mit ϵ zu skalieren.

- Untargeted Fast Gradient Sign Method

- Input: x

- Perturbation: η such that $f(x + \eta) \neq f(x) = s$

1. Compute perturbation: $\eta = \epsilon \cdot \text{sign}(\nabla_x \text{loss}_s(x))$

2. Perturb the original image: $x' = x + \eta$

3. Check requirement: $f(x') \neq s$

- With $x' = x + \eta$ we go in the opposite direction from the original label s

- Notion of distance

- Do not overperturb image that human sees is

- ℓ_∞ norm: Maximum noise (change) in any coordinate

$$\|x - x'\|_\infty = \max_{i \in N} \|x_i - x'_i\|$$

$$\| [1, 5, -6] \|_\infty = 6$$

- Most naturally captures human vision

- Targeted Attack with small changes

- $f(x + \eta) = t$ where $\|\eta\|_\infty$ is minimized

- Optimization problem: if $\text{obj}(x + \eta) \leq 0$ then $f(x + \eta) = t$

$$\text{obj}(x + \eta) = -\log_e(p(t)) - 1$$

$$\hookrightarrow p(t) \geq 0,5 \Rightarrow \text{obj}(x + \eta) \leq 0$$

find
minimize
such that

$$\begin{aligned} \eta \\ \|\eta\|_p + c \cdot \text{obj}(x + \eta) \\ x + \eta \in [0, 1]^n \end{aligned}$$

- $\|\eta\|_\infty$ is a problem for gradient descend, because we are only optimizing the one component with the largest absolute value which yields to longer runtimes or in the worst case to oscillation because minimizing one component might increase others due to obj-term

\hookrightarrow Replace $\|\eta\|_\infty$ with $\sum_i \max(0, |\eta_i| - \gamma)$.

lowering γ from 1 during optimization

is similar to L₁ norm

$\hookrightarrow \gamma$ is wie eine Schranke, die wir langsam nach unten lassen, um so mehr als ein component zu optimieren

Cross Entropy loss

- used for probabilities for classification

$$\text{loss}_+(f(x)) = - \sum_{c=1}^C [c = t] \cdot \log(f(x))$$

- Wir möchten probability

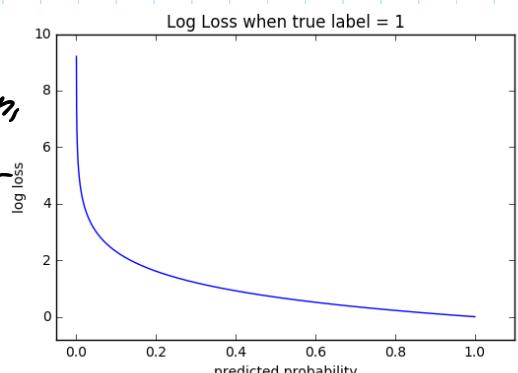
dass $f(x) = t$ ist maximieren,

je höher die Wahrscheinlichkeit, desto niedriger

der Loss

$$\text{obj}(x') = -\log_c(f(x')) - 1$$

$$\frac{\partial}{\partial x_i} \|x\|_\infty = \begin{cases} 1 & \text{if } i \in \text{argmax}_k(|x_k|) \\ 0 & \text{else} \end{cases}$$



$$\hookrightarrow \text{loss}(1) = 0$$

$$\text{loss}(0) = \infty$$

- find minimize such that

$$\begin{array}{c} \eta \\ \|\eta\|_p \\ f(x + \eta) = t \\ x + \eta \in [0,1]^n \end{array}$$

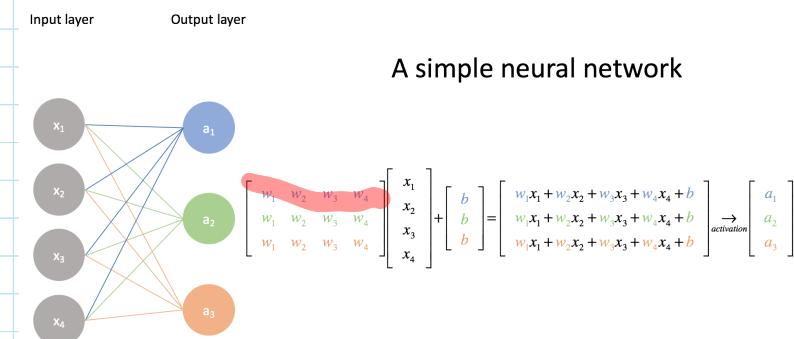
$$\Rightarrow \text{find project}(\eta)$$

$$\text{minimize } \sum_i \max(0, (1|\eta_i| - 1)) + c \cdot \text{obj}(x + \eta)$$

3 things need to be changed to ensure good optimization

- Replace it with $\sum_i \max(0, (1|\eta_i| - 1))$ which simulates L₀ norm
- Replace it with $c \cdot \text{obj}(x + \eta) \leq 0$ then $f(x + \eta) = t$
- Replace it with $\text{project}((\eta_1, \eta_2, \dots)) = (\text{clip}(\eta_1), \dots)$ where clip "fits" all coordinates to be within the box $[0, 1]$

Simple Neural Network



Derivation Rules

- $(c)' = 0$
- $(x)' = 1$
- $(x^2)' = 2x$
- $(c \cdot g(x))' = c \cdot g(x)'$
- $(g(x) \cdot h(x))' = g(x)' \cdot h(x) + g(x) \cdot h(x)'$
- $(g(h(x)))' = g'(h(x)) \cdot h'(x)$

PGD attack (= iterated, projected FGSM)

1. Generate a random starting point inside ℓ_∞ -ball
2. For $i=1 \dots k$
 - 2.1. Perform FGSM with magnitude ϵ_s
 - 2.2. Project back to the ϵ -sized ℓ_∞ -ball
- we start with random point x' inside ϵ -sized ℓ_∞ -ball around x
- we iteratively update x' with new perturbation, but always keep staying in that ball
- The ball ensures, that we aren't perturbing the image too much
- It is possible to never step outside the box and then projection will have no effect
- After running it k times, the loss function is maximized but there is no guarantee, that the image is classified to another label.

Adversarial Accuracy vs. Test Accuracy

- Given: ϵ - ℓ_∞ ball
- 100 examples
 - 5 wrongly classified
 - Additional 15 classified wrongly by an adversarial example inside ϵ - ℓ_∞ -ball
- Adversarial Accuracy: $\frac{100 - 5 - 15}{100} = \frac{80}{100}$
- Test Accuracy: $\frac{100 - 5}{100} = \frac{95}{100}$

↳ Trade between Adversarial vs. Test Accuracy

- Adversarial Accuracy \leq Test Accuracy

↳ If $= \Rightarrow$ Robust Network

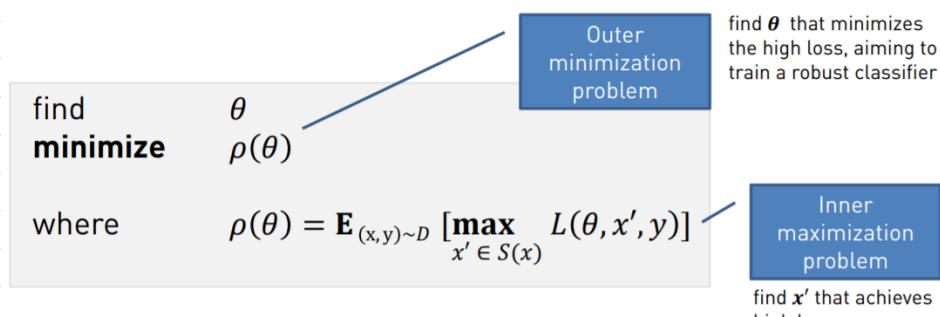
Training with adversarial examples (= PGD training)

- Making networks more robust

1. Select a mini-batch B of training data for adversarial-training
2. Compute for every data-point x an adversarial datapoint x' which maximizes the loss while still being inside the ℓ_∞ -ball
3. Modify existing network to be more robust by using the adversarial data points

$$\hookrightarrow \theta' = \theta - \frac{1}{\|x'\|} \cdot \sum_{(x_i, y) \in X'} \nabla_{\theta} L(\theta, x'_i, y)$$

- Large networks are more defendable
 - Training smaller networks with PGD has negative effects on test accuracy
 - By alternating solving the outer minimization and inner maximization step individually, we can get stuck at a local minimum
- ↳ This method is training the network to be experimentally robust



Soundness/Completeness/Exactness/Optimality

- Sound:
 - If there is a violation, the program always returns that the property is violated
 - e.g. return always false
 - $\forall z. F(\gamma(z)) \subseteq \gamma(F^{\#}(z))$

- Complete:
 - If the property is satisfied, the method can prove it
 - e.g. return always true

- Sound + Complete + not scalable: MILP
- Sound + Incomplete + scalable: Box, zonotope

- Not sound + incomplete: Random guessing

(\hookrightarrow) If a point inside the L_∞-cube gets misclassified, the verifier will always return violation (=sound)

(\hookrightarrow) Because of overapproximation we verify a larger area which can lead to returning violation because there is a violation in the overapproximation, but not in the L_∞-ball area itself (=Incomplete)

- Exactness (exact)

· If we enumerate all points in the set or just use the abstract transformer, do we get the same result (=No overapproximation)

- Optimality (best)

· Sound transformer F is called best transformer

if for all other sound transformer it is the most exact one

Incomplete Method: Box

Certification



Convex approximation

...

Before softmax



Region ϕ :

$$\begin{aligned}x_0 &= 0 \\x_1 &= 0.975 + 0.025\epsilon_1 \\x_2 &= 0.125 \\&\dots \\x_{784} &= 0.938 + 0.062\epsilon_{784}\end{aligned}$$

$\forall i. \epsilon_i \in [0,1]$

Output region ϕ_n

$$\begin{aligned}o_0 &= 0 \\o_1 &= 2.60 + 0.015\epsilon_0 + 0.023\epsilon_1 + 5.181\epsilon_2 + \dots \\o_2 &= 4.63 - 0.005\epsilon_0 - 0.006\epsilon_1 + 0.023\epsilon_2 + \dots \\&\dots \\o_9 &= 0.12 - 0.125\epsilon_0 + 0.102\epsilon_1 + 3.012\epsilon_2 + \dots\end{aligned}$$

$\forall i. \epsilon_i \in [0,1]$

- Interval for each coordinate using $\epsilon_i \in [0,1]$
- In the output region we can check, if o_3 is greater than all the other o_i , this means if the interval doesn't overlap and is the greatest

$$[a, b] +^{\#} [c, d] = [a + c, b + d] \quad \text{Addition transformer}$$

$$-^{\#}[a, b] = [-b, -a] \quad \text{Negation transformer}$$

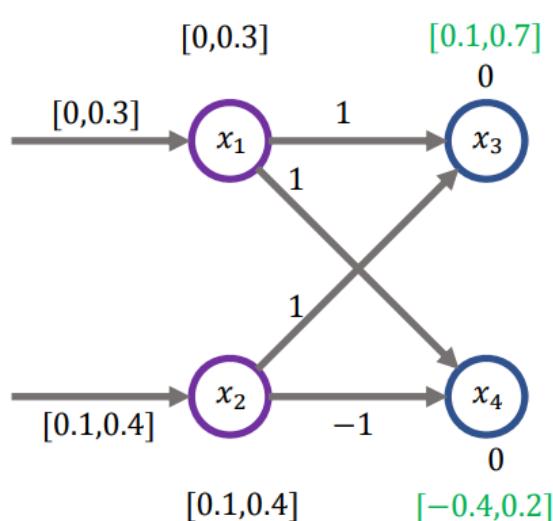
$$\text{ReLU}^{\#}[a, b] = [\text{ReLU}(a), \text{ReLU}(b)] \quad \text{ReLU transformer}$$

$$\lambda^{\#}[a, b] = [\lambda * a, \lambda * b] \quad \text{Multiplication by a constant } \lambda > 0$$

$$\Rightarrow [a, b] - [c, d]$$

$$\Leftrightarrow [a, b] + [-d, -c]$$

$$\Leftrightarrow [a-d, b-c]$$



Incomplete Method: Zonotope

- Replaces expensive ReLU transformation with affine form
- Unlike Box, Zonotope is exact for affine
- Like Box, Zonotope lose precision on ReLU

$$\hat{x}_1 = a_0^1 + \sum_{i=1}^k a_i^1 \epsilon_i \quad \begin{cases} \text{Noise} \\ \in [-1, 1] \end{cases}$$

.....

$$\hat{x}_d = a_0^d + \sum_{i=1}^k a_i^d \epsilon_i \quad \begin{cases} \text{Definition} \\ \text{Base coordinate} \end{cases}$$

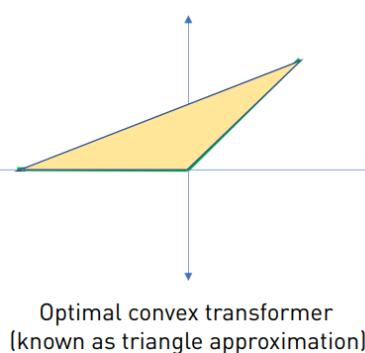
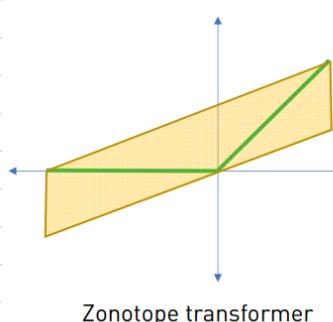
Multiplication of the p 'th neuron by a real-valued constant C :

$$(a_0^p + \sum_{i=1}^k a_i^p \epsilon_i) * C = (C * a_0^p + \sum_{i=1}^k C * a_i^p \epsilon_i)$$

Affine Transform

Adding two neurons p and q is done component-wise:

$$(a_0^p + \sum_{i=1}^k a_i^p \epsilon_i) + (a_0^q + \sum_{i=1}^k a_i^q \epsilon_i) = (a_0^q + a_0^p) + \sum_{i=1}^k (a_i^p + a_i^q) \epsilon_i$$



ReLU Transform

- 1. Case: If $u_x \leq 0 \Rightarrow y = 0$
- 2. Case: If $l_x > 0 \Rightarrow y = x$
- 3. Case: If $l_x < 0$ and $u_x > 0 \Rightarrow$ Crossing Boundary

$$\lambda * \hat{x} \leq y(\hat{x}) \leq \lambda * \hat{x} - \lambda * l_x$$

Lower Line

Upper Line

$$\lambda := \frac{u_x}{u_x - l_x}$$

The final ReLU transformer is:

$$y(\hat{x}) = \lambda * \hat{x} - \epsilon_{\text{new}} * \frac{\lambda * l_x}{2} - \frac{\lambda * l_x}{2}$$

Complete Method: MILP

- Checks, by discovering different paths, by setting $a_i \in \{0, 1\}$, and checking, if all

$$\begin{array}{l} \min c_1x_1 + c_2x_2 + \dots + c_nx_n \quad \text{object} \\ a_{11}x_1 + \dots + a_{1n}x_n \leq b_1 \\ \dots \\ a_{m1}x_1 + \dots + a_{mn}x_n \leq b_m \\ l_i \leq x_i \leq u_i \quad 1 \leq i \leq n \\ x_j \in \mathbb{Z} \end{array}$$

} constraints
 } bounds on variables
 } integer/binary

- ReLU Encoding

$\begin{cases} y \leq x - \epsilon \\ y \geq x \\ y \leq 0 \\ y > 0 \end{cases}$	$\begin{cases} y \leq x \\ y \geq x \\ y \leq u \\ y \geq 0 \end{cases}$
$\leftarrow a=0$ $\rightarrow a=1$	$y = x$ $y \in [0, u]$

} Running Time: $2^{\max(\dots)}$
 (We need a new branch for each max-operator)

- L₀-Ball encoding

$$x_i - \epsilon \leq x'_i \leq x_i + \epsilon$$

- Postcondition encoding

$$\min o_0 - o_1$$

After MILP finishes, we can check, if

- $o_0 - o_1 < 0 \Rightarrow$ Class 1
 - $o_0 - o_1 > 1 \Rightarrow$ Class 0
- } We must check this for every branch that is satisfiable at the end

$$\min o_0 - o_1$$

Plug in the affine and ReLU MILP encodings as defined

$l_i \leq x_i^p \leq u_i$ pre-computed Box bounds on all neurons in each layer p

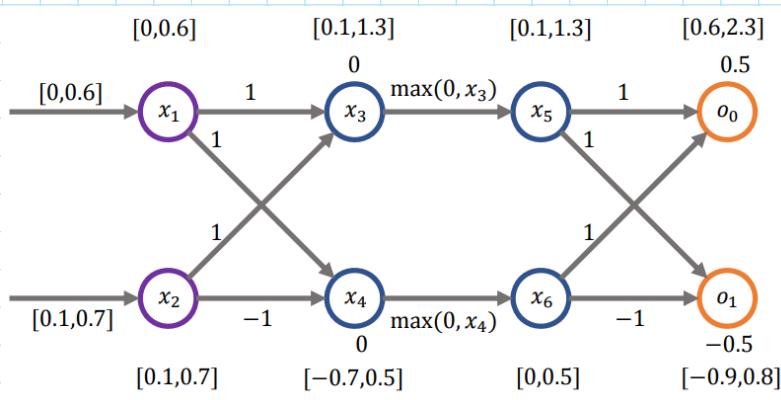
$x_i - \epsilon \leq x'_i \leq x_i + \epsilon$ ϕ : bounds on input neurons

$a_j \in \{0, 1\}$ These are the $a \in \{0, 1\}$ vars from the ReLU encoding

} MILP uses the ϕ constraints to save generating branches for a_i values of ReLU

For example, if Box-Bounds are positive, we can just use $x_5 - \text{ReLU}(x_5)$

$$x_5 = x'_5$$



$$\min o_0 - o_1$$

Affine

$$\begin{aligned}x_3 &= x_1 + x_2 \\x_4 &= x_1 - x_2 \\o_0 &= x_5 + x_6 + 0.5 \\o_1 &= x_5 - x_6 - 0.5\end{aligned}$$

ReLU: $x_5 = \max(0, x_3)$

$$\begin{aligned}x_5 &\leq x_3 - 0.1 * (1 - a_5) \\x_5 &\geq x_3 \\x_5 &\leq 1.3 * a_5 \\x_5 &\geq 0\end{aligned}$$

ReLU: $x_6 = \max(0, x_4)$

$$\begin{aligned}x_6 &\leq x_4 + 0.7 * (1 - a_6) \\x_6 &\geq x_4 \\x_6 &\leq 0.5 * a_6 \\x_6 &\geq 0\end{aligned}$$

Input bounds

$$\begin{aligned}0 \leq x_1 &\leq 0.6 \\0.1 \leq x_2 &\leq 0.7\end{aligned}$$

Pre-computed Box bounds

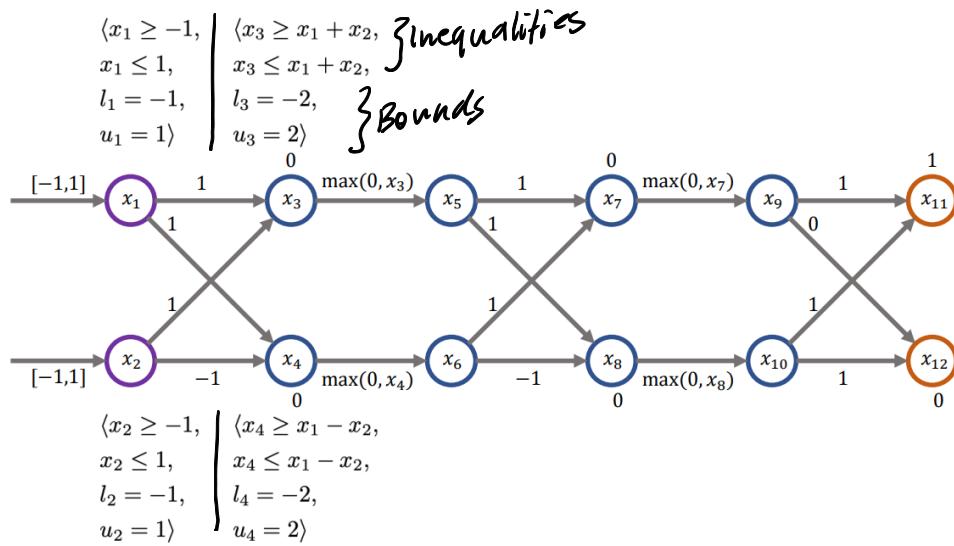
$$\begin{aligned}0.1 \leq x_3 &\leq 1.3 \\-0.7 \leq x_4 &\leq 0.5 \\0.1 \leq x_5 &\leq 1.3 \\0 \leq x_6 &\leq 0.5\end{aligned}$$

Binary integer variables

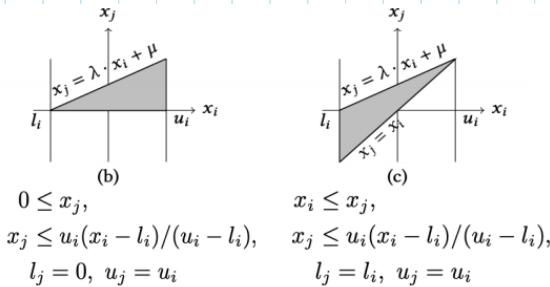
$$a_5, a_6 \in \{0, 1\}$$

Deep Poly (certification)

- We keep two inequalities and two bounds per neuron



ReLU activation



- choose one of the two abstract transformers, depending on which has the smaller area
 ↳ smaller area = more precise

Affine Transformation

. use backsubstitution

to get a more

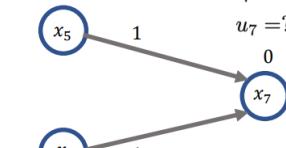
precise bound

- substitute until

you reach x_1/x_2

$$\left. \begin{array}{l} \langle x_5 \geq 0, \\ x_5 \leq 0.5 \cdot x_3 + 1, \\ l_5 = 0, \\ u_5 = 2 \rangle \end{array} \right.$$

$$\left. \begin{array}{l} \langle x_7 \geq 0, \\ x_7 \leq 0.5 \cdot x_3 + 0.5 \cdot x_4 + 2, \\ l_7 = ?, \\ u_7 = ? \rangle \end{array} \right.$$



$$\left. \begin{array}{l} \langle x_6 \geq 0, \\ x_6 \leq 0.5 \cdot x_4 + 1, \\ l_6 = 0, \\ u_6 = 2 \rangle \end{array} \right.$$

Replace the bounds for x_7 using the ones from the previous layer:

$$\begin{aligned} x_7 &\leq x_5 + x_6 \\ &\leq (0.5x_3 + 1) + (0.5x_4 + 1) \\ &= 0.5x_3 + 0.5x_4 + 2 \\ &= 0.5(x_3 + x_4) + 0.5(x_4 - x_2) + 2 \\ &= \underbrace{x_1 + 2}_{x_3 \leq x_1 + 2} + \underbrace{\frac{1}{2}(x_4 - x_2)}_{x_4 \leq x_2} + 2 \\ &= x_1 + 2 + \frac{1}{2}(x_4 - x_2) + 2 \\ &= x_1 + 2 + \frac{1}{2}x_4 - \frac{1}{2}x_2 + 2 \\ &= x_1 + 2 + \frac{1}{2}x_4 - \frac{1}{2}(x_1 + 1) + 2 \\ &= \frac{1}{2}x_1 + \frac{1}{2}x_4 + 2 \end{aligned}$$

$$x_7 \leq x_1 + 2 \quad u_7 = 3$$

- To prove something like $x_{11} - x_{12}$, instead of using the naive approach by using b_{11}, u_{12} , we use backsubstitution to prove robustness

Visualizing CNNs

- Filter: $\mathbb{R}^{C_0 \times h \times w \times C_1}$

- height

- width

- C_0 : #input channels

- C_1 : #output channels

1. First layers as filters

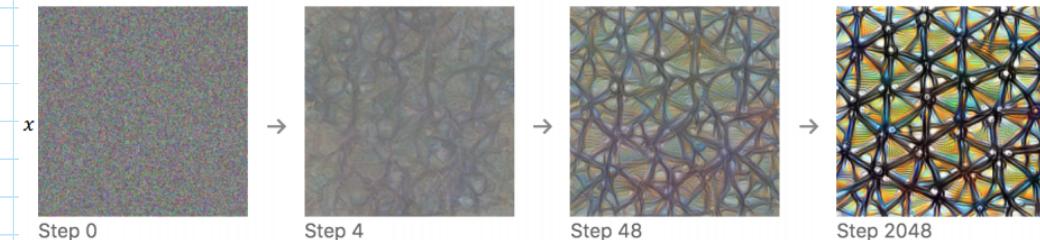
- Take the filters as images to visualize them

- Only the first layers possible, where $C_0 \leq 3$

2. Optimization

- Initialize the input image with white noise

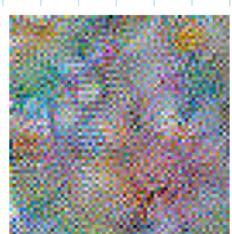
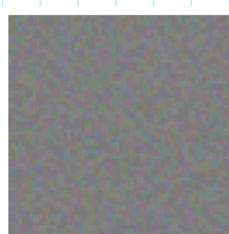
- Maximize gradient ascend to get maximum activation of the particular neuron/filter/channel



- Converges to local minimum depending on input-image

- You can also set the score you want to maximize to multiple neurons/channels/filters by taking the mean over their activation

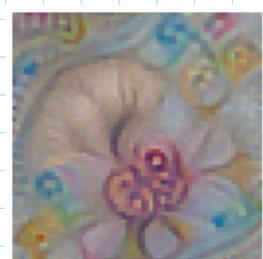
- When there is no regularization (Lasso/Ridge) inside the network, the patterns are unidentifiable



- Use regularization

- Or train with adversarial examples

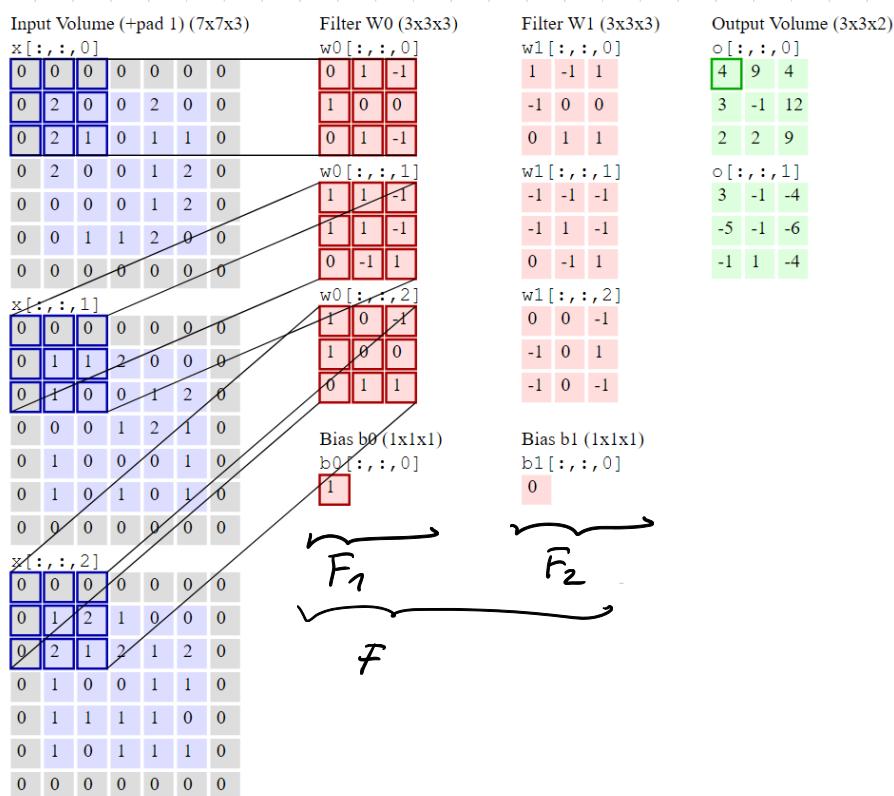
- ↳ Robust Networks



3. Deep Dream

- We take as input a real image
- Run gradient ascend on it to optimize early convolutions
- Early layers produce easy patterns such as lines
- While later layers show more high level concepts

Convolutions



\checkmark x \checkmark $output$
 $- x^{1 \times 5 \times 5 \times 3}$ (Padding: 1 / stride: 1)

$- F_1^{3 \times 3 \times 3 \times 1}$

$- F_2^{3 \times 3 \times 3 \times 1}$

\uparrow $Input$

$- F^{3 \times 3 \times 3 \times 2}$

Maxpooling:

Single depth slice

1	0	2	3
4	6	6	8
3	1	1	0
1	2	2	4



- Dimension reduction
- Get max. value of each kernel

Image Attribution

1. Gradient based

- direction of steepest change \Rightarrow overlay them
- Does not work well if feature is at a saddle point

2. Integrated Gradients

- Use an additional baseline picture x'
 - e.g. black or an adversarial example
- We are going the path from the base image to a combination of both to the original image
- We are integrating the gradient as we are going along and adding the features

Attribution Problem



3. Meaningful Perturbations

- Blur image in such a way, that the probability of the feature gets minimized

4. Shapley Values

- Very expensive
- For each feature we check, if we add a certain subset of the image, if this increases or decreases our activation function

Quering the Network with logic (= find formula example)

$$\theta = \bigwedge_{j=1, j \neq 9}^k \text{NN}(i)[j] < \text{NN}(i)[9] \wedge \|i - \text{deer}\|_\infty < 25 \wedge \|i - \text{deer}\|_\infty > 5 \quad \left\{ \begin{array}{l} i = \text{image} \\ \text{bias}(\text{NN}(i)) = 9 \end{array} \right.$$

↳ "Find an image i which gets classified to 9 (=Truck) where the image is within some distance to the image deer."

1. SAT solver

- Only if neural network is really small

2. Translate to optimization problem

- Translation function $T(\theta)$

- $T(\theta)(i) = 0 \Leftrightarrow i \text{ satisfies } \theta$

Logical Term	Translation
$t_1 \leq t_2$	$\max(0, t_1 - t_2)$
$t_1 \neq t_2$	$[t_1 = t_2]$
$t_1 = t_2$	$T(t_1 \leq t_2 \wedge t_2 \leq t_1)$
$t_1 < t_2$	$T(t_1 \leq t_2 \wedge t_1 \neq t_2)$
$\varphi \vee \psi$	$T(\varphi) \cdot T(\psi)$
$\varphi \wedge \psi$	$T(\varphi) + T(\psi)$

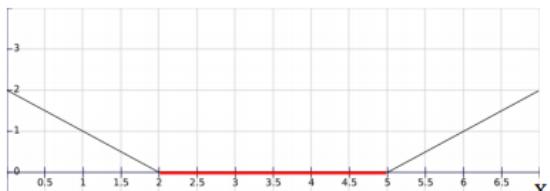
$\left\{ \begin{array}{l} \text{As long as } t_1 \leq t_2, \text{ the value is 0} \\ \text{if } t_1 = t_2 \Rightarrow \text{value is 1} \\ \max(0, t_1 - t_2) + \max(0, t_2 - t_1) = |t_1 - t_2| \\ \max(0, t_1 - t_2) + [t_1 = t_2] \end{array} \right.$

Assume $\theta(T) = 0$

Example

↳ To prove: Base case: Equalities / Inequalities
Step case: Disjunctions / Conjunctions

$$T(x \geq 2 \wedge x \leq 5) = \max(0, 2-x) + \max(0, x-5)$$



$\left\{ \begin{array}{l} \text{The function is 0 when } x \text{ is between 2-5} \end{array} \right.$

$$\neg(a \wedge b) = \neg a \vee \neg b$$

$$\neg(a \vee b) = \neg a \wedge \neg b$$

$$\neg(a \leq b) = a > b$$

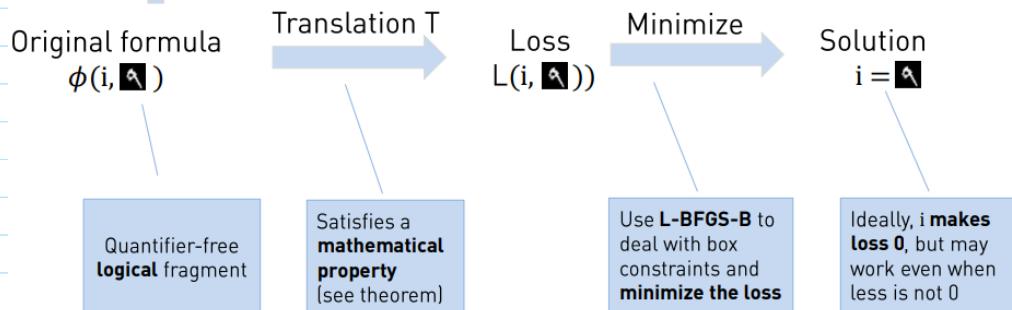
$$\neg(a = b) = a \neq b$$

$NN(i)[1] < NN(i)[9] \wedge$ 
 $NN(i)[2] < NN(i)[9] \wedge$ 
 $NN(i)[3] < NN(i)[9] \wedge$ 
 $NN(i)[4] < NN(i)[9] \wedge$ 
 $NN(i)[5] < NN(i)[9] \wedge$
 $NN(i)[6] < NN(i)[9] \wedge$
 $NN(i)[7] < NN(i)[9] \wedge$
 $NN(i)[8] < NN(i)[9] \wedge$
 $\|i - deer\|_\infty < 25 \wedge$
 $\|i - deer\|_\infty > 5$ 

$\max(0, NN(i)[1] - NN(i)[9]) + [NN(i)[1] = NN(i)[9]]$
 $+ \max(0, NN(i)[2] - NN(i)[9]) + [NN(i)[2] = NN(i)[9]]$
 $+ \max(0, NN(i)[3] - NN(i)[9]) + [NN(i)[3] = NN(i)[9]]$
 $+ \max(0, NN(i)[4] - NN(i)[9]) + [NN(i)[4] = NN(i)[9]]$
 $+ \max(0, NN(i)[5] - NN(i)[9]) + [NN(i)[5] = NN(i)[9]]$
 $+ \max(0, NN(i)[6] - NN(i)[9]) + [NN(i)[6] = NN(i)[9]]$
 $+ \max(0, NN(i)[7] - NN(i)[9]) + [NN(i)[7] = NN(i)[9]]$
 $+ \max(0, NN(i)[8] - NN(i)[9]) + [NN(i)[8] = NN(i)[9]]$
 $+ \max(0, \|i - deer\|_\infty - 25) + [\|i - deer\|_\infty = 25]$
 $+ \max(0, 5 - \|i - deer\|_\infty) + [\|i - deer\|_\infty = 5]$

Take out box constraints and pass them separately

e.g. $i \in [0, 1]$



Training the Network with Logic (=ensure formula)

- How to find counter-examples
- e.g. all images in dataset that are classified as a car have a higher probability for the label truck than the probability for label dog
 - $\hookrightarrow y = \text{car} \Rightarrow \text{NN}(x)[\text{truck}] > \text{NN}(x)[\text{dog}]$
 - $\hookrightarrow \forall z \in L_\infty(x, \varepsilon). y = \text{car} \Rightarrow \text{NN}(z)[\text{truck}] > \text{NN}(z)[\text{dog}]$
 - \hookrightarrow For all images close to x

1. Train a base classifier $\hat{\theta}$ on labeled data
2. Use $\hat{\theta}$ to infer the labels for unlabeled dataset
3. Use adversarial training to learn robust classifier

Problem statement

find θ 3 find weights for Neural Network
minimize $p(\theta)$

where $p(\theta) = \mathbf{E}_{s \sim D} [T(\phi)(z_{\text{worst}}, s, \theta)]$ } Minimize the expected value
 of the worst counter-example

and $z_{\text{worst}} = \underset{z}{\operatorname{argmin}}(T(\neg\phi)(z, s, \theta))$

Find the worst case counter example

high loss

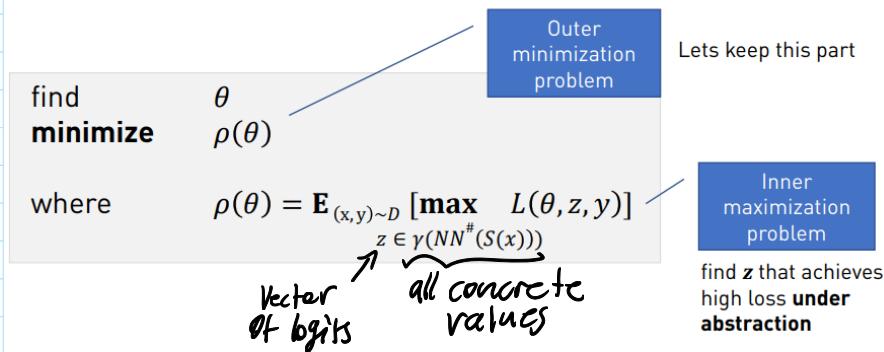
- We want to find the worst possible violation of the formula and then try to find a network that minimizes its effect

- Example formula: $\|\phi(z, x, \theta)\|_\infty \leq \varepsilon \Rightarrow \text{NN}_\theta(z)[3] > \sigma$
- \hookrightarrow Find an image z which is inside the L_∞ ball around the image x where the ball has size ε and z gets classified as 3 is greater than σ

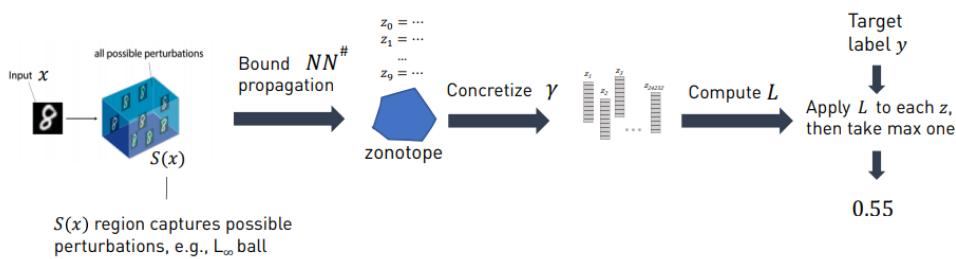
Provably defenses

and training the network with logit

- PGD training only trains the network experimentally
- we want to train the network to be provably robust after training

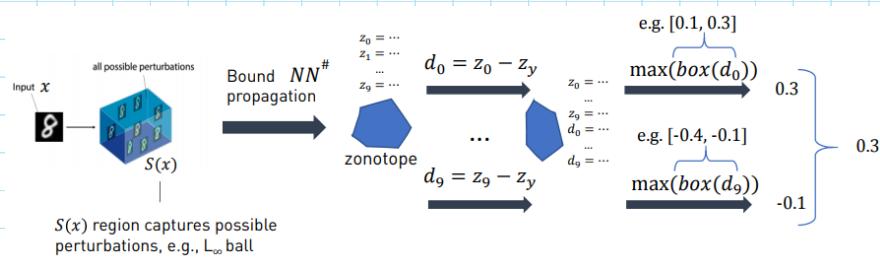


- We start at the beginning and push $S(x)$ through the network to get the area which we get values from
- While in PGD training we look at the end



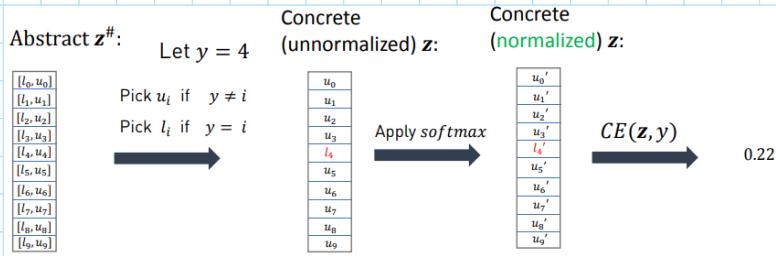
- We could have a large number of logits in the vector, enumerating all possible values impossible

1. Use zonotope to push through, convert the logits to box intervals and look for max loss



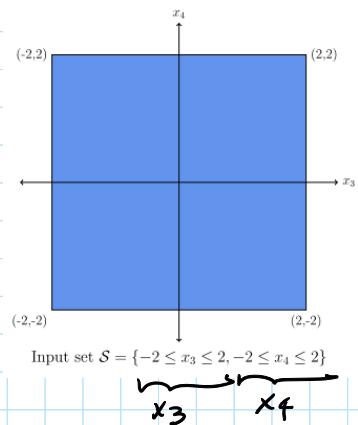
Take that adversarial example to train network

2. Use zonotope to push through, convert the logits to box intervals and use cross-entropy loss

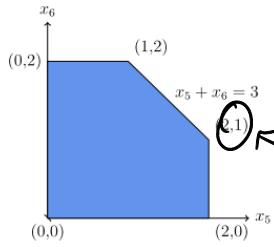


K-Relu

- Normal Relu transformer is neuron-wide

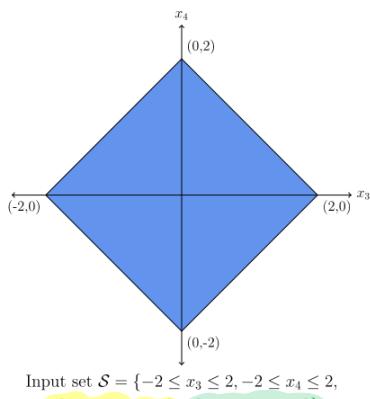


$$\begin{aligned}x_5 &:= \text{ReLU}(x_3) \\x_6 &:= \text{ReLU}(x_4)\end{aligned}$$

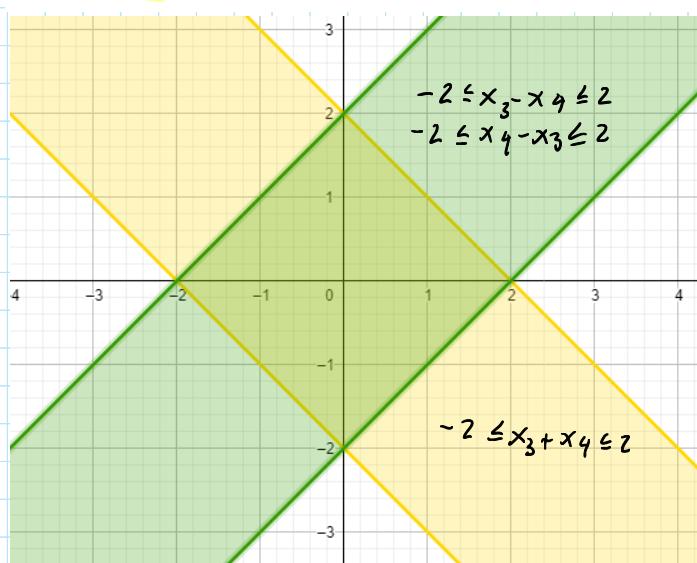


OK Not possible
↓ overapproximation

- Put constraints between two/more Relus



- If $x_3 \leq 0$, then $\max(0, x_3) = 0 = x_5$
 - $S \cap \{x_3 \leq 0, x_4 \leq 0, x_5 = 0, x_6 = 0\}$
 - $S \cap \{x_3 \leq 0, x_4 \geq 0, x_5 = 0, x_6 = x_4\}$
 - $S \cap \{x_3 \geq 0, x_4 \leq 0, x_5 = x_3, x_6 = 0\}$
 - $S \cap \{x_3 \geq 0, x_4 \geq 0, x_5 = x_3, x_6 = x_4\}$
- ↳ enumerate all possibilities



Randomized smoothing

- Guarantees are not absolute (like MILP), but probabilistic
- construct from a given classifier f a smoothed classifier \bar{f}
- \bar{f} returns the class with highest probability under $N(x, \sigma^2 I)$
- $\bar{f}(x) := \operatorname{argmax}_{\epsilon \in \mathbb{R}^d} P[\bar{f}(x + \epsilon) = c_k]$ where $\epsilon \sim N(0, \sigma^2 I)$
 $= \operatorname{argmax}_{\epsilon \in \mathbb{R}^d} \int [f(x + \epsilon) = c_k] \cdot d\epsilon$ ↗ Infeasible, too large
 $\approx \operatorname{argmax}_{c \in \mathcal{Y}} \left(\frac{1}{n} \sum_{i=1}^n [f(x + \epsilon_i) = c] \right)$

```

function PREDICT( $f, \sigma, x, n, \alpha$ )
    counts  $\leftarrow$  SAMPLEUNDERNOISE( $f, x, n, \sigma$ ) → Evaluates  $n$  times  $f(x + \epsilon)$ 
     $\hat{c}_A, \hat{c}_B \leftarrow$  top two indices in counts and counts classes
     $n_A, n_B \leftarrow$  counts[ $\hat{c}_A$ ], counts[ $\hat{c}_B$ ] →  $n_A$ : # class A is chosen
    if BINOPVALUE( $n_A, n_A + n_B, 0.5 \leq \alpha$ ) return  $\hat{c}_A$ 
    else return ABSTAIN
        ↗  $P[n_A \text{ of } (n_A + n_B) \text{ samples came up}]$ 

```

- BinopValue returns the probability that n_A out of $(n_A + n_B)$ times we came up with the wrong class, because we could have sampled from wrong region
- Predict returns wrong class with probability at most α

```

function CERTIFY( $f, \sigma, x, n_0, n, \alpha$ )
    counts0  $\leftarrow$  SAMPLEUNDERNOISE( $f, x, n_0, \sigma$ ) } Get  $\hat{c}_A$ 
     $\hat{c}_A \leftarrow$  top index in counts0
    counts  $\leftarrow$  SAMPLEUNDERNOISE( $f, x, n, \sigma$ ) } Sample again to avoid introducing
     $p_A \leftarrow$  LOWERCONFBOUND(counts[ $\hat{c}_A$ ],  $n, 1 - \alpha$ ) ↗ Biases
    if  $p_A > \frac{1}{2}$  return prediction  $\hat{c}_A$  and radius  $\sigma \Phi^{-1}(p_A)$ 
    else return ABSTAIN

```

↗ With probability at least $1 - \alpha$ class C_A is the real class if we sample $\bar{f}(x + \sigma) = C_A$, $\sigma \in R$

- Adversarial attacks using randomized smoothing

- argmax in $\hat{f}(x) = \underset{c_k \in Y}{\operatorname{argmax}} (P_\epsilon[f(x+\epsilon)=c_k])$ is

a problem for gradient descend

(\rightarrow we assume softmax layer at the end)

- $x_{\text{adv}} = \underset{\substack{\|x'-x\| \leq p \\ }}{\operatorname{argmax}} (-\log E_\epsilon [F(x'+\epsilon)_y])$

we find adversarial example x' within the L_∞ -ball

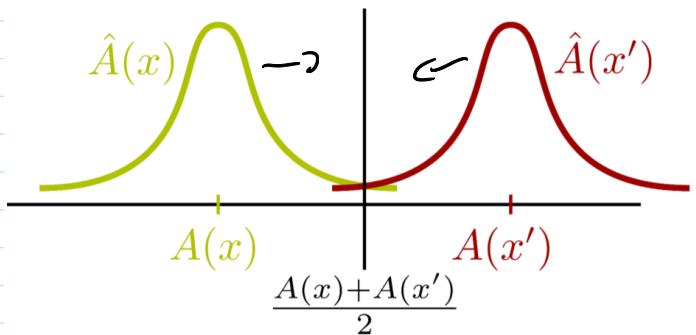
$$= \nabla_{x'} (-\log E_\epsilon [F(x'+\epsilon)_y])) \quad \text{Infeasible, too large}$$

$$\approx \nabla_{x'} (-\log \left(\frac{1}{m} \sum_{i=1}^m F(x'+\epsilon)_y \right)) \quad \begin{array}{l} \text{Use Monte-Carlo to} \\ \text{approximate gradient} \end{array}$$

- instead of argmaxing $P[f(x+\epsilon)=c_k]$ we are
argmaxing the highest loss a perturbed input x'
can make $-\log E[F(x'+\epsilon)_y]$

Differential Privacy

- Good composition properties: $B(x) = (\hat{A}_1(x), \dots, \hat{A}_n(x))$, $\sum_{i=1}^n \epsilon_i = \epsilon$ DP
- Preserved under post processing: $g(\hat{A}(x))$
- The probability distribution of the database with the user and without it is nearly the same
- We want to bring $\hat{A}(x)$ and $\hat{A}(x')$ as close together as possible, because $[-\infty, \frac{A(x)+A(x')}{2}]$ is probably the real dataset x and not the perturbed x' dataset



- ϵ -differential privacy:
 - $\Pr[\hat{A}(x) \in \Phi] \leq \exp(\epsilon) \cdot \Pr[\hat{A}(x') \in \Phi]$
 - $(x, x') \in \text{Neighborhood}$
 - The two distributions are roughly the same
 - ϵ captures the notion of distance between $\hat{A}(x)$ and $\hat{A}(x')$
 - $\epsilon=0 : \Pr[\hat{A}(x) \in \Phi] = \Pr[\hat{A}(x') \in \Phi]$
 - Full privacy
 - Small $\epsilon : \Pr[\hat{A}(x) \in \Phi] \approx \Pr[\hat{A}(x') \in \Phi]$

Proof:

$$\Pr[\hat{A}(x) \in \Phi] \leq \exp(\epsilon) \Pr[\hat{A}(x') \in \Phi]$$

\Rightarrow (by symmetry of Neig)

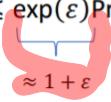
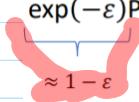
$$\Pr[\hat{A}(x') \in \Phi] \leq \exp(\epsilon) \Pr[\hat{A}(x) \in \Phi]$$

$\Rightarrow (\cdot \exp(-\epsilon))$

$$\exp(-\epsilon) \Pr[\hat{A}(x') \in \Phi] \leq \Pr[\hat{A}(x) \in \Phi]$$

\Rightarrow (combine)

$$\exp(-\epsilon) \Pr[\hat{A}(x') \in \Phi] \leq \Pr[\hat{A}(x) \in \Phi] \leq \exp(\epsilon) \Pr[\hat{A}(x') \in \Phi]$$



- Neighborhood:
 - Changing the data of one person
 - Adding/removing one person from/to x'
 - For images $\|x - x'\|_p \leq R$
 - \hookrightarrow Symmetric $(x, x') \Leftrightarrow (x', x)$

Tradeoff between Privacy > Precision

- Trivial choice $\hat{A}(x) = 0$
- $\epsilon = 0$, since we always return 0 \Rightarrow Fully private
- $\|\hat{A}(x) - A(x)\|$ can be large \Rightarrow Inprecise

Laplace mechanism

- $A(x) = \sum_{i=1}^d x_i$
- Sensitivity of A : $\Delta(A) = \max \|A(x) - A(x')\|_1 = 1$
- $\hat{A}(x) = \sum_{i=1}^d x_i + \text{Lap}(0, \frac{1}{\epsilon}) = \text{Lap}\left(\sum_{i=1}^d x_i, \frac{1}{\epsilon}\right)$

· We need to check: $\frac{\Pr[\hat{A}(x) \in \phi]}{\Pr[\hat{A}(x') \in \phi]} \leq \exp(\epsilon)$

$$= \frac{\Pr[\hat{A}(x) = t]}{\Pr[\hat{A}(x') = t]}$$

$$= \frac{\exp\left(-\frac{|t - A(x)|}{1/\epsilon}\right)}{\exp\left(-\frac{|t - A(x')|}{1/\epsilon}\right)}$$

$$= \exp\left(-\frac{|t - A(x)|}{1/\epsilon} + \frac{|t - A(x')|}{1/\epsilon}\right)$$

$$\leq \exp\left(-\frac{|A(x') - A(x)|}{1/\epsilon}\right) \quad |\Delta(A)| = 1$$

$$= \exp\left(-\frac{1}{1/\epsilon}\right)$$

$$= \exp(-\epsilon)$$

$$\leq \exp(\epsilon)$$

- Generalizing differential privacy

• (ϵ, σ) -DP: $\Pr[\hat{A}(x) \in \phi] \leq \exp(\epsilon) \cdot \Pr[\hat{A}(x') \in \phi] + \sigma^2$

• ϵ -DP: $(\epsilon, 0)$ -DP

$$\hat{A}(x) = A(x) + \eta$$

- Laplacean mechanism: works well with $\|x\|_1 = |3/r|q| = 7$

$$\eta_i \sim \text{Lap}(0, 1/\epsilon)$$

- Gaussian mechanism: works well with $\|x\|_2 = \sqrt{s^2 + q^2} = 5$

$$\eta \sim N(0, \sigma^2)$$

- ϵ -DP counterexamples

- Prove that $\hat{A}(x)$ not ϵ -DP

$$\log\left(\frac{\Pr[\hat{A}(x) \in \phi]}{\Pr[\hat{A}(x') \in \phi]}\right) \geq \epsilon$$

- Connection to randomized smoothing

• $\hat{f}(x) = f(x + \eta)$ where $\eta_i \sim \text{Lap}(0, R/\epsilon)$

• $(x, x') \in \text{Neigh} \iff \|x - x'\|_1 < R$

• \hat{f} is ϵ -DP

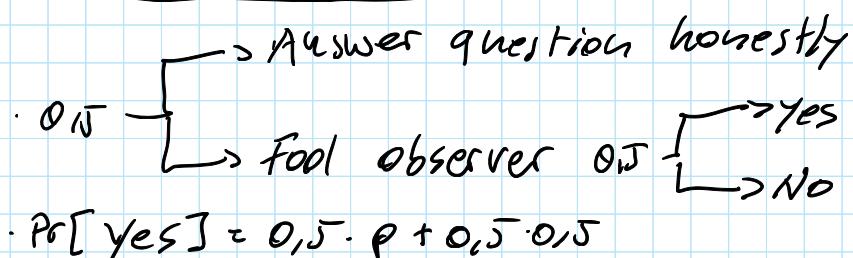
1. Consider $\hat{A}(x) = x + \eta$ and view f as post-processing

2. Generalize from $\|x - x'\|_1 < 1$ to $\|x - x'\|_1 < R$

by $\text{Lap}(0, R/\epsilon)$

$$\hookrightarrow x' = x + \eta$$

- Randomized Response



$$\Pr[\text{Yes}] = 0.5 \cdot p + 0.5 \cdot 0.5$$

$$\frac{\Pr[\text{Answer yes}|\text{Truly yes}]}{\Pr[\text{Answer yes}|\text{Truly no}]} = \frac{\frac{1}{2} \cdot 1 + \frac{1}{2} \cdot \frac{1}{2}}{\frac{1}{2} \cdot 0 + \frac{1}{2} \cdot \frac{1}{2}} = \frac{\frac{3}{4}}{\frac{1}{4}} = 3 \leq \exp \epsilon \Leftrightarrow \epsilon \geq 1.10 \Rightarrow 1,10-\text{DP}$$