

Chapter 1 - Introduction to adversarial robustness

- ✓ The model: $h_\theta: \underline{x} \rightarrow \mathbb{R}^K$ mapping from input space to the output space
 \underline{x} : input space
 K : number of classes being predicted
 θ : model's parameters
- ✓ Loss function: $\ell: \mathbb{R}^K \times \mathbb{Z}_+ \rightarrow \mathbb{R}_+$ mapping from model predictions and true labels to a non-negative number
 $\ell(h_\theta(\underline{x}), y) = \begin{cases} h_\theta(\underline{x})_y & \text{the model output (logits)} \\ y & \text{the index of the true class} \end{cases}$
- ✓ A common form of loss used in deep learning: Cross entropy loss

To maximize the log of the true probability of the true class label using softmax operator:

$$\log \left(\frac{\exp(h_\theta(\underline{x})_y)}{\sum_{j=1}^K \exp(h_\theta(\underline{x})_j)} \right) = h_\theta(\underline{x})_y - \log \left(\sum_{i=1}^K \exp(h_\theta(\underline{x})_i) \right) \Rightarrow$$

$$\text{minimize } \ell(h_\theta(\underline{x}), y) = \log \left(\sum_{i=1}^K \exp(h_\theta(\underline{x})_i) \right) - h_\theta(\underline{x})_y$$

- ✓ To train a classifier: $\underset{\theta}{\text{minimize}} \frac{1}{m} \sum_{i=1}^m \ell(h_\theta(\underline{x}_i), y_i)$

we solve it by stochastic gradient descent for some minibatch $B \subseteq \{1, \dots, m\}$

$$\theta := \theta - \frac{\alpha}{|B|} \sum_{i \in B} \nabla_{\theta} \ell(h_\theta(\underline{x}_i), y_i) \quad \alpha: \text{step size}$$

- ✓ To form an adversarial example: $\underset{\hat{x}}{\text{maximize}} \ell(h_\theta(\hat{x}), y)$

We need to ensure that \hat{x} is close to our original input x .

$$\underset{\delta \in \Delta}{\text{maximize}} \ell(h_\theta(x + \delta), y)$$

chapter 2 - Linear models

✓ minimize $\frac{1}{|D|} \sum_{x,y \in D} \max_{\|\delta\| \leq \epsilon} l(\hat{w}(x) + b, y)$

$\hat{w}(x)$

& Inner maximization can be solved exactly for the case of binary classification
 Minimization problem is still convex \rightarrow robust training procedure can also be solved optimally.
 However: deep network neither the inner maximization problem nor the outer minimization problem can be solved globally.

✓ Binary classification: Binary class entropy:

$$l(h_\theta(x), y) = \log(1 + \exp(-y \cdot h_\theta(x))) \equiv L(y \cdot h_\theta(x))$$

Proof: Probability of predicting class 2:

$$\frac{\exp(h_\theta(x)_2)}{\exp(h_\theta(x)_1 + h_\theta(x)_2)}$$

Probability of predicting class 1:

$$\frac{\exp(h_\theta(x)_1)}{\exp(h_\theta(x)_1 + h_\theta(x)_2)}$$

$$h'_\theta(x) = h_\theta(x)_1 - h_\theta(x)_2$$

$$\Rightarrow P(y|x) = \frac{1}{1 + \exp(-y \cdot h'_\theta(x))} \stackrel{-\log}{\longrightarrow} \log(1 + \exp(-y \cdot h'_\theta(x))) \text{ Logistic Loss}$$

✓ Solving the inner maximization:

Function of $L(z) = \log(1 + \exp(-z))$ is monotonic decreasing \Rightarrow

$$\max_{\|\delta\| \leq \epsilon} l(y \cdot (\hat{w}^T(x) + b)) = L(\min_{\|\delta\| \leq \epsilon} y \cdot (\hat{w}^T(x) + b)) = L(y \cdot (\hat{w}^T x + b) + \min_{\|\delta\| \leq \epsilon} y \cdot \hat{w}^T \delta)$$

$$\Rightarrow \min_{\|\delta\| \leq \epsilon} y \cdot \hat{w}^T \delta$$

$$\begin{aligned} &\text{if } y=+1, \|\delta\|_\infty \leq \epsilon \Rightarrow \begin{cases} \delta_i = -\epsilon & \text{for } w_i > 0 \\ \delta_i = \epsilon & \text{for } w_i < 0 \end{cases} \Rightarrow \delta^* = y \cdot \text{sign}(w) \\ &\text{if } y=-1 \text{ Flip} \end{aligned}$$

$$y \cdot \hat{w}^T \delta^* = -y^2 \epsilon \sum_i |w_i| = -\epsilon \|\hat{w}\|_1 \Rightarrow \underset{w,b}{\text{minimize}} \frac{1}{|D|} \sum_{(x,y) \in D} l(y \cdot (\hat{w}^T x + b) - \epsilon \|\hat{w}\|_1)$$

Convex in w, b

Intuitively: if a point is far from the decision boundary, we don't penalize the norm of the parameters, but we do penalize the norm of the parameters for a point where we close to the decision boundary.

✓ A common perturbation set to use is the ℓ_∞ ball:

$$\Delta = \{ \delta : \|\delta\|_\infty \leq \epsilon \}, \quad \|\zeta\|_\infty = \max_i |\zeta_i|$$

Since each step is followed by a projection back onto the ℓ_∞ ball, this procedure is known as **Projected gradient descent (PGD)**

✓ targeted attacks

We maximize the loss of the correct class while also minimize the loss of the target class.

$$\begin{aligned} & \underset{\delta \in \Delta}{\text{maximize}} (l(h_{\theta}(\mathbf{x} + \delta), y) - l(h_{\theta}(\mathbf{x} + \delta), y_{\text{target}})) = \underset{\delta \in \Delta}{\text{maximize}} \\ & (h_{\theta}(\mathbf{x} + \delta)y_{\text{target}} - h_{\theta}(\mathbf{x} + \delta)y) \end{aligned}$$

✓ The risk of a classifier is its expected loss under the true distribution of samples

$$R(h_{\theta}) = E_{(\mathbf{x}, y) \sim D} [l(h_{\theta}(\mathbf{x}), y)] \quad D: \text{true distribution over samples}$$

In practice:

$$\hat{R}(h_{\theta}, D) = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} l(h_{\theta}(\mathbf{x}), y) \Rightarrow \text{minimize } \hat{R}(h_{\theta}, D_{\text{train}})$$

adversarial risk:

$$\hat{R}_{\text{adv}}(h_{\theta}, D) = \frac{1}{|D|} \sum_{(\mathbf{x}, y) \in D} \max_{\delta \in \Delta(\mathbf{x})} l(h_{\theta}(\mathbf{x} + \delta), y)$$

If we are truly operating in an adversarial environment, adversarial risk would provide a more accurate estimate of the expected performance of the classifier.

✓ To train adversarially robust classifiers:

$$\underset{\theta}{\text{minimize}} \quad \frac{1}{|D_{\text{train}}|} \sum_{(\mathbf{x}, y) \in D_{\text{train}}} \max_{\delta \in \Delta(\mathbf{x})} l(h_{\theta}(\mathbf{x} + \delta), y)$$

$$\theta := \theta - \frac{\alpha}{|B|} \sum_{(\mathbf{x}, y) \in B} \nabla_{\theta} \max_{\delta \in \Delta(\mathbf{x})} l(h_{\theta}(\mathbf{x} + \delta), y)$$

How to solve it? **Danskin's theorem**

$$\delta^* = \arg \max_{\delta \in \Delta(\mathbf{x})} l(h_{\theta}(\mathbf{x} + \delta), y) \Rightarrow \nabla_{\theta} \max_{\delta \in \Delta(\mathbf{x})} l(h_{\theta}(\mathbf{x} + \delta), y) = \nabla_{\theta} (h_{\theta}(\mathbf{x} + \delta^*), y)$$

We iteratively compute adversarial examples and then update the classifier based not upon the original data points, but upon these adversarial examples.

chapter 3: Adversarial examples solving the inner maximization.

- ✓ Neural network is much more sensitive to adversarial attacks than even the naive linear models and much harder to protect against

$$\begin{aligned} & \text{Maximize } l(h_\theta(x), y) \\ & \|s\| \leq \epsilon \end{aligned}$$

$h_\theta(x)$ represents a neural network:

$$z_1 = x$$

$$z_{i+1} = f_i(w_i z_i + b_i) \quad i, \dots, d$$

$$h_\theta(x) = z_{d+1}$$

ReLU operator $f_i(z_i) = \max\{0, z_i\}$ for $i=1, \dots, d-1$

Linear operator $f_d(z_d)$ for layer d

- ✓ In the high dimensional setting that we will consider with deep networks, there is a high likelihood that at most any point in input space, there will be some direction along the loss surface that is very steep.

* moving a small distance in input space \Rightarrow a big increase in loss \Rightarrow NNs are prone to adversarial examples.

* The cost surface for NN is not convex \Rightarrow is prone to Local optima

* Many directions of high cost increase \Rightarrow an adversarial example
No optimal adversarial example \Rightarrow training robust networks will be a big problem

- ✓ strategies for the inner maximization: Danskin's theorem X

1: lower bounds

2: exact solutions

3: upper bounds

- ✓ Three different architectures here: 1: two-layer fully-connected
2: four-layer fully-connected
3: four convolutional layers + one fully-connected

✓ Lower bounding the inner maximization:

1: using backpropagation 2: δ needs to stay within the norm bound ϵ

✓ Fast Gradient Sign method (FGSM)

$$g := \nabla_{\delta} l(h_{\theta}(x+\delta), y) \quad \delta := \delta + \alpha g \\ \| \delta \|_{\infty} \leq \epsilon$$

If we want to make increase the loss as much as possible, it makes sense to take as large a step as possible. & for $\delta=0$

if $\alpha \gg \epsilon \Rightarrow$ the relative sizes of the entries of g won't matter \Rightarrow

δ can be considered to be either $+\epsilon$ or $-\epsilon \Rightarrow \delta := \epsilon \cdot \text{sign}(g)$

Results: Fully connected networks are particularly susceptible to minor changes and ConvNets slightly less but of course they are still very sensitive.

Points: FGSM is an attack under an l_{infinity} norm bound. FGSM is exactly the optimal attack against a linear binary

NNs are not in fact linear even over a relatively small region.

✓ Projected Gradient descent

Repeat:

$$\delta := P(\delta + \alpha \nabla_{\delta} l(h_{\theta}(x+\delta), y))$$

P: Projection on to the ball of interest

We need to determine actual step size and the number of iterations.

In this method, step size α is too large. $\xrightarrow{\text{To solve this issue}}$

Normalized steepest descend method

traditional GD: $z := z - \alpha \nabla_z f(z)$

$$z := z - \alpha \frac{\nabla_z f(z)}{\| \nabla_z f(z) \|_2} \quad \text{for } l_2$$

Normalized steepest descent method: $z := z - \arg \max_{\| z \| \leq \alpha} z^T \nabla_z f(z) \Rightarrow \| z \| \leq \alpha$

$$z := z - \alpha \text{Sign}(\nabla_z f(z)) \quad \text{against linear func} \quad \text{for } l_{\infty}$$

Result: PGD is not much more effective than FGSM

Normalized steepest descent is doing a lot better than the FGSM algorithm.

- * In normalized steepest descent, α is on the same scale as the total perturbation bound ϵ .
- * We can use randomization to improve the performance of algorithm.
Randomization mitigates the problem of local optima
we run it multiple times from different random locations. \Rightarrow increases runtime \Rightarrow may not be practical.

✓ Targeted attacks

We can modify our objective to maximize the target class logit:

$$\text{maximize}_{\delta} \left(h_{\theta}(x+\delta)_{y_{\text{targ}}} - \sum_{y' \neq y_{\text{targ}}} h_{\theta}(x+\delta)_{y'} \right)$$

This func. is more difficult than the previous one, we aren't able to fool the classifier as much, but it is more able to predict the target class.

✓ Non- ℓ_∞ norm

$$\delta := P_E \left(\delta - \alpha \frac{\nabla_{\delta} l(h_{\theta}(x+\delta), y)}{\| \nabla_{\delta} l(h_{\theta}(x+\delta), y) \|} \right)$$

Projection on to the ℓ_2 ball of radius ϵ $\leftarrow P_E(z) = \epsilon \frac{z}{\max\{\epsilon, \|z\|_2\}}$

* ϵ for ℓ_2 norm is larger than ϵ for ℓ_∞ norm

✓ ℓ_∞ attacks lead to small noise everywhere in image

✓ ℓ_2 is more localized in the image

✓ ℓ_1 encourages sparsity in the $\delta \Rightarrow$ a few pixels locations are adjusted.

✓ It is more useful to refer to attacks more formally by 1) the norm ball perturbation they consider 2) the method they use for optimizing over that norm ball.

✓ Exactly solving the inner maximization

It is for small NNs.

$$z_1 = x$$

$$z_{i+1} = f_i(w_i z_i + b_i) \quad i=1, \dots, d$$

$$h_\theta(x) = z_{d+1}$$

$$f_i(z) = \text{ReLU}(z) \text{ for } i=1, \dots, d-1 \quad \& \quad f_d(x) = x$$

we focus on l_∞

we can rewrite the problem as

$$\text{minimize } (e_y - e_{y\text{targ}})^T z_{d+1}$$

$$\text{Subject to } \|z_i - x\|_\infty \leq \epsilon$$

non-convex & is not handled by most optimization solvers

$$\leftarrow z_{i+1} = \max\{0, w_i z_i + b_i\}, \quad i=1, \dots, d-1$$

$$z_{d+1} = w_d z_d + b_d$$

e_i = a vector with a one in the i th position and zeros everywhere else.
A mixed integer programming formulation (MILP) is an extremely well-studied area for small problems.

Linearization: we can express $z_{i+1} = \max\{0, w_i z_i + b_i\}$ using linear constraints

If l_i, u_i are lower bound and upper bound for $w_i z_i + b_i$

$$z_{i+1} = \max\{0, w_i z_i + b_i\} \Rightarrow z_{i+1} \geq w_i z_i + b_i$$

$$z_{i+1} \geq 0 \Rightarrow \begin{cases} z_{i+1} = w_i z_i + b_i & \text{if } w_i z_i + b_i > 0 \\ z_{i+1} = 0 & \text{if } w_i z_i + b_i = 0 \end{cases}$$

$$w_i z_i + b_i \geq z_{i+1} + (1 - z_i) l_i$$

$$z_i \in \{0, 1\}^{V_i}$$

ReLU operation

we choose a set of much looser bounds for l_i and u_i

For a single entry $(wz+b)_i = \sum_j w_{ij} z_j + b_i$

if we want to minimize

$$\sum_j w_{ij} z_j + b_i \Rightarrow \begin{cases} z_j = \hat{u}_j & \text{if } w_{ij} < 0 \\ z_j = \hat{l}_j & \text{if } w_{ij} > 0 \end{cases}$$

to maximize vice versa

$$\sum_j w_{ij} z_j + b_i \Rightarrow \begin{cases} z_j = \hat{u}_j & \text{if } w_{ij} < 0 \\ z_j = \hat{l}_j & \text{if } w_{ij} > 0 \end{cases}$$

$$\Rightarrow \max\{w, 0\} \hat{l} + \min\{w, 0\} \hat{u} + b \leq (wz+b) \leq \max\{w, 0\} \hat{u} + \min\{w, 0\} \hat{l}$$

if we begin with $\hat{l} = \max\{x - \epsilon, 0\}$, $\hat{u} = \min\{x + \epsilon, 0\}$

These bounds tell us that for a perturbation delta of ϵ , the logit corresponding is somewhere in the range $[\hat{l}, \hat{u}]$

A final integer Programming Formulation

$$\underset{z_1, \dots, z_{d+1}, \tau_1, \dots, \tau_d}{\text{minimize}} (e_y - e_{y+\text{target}})^T z_{d+1}$$

$$\text{subject to } z_{i+1} \geq w_i z_i + b_i, i=1, \dots, d-1$$

$$z_{i+1} \geq 0, i=1, \dots, d-1$$

$$w_i \cdot z_i \geq z_{i+1}, i=1, \dots, d-1$$

$$w_i z_i + b_i \geq z_{i+1} + (1 - \tau_i) \ell_i, i=1, \dots, d-1$$

$$z_i \in \{0, 1\}^{128}, i=1, \dots, d-1$$

$$z_1 \leq x + \delta$$

$$z_1 \geq x - \delta$$

$$z_{d+1} = w_d z_d + b_d$$

We can solve this problem using the `cvxpy` library, by Eurobi-solver when the resulting objective is negative, we are able to find a perturbation that makes the class logit for the target class is larger than the class logit for the original class.

Certifying robustness

If we want to determine, exactly, whether any adversarial examples exists for a given example, we can run the integer Programming Solution. If any of these optimization objectives have a negative solution, then there exists an adversarial example. In contrast, if none of the optimization objectives is negative for any target class, then the classifier has been certified to be robust on this example.

Upper bounding the inner maximization (Convex relaxations)

If we want to provide formal guarantees of robustness, it is important to be able to obtain fast upper bounds on the inner maximization Problem.

Method 1: $0 \leq \tau_i \leq 1$

either zero or one.

to relax the constraint that each element of τ_i must be

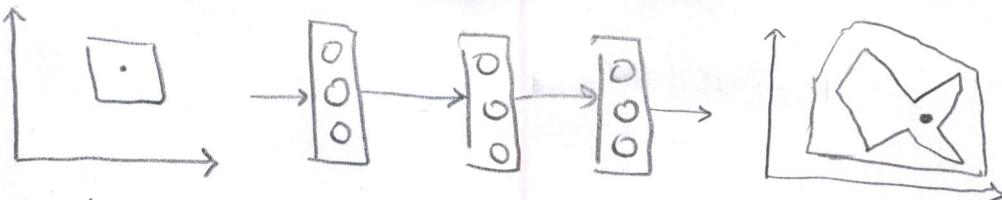
Solving the relaxed version of
a targeted attack problem
and objective
is still positive



this attack
will not work

adv examples
will not exist
in the original
set

relaxed set
is larger than the
original set



Input x and
allowable perturbations

Deep Network

Final layer and
Convex outer bound on the
adversarial polytope

If no adversarial example exists in the larger set, then none exists in the true adversarial polytope either.

✓ Faster solutions of the convex relaxation

A procedure known as Conven duality can compute a provable lower bound on the relaxed objective.

Interval-Propagation-based bounds

Using bound propagation techniques with more complex network is more worthwhile than to use a simpler network and the tighter but more costly bounds based upon conven relaxations

If we have an interval bound on the second-to-last layer $\hat{l} \leq z_d \leq \hat{u}$
we can simply solve the optimization problem

$$\begin{aligned} & \underset{z_d, z_{d+1}}{\text{minimize}} \quad c^T z_{d+1} \\ & \equiv \\ & \text{subject to } z_{d+1} = w_d^T z_d + b_d \quad \hat{l} \leq z_d \leq \hat{u} \\ & \quad \hat{l} < z_d < \hat{u} \end{aligned}$$

$$\text{The result: } \max\{c^T w_d c, 0\} \hat{l} + \min\{c^T w_d c, 0\} \hat{u} + c^T b_d$$