

1 Adversarial Attacks

Targeted FGSM: $x' = x - \varepsilon \text{sgn}(\nabla_x \mathcal{L}_t(x))$, where t is the target label.

Untargeted FGSM: $x' = x + \varepsilon \text{sgn}(\nabla_x \mathcal{L}_y(x))$, where y is the original label.

Carlini-Wagner:

Solve $\text{argmin}_\eta \|\eta\|_p + \lambda \mathcal{L}_t(x + \eta)$ s.t. $x + \eta \in [0, 1]^n$. \mathcal{L} satisfies $\mathcal{L}_t(x + \eta) \leq 0 \Rightarrow f(x + \eta) = t$. Optimized with LBFGS-B.

PGD: Iterate FGSM, each time with a small step. Typically used for adversarial training.

Adversarial Accuracy: Correct classification and no adversarial example.

Defense: $\min_\theta \mathbb{E}_{x,y} [\max_{x' \in S(x)} \mathcal{L}(x', y; \theta)]$.

Norms:

$$\|x\|_\infty \leq \|x\|_1 \leq d\|x\|_\infty$$

$$\|x\|_\infty \leq \|x\|_2 \leq \sqrt{d}\|x\|_\infty$$

$$\|x\|_2 \leq \|x\|_1 \leq \sqrt{d}\|x\|_2$$

$$\text{Hölder: } \|fg\|_1 \leq \|f\|_p \|g\|_q \text{ if } \frac{1}{p} + \frac{1}{q} = 1$$

2 Certification

Goal: given a neural network N , a constraint over input ϕ and a property over outputs ψ , prove that $i \models \phi \Rightarrow N(i) \models \psi$ or return a violation. We want sound, complete and scalable algorithms, i.e., to either scale sound and complete methods or to make sound but incomplete methods more complete.

- Soundness: If the property does not hold, the program always states the property does not hold.
- Completeness: If the property holds, the program always states the property holds.

Convex Relaxations

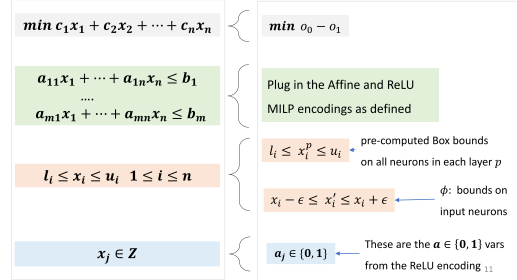
Incomplete but sound

- Box. Use hypercubes as approximation.
- DeepPoly. Bound every variable by $\sum_j w_j^l x_j + v^l \leq x_i \leq \sum_j w_j^u x_j + v_j$ and $l_i \leq x_i \leq u_i$. First do a forward propagation and then a backward propagation. (1) For affine layer $y = Wx + b$, use $Wx + b \leq y \leq Wx + b$. (2) For ReLU layers, if it is strictly positive (negative), use $x \leq y \leq x$ ($0 \leq y \leq 0$). Otherwise, apply min-area heuristic to choose one from $0 \leq y \leq \lambda(x - l)$ and $x \leq y \leq \lambda(x - l)$ (where $\lambda = \frac{u}{u-l}$), i.e., choose the first one if $u \leq -l$ o.w. the second one. When forward-propagating, we compute the interval bound by backsubstituting to the first layer.

Relaxations of methods not included in one another cannot be compared.

MILP

Generic MILP problem

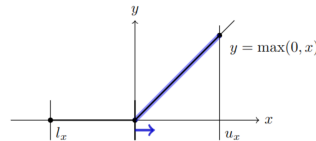


Complete and sound for ReLU network, but is NP-complete.

The encoding for **affine** layer is: $Wx + b \leq y \leq Wx + b$. The encoding for **ReLU** is: $x \leq y \leq x - l(1 - a)$, $y \leq ua$, $y \geq 0$ and $a \in \{0, 1\}$. The variable a serves as a “state” indicator to encode a piece-wise linear function. The bounds l, u are pre-computed by Box. MILP benefits from Box bounds in that it could stop further exploration of infeasible combinations of integer variables, e.g., if the objective lower bound is positive when $a_1 = 0$, then there is no need to explore a_2 for $a_1 = 0$.

Branch and Bound

Can split variables in cases (e.g. ReLU output can be positive or negative) and certify both splits separately. Constraints can be enforced using KKT condition.



$$\begin{aligned} \max_{x \in \mathcal{X}} a^\top x + c &\leq \max_{x \in \mathcal{X}} \min_{\beta \geq 0} a^\top x + c + \beta x_i \\ \text{s.t. } x_i &\geq 0 \end{aligned}$$

Δ Adjust output bounds when splitting.

Initialize queue (without splits). While queue is not empty and not timed out: 1. Get sub-problem from queue. 2. Compute bound of interest. 3. If not verified, pick neuron to split according to heuristic. 4. Add both new sub-problems to queue.

3 Certified Training

Goal: $\min_\theta \mathbb{E}_{(x,y) \sim D} [\max_{x' \in B_\varepsilon(x)} \mathcal{L}(\theta, x', y)]$

DiffAI / IBP

IBP (over)approximates inner max by propagating the box $B_\varepsilon(x)$ and computing the max loss over the output set. Can use any abstract transformer (Box, DeepPoly, etc.) to get bounds on logits.

SABR

SABR uses the loss

$$\mathcal{L}_{\varepsilon, \tau}(x, y) = \max_{x' \in B_{\varepsilon - \tau}(x)} \max_{z \in B_\tau(x')} CE(f(z), y),$$

i.e. it uses a box $B_\tau(x')$ of size $\tau \leq \varepsilon$ around an adversarial sample x' . The outer max is (under)approximated by PGD search and the max inside is (over)approximated using IBP.

4 Randomized Smoothing

Given a classifier f and $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$, define a smoothed classifier g as

$$g(x) = \text{argmax}_c \mathbb{P}[f(x + \epsilon) = c].$$

THM: Suppose $\exists p_A, p_B \in [0, 1]$, s.t. $\mathbb{P}[f(x + \epsilon) = c_A] \geq p_A \geq p_B \geq \max_{c \neq c_A} \mathbb{P}[f(x + \epsilon) = c]$. Then $g(x + \delta) = c_A$ for all $\|\delta\|_2 < R$, where the radius $R = \frac{\sigma}{2}(\Phi^{-1}(p_A) - \Phi^{-1}(p_B))$ and Φ^{-1} is the Gaussian quantile function.

If we have $p_A > \frac{1}{2}$, then $p_B \leq 1 - p_A < p_A$ and thus $R > \sigma \Phi^{-1}(p_A)$. To compute probability, we use sampling. To avoid selection bias, we first sample to get the top label and then use an independent sampling to estimate the probability:

```
function CERTIFY(f, σ, x, n0, n, α)
  counts0 ← SAMPLEUNDERNOISE(f, x, n0, σ)
  ĉ_A ← top index in counts0
  counts ← SAMPLEUNDERNOISE(f, x, n, σ)
  p_A ← LOWERCONFBOUND(counts[ĉ_A], n, 1 - α)
  if p_A > 1/2 return prediction ĉ_A and radius σ Φ⁻¹(p_A)
  else return ABSTAIN
```

By the above theorem, we have $g(x) = \hat{c}_A$ for all $\|\delta\|_2 < R$ with probability $1 - \alpha$ if \hat{c}_A, R is returned.

For inference, to compute $g(x)$, we use:

```
function PREDICT(f, σ, x, n, α)
  counts ← SAMPLEUNDERNOISE(f, x, n, σ)
  ĉ_A, ĉ_B ← top two indices in counts
  n_A, n_B ← counts[ĉ_A], counts[ĉ_B]
  if BINOMPVVALUE(n_A, n_A + n_B, 0.5) ≤ α return ĉ_A
  else return ABSTAIN
```

The probability of wrong prediction is $\mathbb{P}[\hat{c}_A \neq c_A, \text{no abstain}] = \mathbb{P}[\text{no abstain} | \hat{c}_A \neq c_A] \mathbb{P}[\hat{c}_A \neq c_A] \leq \alpha$.

Pros: (1) it can scale to large networks and (2) is model-agnostic.

Cons: (1) it requires sampling at inference time and (2) many samples could be needed to increase the certified radius.

5 Federated Learning

Communication Round: (1) Server selects clients to participate. (2) Clients do local training update on private data and send information to the server. (3) Server aggregates updates into updated model.

FedSGD

Clients compute gradient wrt. (minibatch of) local data on global weights. Server takes average for global update. Pros: convergence guarantee. Cons: requires lots of communication. Gradient inversion:

$$\text{argmin}_{x^*} d(\nabla_{\theta} \mathcal{L}(f_{\theta}(x^*), y^*), g_k) + \alpha_{\text{reg}} \cdot \mathcal{R}(x^*)$$

Distance between reconstructed gradient on f_{θ} and the true gradient g_k Domain specific prior

FedAvg

Do several iterations of local SGD on the global weights and average the local weights as the global weight. Pros: Less communication. Cons: No convergence guarantee.

6 Privacy

Privacy Attacks

- Model stealing (white/black box)
- Model inversion: Client queries the model to find representative samples.
- Data extraction: Client queries the model to find exact training samples.
- Membership inference: Client wants to determine if the data point was used to train the model or not (train shadow models)

Differential Privacy

A randomized algorithm M is ϵ, δ -DP if for all "neighbouring" datasets D, D' and all attacks (sets) S ,

$$\mathbb{P}[M(D) \in S] \leq e^\epsilon \mathbb{P}[M(D') \in S] + \delta.$$

DP Rules

Post-processing: If M is ϵ, δ -DP, then $g \circ M$ is ϵ, δ -DP.

Composition: If M_1, M_2 are ϵ_1, δ_1 -DP and ϵ_2, δ_2 -DP, then $M := (M_1, M_2)$ is $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -DP.

Privacy Amplification: Applying an ϵ, δ -DP

mechanism on a random fraction $q \leq 1$ subset yields a $\tilde{q}\epsilon, q\delta$ -DP mechanism, where $\tilde{q} \approx q$.

Laplace Mechanism

Let $\Delta_1 = \max_{D \sim \text{neigh } D'} \|f(D) - f(D')\|_1$. Then $f + \text{Laplace}(0, \Delta_1/\epsilon)$ is ϵ -DP.

Gaussian Mechanism

Let $\Delta_2 = \max_{D \sim \text{neigh } D'} \|f(D) - f(D')\|_2$. Then if $\sigma^2 = \frac{2\log(1.25)}{\delta\epsilon^2} \Delta_2^2$ we get that $f + \mathcal{N}(0, \sigma^2 I)$ is ϵ, δ -DP.

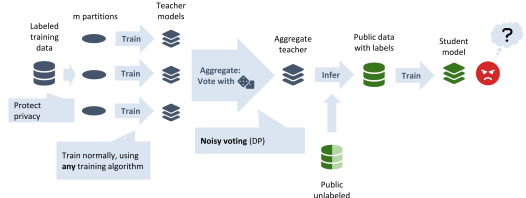
DP-SGD

ϵ, δ -DP adaptation of SGD.

1. Clip gradients: $g'_i = \min\left(\frac{R}{\|g\|_2}, 1\right) g_i$
2. Aggregate: $\bar{g} = \frac{1}{n} \sum_{i=1}^n g'_i$
3. Add noise: $\hat{g} = \bar{g} + \mathcal{N}(0, \sigma^2 I)$ where $\sigma^2 = \frac{2\log(1.25)}{\delta\epsilon^2} \frac{R^2}{n^2}$
4. Update: $\theta \leftarrow \theta - \eta \hat{g}$

DP-SGD is $O\left(q\epsilon\sqrt{T\log\frac{1}{\delta}}\right), O(qT\delta)$ -DP. It is $O(q\epsilon\sqrt{T}), \delta$ -DP when allowing data-dependent σ .

PATE



Noisy voting: Add $\text{Laplace}(0, \frac{2}{\epsilon})$ -noise to each class count $n_c(x)$ before argmax (n_c has sensitivity 2). Results in $N\epsilon$ -DP, where N is the number of classifications we compute on the public unlabeled data.

7 Synthetic Data

Goal: Dataset with same statistical properties & provable DP guarantees

Marginal: #Elements per attribute subset split by attribute domains (L_2 -sensitivity: 1)

Exponential Mechanism: $\mathbb{P}[\mathcal{M}(\mathcal{D}) = r] \propto \exp(\epsilon \cdot q(\mathcal{D}, r))$, q : quality score

ProgSyn: 1. Sample random noise $z \sim \mathcal{N}(0, I_p)$. 2. pass z through generative model $\Theta(z)$. 3. Obtain synthetic dataset $g_\Theta(z)$. 4. Adapt params $\Theta \rightarrow g_\Theta(z)$ close to X . 5. Fine-tune params $\Theta \rightarrow g_\Theta(z)$ satisfy additional constraints

8 Logical Properties

DL2

Idea: Translate properties ϕ into loss function on T and optimize. For example, ϕ could be: $\|z - x\|_\infty < \epsilon \Rightarrow \|NN(z) - NN(x)\|_\infty < \delta$.

ϕ	$T(\phi)$
$t_1 \leq t_2$	$\max(0, t_1 - t_2)$
$t_1 \neq t_2$	$\mathbb{1}_{\{t_1=t_2\}}$
$t_1 = t_2$	$T(t_1 \leq t_2 \wedge t_2 \leq t_1)$
$t_1 < t_2$	$T(t_1 \leq t_2 \wedge t_1 \neq t_2)$
$\phi \vee \psi$	$T(\phi)T(\psi)$
$\phi \wedge \psi$	$T(\phi) + T(\psi)$

$\forall x: T(\phi)(x) = 0 \Leftrightarrow x$ satisfies ϕ .

Training for Logical Properties

Idea: Get the worst violation of ϕ and find θ that minimizes its effect. I.e. solve $\min_\theta \mathbb{E}_{s \sim D} [T(\phi)(z^*, s, \theta)]$ where $z^* \in \text{argmin}_z T(\neg\phi)(z, s, \theta)$.

In practice, restrict z to a convex set with efficient projections (closed form). One can then remove the constraint from ϕ that restricts z on the convex set and do PGD while projecting z onto the convex set.

9 Individual Fairness

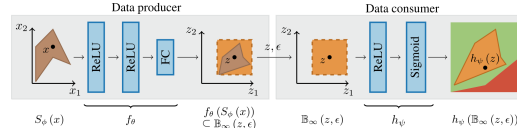
DEF: A model M is individually fair if for any similar x, x' it produces similar outputs $M(x), M(x')$, e.g.

- Lipschitz continuity: $\mu(M(x), M(x')) \leq L\phi(x, x')$.
- Binary similarity: $\phi(x, x') \Rightarrow M(x) = M(x')$.

Fair Representation Learning:

- Data Regulator: Determines fairness criteria and data source(s), audits results.
- Data Producer: Computes the fair representation given data regulator criteria.
- Data Consumer: Trains the ML model given the sanitized data. Split models into an encoder f_θ (data producer) and a predictor h_ψ (data consumer).

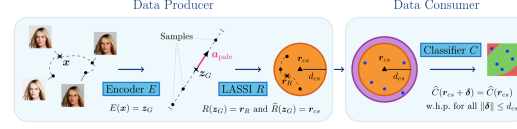
LCIFR



Consider binary ϕ, μ . 1. Encoder f_θ is trained s.t. $\phi(x, x') \Rightarrow \|f(x) - f(x')\|_\infty \leq \delta$ (with a combined loss ensuring informativeness of the representations) using DL2 (or δ is computed when encoded as MILP). 2. Use L_∞ -

robustness training on encoded data for perturbations of magnitude δ and certify robustness. Gives certified individual fairness. Problem: Relies on NN verifiers so not scalable to real-world models.

LASSI



1. Data producer uses center smoothing to get an encoder that provably maps all similar points close together with high probability.
2. Data consumer uses randomized smoothing to certify that all points within a certain radius get classified the same with high probability. Problem: Depends on quality of generative model.

10 Group Fairness

Demographic Parity

$$\mathbb{P}[\hat{Y} | G = 0] = \mathbb{P}[\hat{Y} | G = 1]$$

Equalized Odds

$$\mathbb{P}[\hat{Y} | G = 0, Y = 0] = \mathbb{P}[\hat{Y} | G = 1, Y = 0]$$

$$\mathbb{P}[\hat{Y} | G = 0, Y = 1] = \mathbb{P}[\hat{Y} | G = 1, Y = 1]$$

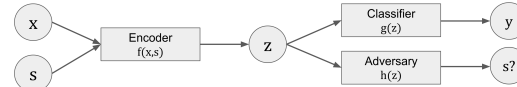
Equal Opportunity

$$\mathbb{P}[\hat{Y} | G = 0, Y = 1] = \mathbb{P}[\hat{Y} | G = 1, Y = 1]$$

Approaches:

- Pre-processing: Debias the data, such that standard training yields a fair model. E.g. fair representation learning.
- In-training: Change the training pipeline to learn a fair model (on biased data). E.g. add soft fairness constraint to the loss.
- Post-processing: Modify predictions at inference time. E.g. different thresholds for different groups.

LAFTR



Jointly trains an encoder f , classifier g , and adversarial classifiers h by solving $\min_{f, g} \max_{h \in \mathcal{H}} \mathcal{L}_{\text{clf}}(f, g) - \lambda \mathcal{L}_{\text{adv}}(f, h)$.

Balanced Accuracy

$$\text{BA}(h) = \frac{1}{2} \mathbb{E}_{Z_0} [1 - h(z)] + \frac{1}{2} \mathbb{E}_{Z_1} [h(z)]$$

$$h^* = \text{argmax}_h \text{BA}_{Z_0, Z_1}(h) = \mathbb{1}_{\{p_1(z) \geq p_0(z)\}}$$

DP-Distance (Relaxation of DP)

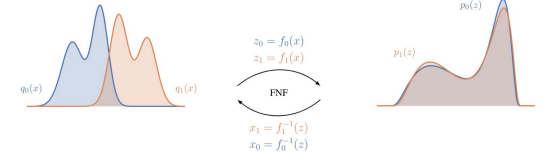
$$\Delta_{Z_0, Z_1}^{\text{DP}}(g) = |\mathbb{E}_{Z_0} g(z) - \mathbb{E}_{Z_1} g(z)|$$

$$\Delta_{Z_0, Z_1}^{\text{DP}}(g) \leq 2 \text{BA}_{Z_0, Z_1}(h^*) - 1$$

LAFTR estimates $\Delta_{Z_0, Z_1}^{\text{DP}}(g)$ given an adversary h . Fairness is overestimated (there may be better adversaries).

Fair Normalizing Flows

Normalizing flows for a distribution $p(z)$ parameterize $z = f(x)$ where x is sampled from a simple distribution $q(x)$ through an invertible transformation: $p(z) = q(f^{-1}(z)) |\det \partial_z f^{-1}(z)|$.



FNF learns two flows f_0, f_1 as encoders for Z_0, Z_1 . 1. Estimate densities q_0, q_1 for both groups. 2. Apply the encoder to get $z = f_0(x)$ (or $f_1(x)$), and use the flows and q_0, q_1 to get $p_0(z)$ and $p_1(z)$. 3. Use $p_0(z)$ and $p_1(z)$ to estimate the optimal adversary h^* and then upper bound its balanced accuracy with probability $1 - \epsilon$ (Hoeffding). 4. Use the inequality from LAFTR above to upper bound the DP-distance of any downstream classifier trained on representations z with probability $1 - \epsilon$. To get tight bounds: train the flows to promote low accuracy of h^* by minimizing $D_{KL}(p_0 || p_1) + D_{KL}(p_1 || p_0)$. \triangle The guarantees hold for the estimated q_0, q_1 .

FARE

Restrict representations z to a finite set. This allows bounding $\text{BA}_{Z_0, Z_1}(h^*)$ whp. directly. Use fairness-aware decision tree as encoder which is obtained by training a tree to have unbalanced leaves wrt. y and uniform leaves wrt. s . Gives provable unfairness upper bound with no restrictive assumptions

11 Miscellaneous

- Laplace: $f(x) = \frac{1}{2\sigma} e^{-\frac{|x-\mu|}{\sigma}}$
- Normal: $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$
- $x \Rightarrow y$ is equivalent to $\neg x \vee y$
- DP: noise $\oplus \rightarrow$ privacy \oplus & utility \ominus
- Carlini-Wagner: $\|\cdot\|_\infty \rightarrow \sum_i \max(\eta_i - \tau, 0)$
- Adversarial Accuracy: $\frac{\# \text{correct} - \# \text{adv_examples}}{\# \text{samples}}$
- Cross entropy loss: $-\sum_i \sum_c y_{i,c} \log(p_{i,c})$