

Adversarial Attacks (White-Box)

Targeted FGSM: $f(x - \epsilon \cdot \text{sgn}(\nabla_x \text{loss}_t(x)))$ approaches t
Untargeted FGSM: $f(x + \epsilon \cdot \text{sgn}(\nabla_x \text{loss}_s(x)))$ strays from S
 \rightarrow both guarantee $x' \in [x - \epsilon, x + \epsilon]$ box
 $\min_{\eta} \|x\|_p + c \cdot \text{obj}_t(x + \eta)$ s.t. $x + \eta \in [0, 1]^n$
 with $\text{obj}_t(x + \eta) \leq 0$ if $f(x + \eta) = t$, e.g. $\text{obj}_t(x) = \max(0, 0.5 - P(f(x) = t))$
 CW: Use LBFGS-B optimizer or $\text{obj}_t(x) = -\log_2(P(f(x) = t))$
 PGD: Repeat FGSM, projecting to $[x - \epsilon, x + \epsilon]$

Adversarial Defenses

Adversarial Accuracy = $\frac{\text{correctly classified, even after AGO-attack}}{\text{tested samples}}$
 Training aims to solve $\min_{\theta} \mathbb{E}_{(x,y) \sim D} [\max_{x' \in S(x)} L(\theta, x', y)]$
 Approximate inner maximization with PGD
 GCG attacks on LLMs: Use "Sure" as next-token target x_{n+2} .
 $x_{n+1} := \arg\min_{x_{n+1} \in S} L(x_{n+1})$ for $S := \arg\text{TopK}(\nabla_{\text{model}(x_{n+1})} L(x_{n+1}))$

NN-Certification

Rice's Theorem: no sound & complete general certification algorithms
Box ($O(n^2L)$): $\# [a, b] = [b, -a]$, $\text{ReLU}^*[a, b] = [\text{ReLU}(a), \text{ReLU}(b)]$,
 $[a, b] \# [c, d] = [\min(ac, ad, bc, bd), \max(ac, ad, bc, bd)]$, $+$ and \wedge trivial
DeepPoly ($O(n^3L^2)$):
 Affine: encode exactly, being careful with signs for b_y and u_y !!
 ReLU: if $u_x \leq 0$: $0 \leq y \leq 0$, $b_y = 0$, $u_y = 0$
 if $b_x > 0$: $x \leq y \leq x$, $b_y = b_x$, $u_y = u_x$
 otherwise: $\lambda := \frac{u_x}{u_x - b_x}$, $\lambda x \leq y \leq \lambda(x - b_x)$, $b_y = 0$, $u_y = u_x$
 with $\lambda \in [0, 1]$, min. area says $d=0$ iff $u_x \leq -b_x$
 Ψ : encode using single output neuron σ , proving $b_0 > 0$ etc.
Branch and Bound: Split ReLU at $x=0$ and use input constraints to split.
 $(\max_x P(x) \text{ s.t. } g(x) \leq 0) \leq \max_x \min_{\beta} (f(x) - \beta g(x))$. $\frac{\min_{\beta} \max_x f(x) - \beta g(x)}{\max_{\beta} \min_x f(x) - \beta g(x)}$
MILP (NP-complete):
 Affine: $Wx + b \leq y \leq Wx + a$
 ReLU: $y \geq 0$, $y \leq x$, $y \leq u_x a$, $y \leq x - b_x(1-a)$, $a \in \{0, 1\}$
 Φ : $x_i - \epsilon \leq x'_i \leq x_i + \epsilon$, Ψ : objective $\min c_0 - a_1$
 Box-constraints (to accelerate): $\ell_i \leq x_i^p \leq u_i$
PRIMA: Abstract Neurons jointly, computing convex hull via dual problem of intersecting halfspaces (under-approximate)

Certified Defenses

Training aims to solve $\min_{\theta} \mathbb{E}_{(x,y) \sim D} [\max_{z \in \gamma(MN^*(S(x)))} L(\theta, z, y)]$
 with $L(z, y) = \max_{q \neq y} (z_q - z_y)$, compute $\max_c (\max(\text{box}(z_c - \# z_y))$
 $L(z, y) = CE(z, y)$, compute $CE(\text{softmax}([u_1, \dots, u_y, \dots, u_n]), y)$
COLT: Only compute symbolically up to n th layer, then PGD.
 Problem: Projection onto DeepPoly shape is hard.

Randomized Smoothing for Robustness

Given classifier f , make $g(x) := \arg\max_{\hat{c}} P_{\epsilon}[f(x + \epsilon) = \hat{c}]$ with $\epsilon \sim \mathcal{N}(0, \sigma^2 I)$ and $P_{\epsilon}(x) := P_{\epsilon}[f(x + \epsilon) = g(x)]$, $P_{\epsilon}(x) := \max_{\hat{c}} \dots$
Robustness guarantee: $P_{\epsilon}(x) \geq P_{\text{acc}} \Rightarrow P_{\text{acc}} \geq P_{\epsilon}(x) \Rightarrow g(x + \delta) = c_A$ for all $\|\delta\|_2 \leq R_x := \frac{1}{\sqrt{2}} (\Phi^{-1}(P_{\text{acc}}) - \Phi^{-1}(P_{\text{acc}}))$ certification radius
 If $P_{\text{acc}} \geq \frac{1}{2}$, then $R_x \geq \frac{1}{\sqrt{2}} \Phi^{-1}(P_{\text{acc}})$, allows for efficient certification.
 \rightarrow Guess \hat{c}_A via Monte-Carlo Integration $P_{\epsilon}(x) = \int_{\epsilon} \mathbb{I}_{f(x+\epsilon)=\hat{c}_A} \text{PDF}_{\mathcal{N}(0, \sigma^2 I)}(d\epsilon)$
 Repeat Monte-Carlo with large n , counting k hits
 Estimate $P_{\text{acc}} := \frac{k}{n} \binom{n}{k}^{-1}$ via Copper-Pearson
 Then, with probability $> 1 - \delta$, $g(x + \delta) = \hat{c}_A$ for all $\delta \leq R_x$
 If ABSTAIN, true $P_{\text{acc}} \leq 0.5$ OR guess was wrong OR lower bound too loose.
 At inference, get top-2-classes via Monte-Carlo Integration and do Binomial P-value test on $n_A, n_A + n_B, \dots, 0.5$.

Privacy

Attacks: Model Stealing, Inversion (exact/representative), Membership Inference
Black-Box ML: Train n models on the same data distribution, some with some without x . If logits given, train classifier mapping logits to x -membership. Otherwise, train classifier mapping adversarial robustness to x -membership.
Federated Learning (FedSGD):
 Client k computes $g_k := \nabla_{\theta} L(\theta, x_k, y_k)$ for minibatch $(x_k, y_k) \sim D_k$
 Server updates model $\theta_{t+1} := \theta_t - \gamma g_k$
 Attack: batch size 1 \Rightarrow exact reconstruction by gradient inversion for piecewise-linear NN (ReLU-based networks)
 batch size $> 1 \Rightarrow$ reconstruct linear combination of inputs
 $\arg\min_{x^*} \text{dist}(\nabla_{\theta} L(\theta, x^*, y^*), g_k) + \text{dreg } R(x^*)$
 Regularization prior, e.g. text perplexity, image variation
FedAVG:
 client runs E epochs of SGD. Server updates $\theta_{t+1} := \frac{1}{k} \sum_{k=1}^k \theta_{t,k}$.
 Attack: Simulate client to θ_k , compute $\text{dist}(\theta_k, \theta_k)$.
 $\arg\min_{x^*} \text{dist}(\theta_k, \theta_k) + \text{dreg } \frac{1}{E} \sum_{k=1}^k R(g_k(\tilde{x}_{k,1}, \tilde{x}_{k,2}, \tilde{x}_{k,3}))$
 Regularization over average distance between average images of each pair

Differential Privacy

Mechanism M is ϵ -DP iff $\forall (a, a') \in \text{Neigh}$. $\forall S$. $P(a) := P[M(a) \in S] \leq e^{\epsilon} P(a')$
Laplace Mechanism: $P(a) + \text{Lap}(0, \frac{\Delta_1}{\epsilon})$ with $\Delta_1 := \max_{a' \in \text{Neigh}} \|f(a) - f(a')\|_1$
 M is (ϵ, δ) -DP iff $\forall (a, a') \in \text{Neigh}$. $\forall S$. $P(a) \leq e^{\epsilon} P(a') + \delta$
Gaussian Mechanism: $P(a) + \mathcal{N}(0, \sigma^2 I)$ with $\sigma := \frac{\Delta_2}{\epsilon}$ $\Delta_2 := \sqrt{2 \log(1.25)/\delta}$
 \rightarrow useful if $\exists (a, a') \in \text{Neigh}$. $P(a) = 0 \wedge P(a') \neq 0$
 M_1 is (ϵ_1, δ_1) -DP and M_2 is (ϵ_2, δ_2) -DP $\Rightarrow (M_1, M_2)$ and $M_1 \text{ OR } M_2$ are $(\epsilon_1 + \epsilon_2, \delta_1 + \delta_2)$ -DP
 M is (ϵ, δ) -DP $\Rightarrow P \circ M$ is (ϵ, δ) -DP
 M_i is (ϵ_i, δ_i) -DP $\Rightarrow \prod_i M_i(a_i)$ is $(\sum \epsilon_i, \sum \delta_i)$ -DP

DP-SGD: Project within-batch gradients onto L_2 -ball of radius C . For batch size L , add $\mathcal{N}(0, \sigma^2 I)$ to aggregate gradient, where $\sigma = \sqrt{2 \log(1.25)/\delta} \frac{C}{\epsilon}$.

Privacy Amplification: Applying an (ϵ, δ) -DP mechanism on a random fraction $q = \frac{1}{N}$ of data yields a $(\tilde{\epsilon}, \tilde{\delta})$ -DP mechanism, where $\tilde{\epsilon} \propto \epsilon q$.

PATE: Split data and train teachers T . Label public dataset by noisy voting: $\arg\max_j |\{t(x) = j \mid t \in T\}| + \text{Lap}(0, \frac{2}{\epsilon})$.
 Method is $(\epsilon, 0)$ -DP for one label, $(\epsilon T, 0)$ -DP for T labels.

AI Regulation

Fairness, Explainability, Data Minimization, Unlearning, Copyright, Emergency Response.

Private Synthetic Data

Select Marginal Queries, Measure using DP, Generate Data Marginal on attributes C is $M_C(\mathcal{D}) := \mu \in \mathbb{R}^{|C|}$ where $\mu_i = \sum_{x \in \mathcal{D}} \mathbb{I}_{x_C = i}$. 1-way marginals are histograms. $\Delta_2 = 1$.
 Mutual Information is $I(X; Y) = \sum_{x,y} \frac{P(x,y)}{P(x)P(y)}$ and is used as edge weight in Chow-Liu algorithm. MST is the optimal 2nd-order approximation. Compute via belief propagation.

Logic in Deep Learning

Test satisfiability of logical formulas over an NN's output.
 $(\text{class}(NN(i)) = c^k) \equiv \bigwedge_{c \neq c^k} NN(i)[c] < NN(i)[c^k]$
 SAT solvers time out on large networks, so translate:
 For all x , $T(\phi)(x) = 0 \Leftrightarrow x \models \phi$ with differentiable T :
 $t_1 \leq t_2 \Leftrightarrow \max(0, t_1 - t_2) \quad t_1 \neq t_2 \Leftrightarrow \mathbb{I}_{t_1 = t_2}$
 $t_1 = t_2 \Leftrightarrow T(t_1 \leq t_2 \wedge t_2 \leq t_1) \quad \psi \vee \phi \Leftrightarrow T(\psi) \cdot T(\phi)$
 $t_1 < t_2 \Leftrightarrow T(t_1 \leq t_2 \wedge t_1 \neq t_2) \quad \psi \wedge \phi \Leftrightarrow T(\psi) + T(\phi)$
 Translate negation with deMorgan $\neg(\psi \wedge \phi) \Rightarrow \neg\psi \vee \neg\phi$.
 By construction, $T(\phi)(x) \geq 0$ for all ϕ, x .
 Box constraints hard to optimize, so use L-BFGS-B solver.
 Training with logic as maximization
 $\max_{\theta} \mathbb{E}_{S \sim D} [\sum_z \phi(z, s, \theta)]$
 Generalized adversarial training beyond robustness
 $\min_{\theta} \mathbb{E}_{S \sim D} [T(\phi)(\arg\min_z T(\gamma \phi(z, s, \theta))), s, \theta]$
 Restrict z to a convex set with efficient projections.

Fairness

Fairness through unawareness (algorithm does not explicitly use sensitive data) does not work.

Individual Fairness: For every $x, y \in \mathcal{X}$, $D(M(x), M(y)) \leq d(x, y)$ where (D, d) are two similarity metrics. M is (d, d) -Lipschitz. The challenge is finding suitable d and D .

Fairness as Robustness:

For $R: \mathbb{R}^d \rightarrow [0, 1]$, $x \mapsto \mathbb{E}_{y \sim \mathcal{Y}(x)} [E[R(x+y)]]$ is 1-Lipschitz. Choosing $d(x, x') := (x - x')^T S (x - x')$ for symmetric positive definite covariance matrix S and $D(M(x), M(x')) := \mathbb{I}_{M(x) \neq M(x')}$ allows reformulation of Lipschitzness for all $\| \delta \|_S < \frac{1}{L}$, $M(x) = M(x + \delta)$, where $\|x\|_S := \sqrt{x^T S x}$ (Mahalanobis dist.).

Fair Representation Learning:

Data Regulator: Define fairness and data sources

Data Producer: Generate fair encoder $f_\theta: \mathbb{R}^n \rightarrow \mathbb{R}^k$.

Data Consumer: Use encoded data to build consumer

Motivation comes at the cost of no certification and unclear fairness-performance tradeoff.

LCIFR:

Define D and d with logic accepted by MLPr or DL2, e.g.

$$d(x, x') = \bigwedge_{i: x_i \neq x'_i} (x_i = x'_i) \bigwedge_{j: x_j = x'_j} |x_j - x'_j| \leq d$$

For each x , obtain $S_d(x) = \{x' \mid d(x, x') \leq d\}$.

Train encoder using DL2 s.t. $\forall x' \in S_d(x) \|f_\theta(x) - f_\theta(x')\|_\infty \leq \delta$ (see above), appending general classifier to preserve utility.

Now, consumers must be δ -robust (e.g. by randomized smoothing) to certify end-to-end fairness.

LaSSL:

For high-dimensional data (e.g. images), use semantic feature space from generative model in similarity formulas. Use center smoothing to produce end-to-end fair model with probability $1 - \delta$ -dcs.

Group Fairness:

Demographic parity $\mathbb{P}[\hat{Y}=1 \mid G=0] = \mathbb{P}[\hat{Y}=1 \mid G=1]$

Equal Opportunity $\mathbb{P}[\hat{Y}=1 \mid Y=1, G=0] = \mathbb{P}[\hat{Y}=1 \mid Y=1, G=1]$

Equalized odds and $\mathbb{P}[\hat{Y}=1 \mid Y=0, G=0] = \mathbb{P}[\hat{Y}=1 \mid Y=0, G=1]$

Group Fairness \nRightarrow Individual Fairness

Post-Processing Approach: In a binary classifier, use different threshold depending on sensitive attribute.

In-training approach: Add relaxed fairness constraints solved with DL2, e.g. $-\epsilon \leq \mathbb{P}[\hat{Y}=1 \mid S=0] - \mathbb{P}[\hat{Y}=1 \mid S=1] \leq \epsilon$.

Preprocessing with guarantees: Jointly train encoder f , classifier g , adversary R trying to predict sensitive attributes from latent data.

$$\min_{f, g} \max_R (\mathbb{E}_{x \sim \mathcal{X}} [\mathbb{E}_{y \sim \mathcal{Y}(x)} [f(x, y)]] - \gamma \mathbb{E}_{x \sim \mathcal{X}} [\mathbb{E}_{y \sim \mathcal{Y}(x)} [f(x, y), R]]) \quad (\text{LAIFR})$$

Bounding Unfairness via the optimal adversary:

$$\text{Soft formulation of demographic parity: } \Delta := \left| \mathbb{E}_{z \sim \mathcal{Z}_0} [g(z)] - \mathbb{E}_{z \sim \mathcal{Z}_1} [g(z)] \right|$$

$$\text{Balanced accuracy of adversary: } BA_{\mathcal{Z}_0, \mathcal{Z}_1}(R) := \frac{1}{2} (\mathbb{E}_{z \sim \mathcal{Z}_0} [1 - R(z)] + \mathbb{E}_{z \sim \mathcal{Z}_1} [R(z)]) \\ = \frac{1}{2} \int (\rho_0(z)(1 - R(z)) + \rho_1(z)R(z)) dz$$

$$\Delta_{\mathcal{Z}_0, \mathcal{Z}_1}(g) \leq 2 BA_{\mathcal{Z}_0, \mathcal{Z}_1}(R^*) - 1 \text{ for optimal } R^*(z) := \mathbb{I}_{\rho_1(z) > \rho_0(z)}$$

LAIFR assumes family of $R \sim \mathcal{E} \mathcal{E}$ fairness is overestimated

Fair Normalizing Flows: Estimate densities q_0 and q_1 , learn invertible encoder $z = f(x)$. Find density of the output distributions by $\log p(x) = \log q(f^{-1}(z)) + \log |\det \frac{\partial f^{-1}(z)}{\partial z}|$. Optimize for low KL-divergence between p_0 and p_1 and utility with downstream classifier g . Hoeffding's inequality bounds BA with probability $1 - \epsilon$. FNV guarantees only w.r.t. estimated densities.

Fairness with Restricted Encoders: Restrict latent space to finite set (e.g. k-means clusters), then compute direct bounds on each output.

Basics

$$l_p\text{-norm } \|x\|_p := \left(\sum_i |x_i|^p \right)^{1/p}, \|x\|_\infty := \max_i |x_i|$$

$$\mathcal{B}_E^1 \subseteq \mathcal{B}_E^2 \subseteq \mathcal{B}_E^\infty \subseteq \mathcal{B}_{E, \text{opt}}^2 \subseteq \mathcal{B}_{E, \text{opt}}^1$$

$$\text{CE-loss } CE(x, y) := - \sum_i y_i \log_2(p_f(x)_i)$$

for classification $y := \text{onehot}(t)$, so $\text{loss}_t(x) = -\log_2(p_f(x)_t)$

$$\text{softmax } p_f(x)_t := \frac{\exp(x_t)}{\sum_i \exp(x_i)}$$

$$\text{Hölder: } \|u^T v\|_1 \leq \|u\|_p \|v\|_q \text{ for } \frac{1}{p} + \frac{1}{q} = 1$$

$$\text{Weak Duality: } \max_a \min_b f(a, b) \leq \min_b \max_a f(a, b)$$

$$\text{Subadditivity: } \sqrt{x+y} \leq \sqrt{x} + \sqrt{y}$$

$$\text{Cauchy-Schwarz: } x^T y \leq \|x\|_2 \|y\|_2$$

Project z onto $\mathcal{B}_E^p(x)$

$$p=00: z_i := \max(\min(z_i, x_i + \epsilon), x_i - \epsilon)$$

$$p=2: z' := x + \frac{z - x}{\max(1, \|z - x\|_2)}$$

$p=1$: Soundly approximate using the $p=2$ approach.

Normal CDF:

$$x \sim \mathcal{N}(\mu, \sigma^2): \mathbb{P}(x \leq z) = \Phi(z), \mathbb{P}(x \leq \Phi^{-1}(z)) = z$$

$$x \sim \mathcal{N}(\mu, \sigma^2): \mathbb{P}(x \leq z) = \Phi\left(\frac{z - \mu}{\sigma}\right), \mathbb{P}(x \leq \mu + \sigma \Phi^{-1}(z)) = z$$

$$d\mathcal{N}(\mu_1, \sigma_1^2) + d\mathcal{N}(\mu_2, \sigma_2^2) = d\mathcal{N}(\alpha\mu_1 + \mu_2, \alpha^2\sigma_1^2 + \sigma_2^2)$$

$$\text{Chebyshev: } \mathbb{P}(|x - \mathbb{E}[x]| \geq \epsilon) \leq \frac{\text{Var}[x]}{\epsilon^2}$$

$$\mathbb{E}[X] := \int_{-\infty}^{\infty} x f(x) dx \text{ with PDF } f.$$

$$\text{Var}[X] := \mathbb{E}[(X - \mathbb{E}[X])^2] = \mathbb{E}[X^2] - \mathbb{E}[X]^2$$

$$\mathbb{P}[A|B] := \frac{\mathbb{P}[A \cap B]}{\mathbb{P}[B]}$$

$$\text{Bayes } \mathbb{P}[A|B] = \frac{\mathbb{P}[B|A] \mathbb{P}[A]}{\mathbb{P}[B]}$$

Derivatives:

$$(fg)' = f'g + fg', (f/g)' = (fg - fg')/g^2$$

$$(f \circ g)' = (f' \circ g)g', \ln(x) = \frac{1}{x}$$

$$\nabla_x b^T x = \nabla_x x b^T = b \quad \nabla_x x^T x = \nabla_x \|x\|_2^2 = 2x$$

$$\nabla_x x^T A x = (A^T + A)x \quad \nabla_x \|x - b\|_2 = \frac{x - b}{\|x - b\|_2}$$

$$\nabla_x \|x\|_1 = \frac{x}{|x|} \quad \nabla_x \|Ax - b\|_2^2 = 2(A^T A x - A^T b)$$