

# Security, Privacy, and Fairness Analysis for HW4

Taha Majlesi

Student ID: 810101504

Trustworthy Artificial Intelligence  
University of Tehran

**Abstract**—This report presents a complete, reproducible implementation of HW4 with emphasis on theoretical correctness and empirical interpretability. For security, the real poisoned checkpoint is analyzed via Neural Cleanse, attacked-label detection is performed by lower-tail MAD, and one-epoch unlearning is evaluated by clean accuracy and ASR before/after mitigation. For privacy, Laplace mechanism behavior is derived from first principles and evaluated under base, sequential-composition, and unbounded-adjacency assumptions. For fairness, baseline and assignment-required mitigation are compared with two bonus methods (reweighing and group thresholds), and results are decomposed into both aggregate metrics and group-level behavior. Every value in this report is generated by executable code.

## I. INTRODUCTION

Trustworthy AI is a multi-objective design problem: models should resist adversarial manipulation, leak limited information about individuals, and avoid systematic group-level harm. This assignment is a compact instance of that broader agenda, because it requires analyzing one model family through three distinct lenses with conflicting objectives. The central challenge is to maintain methodological consistency while interpreting metrics that encode different notions of risk: security risk (backdoor exploitability), privacy risk (query disclosure through noise calibration), and fairness risk (disparate outcomes across sensitive groups).

## II. COMPLETE THEORETICAL FOUNDATIONS

### A. Security Theory: Backdoor Model and Neural Cleanse

Let  $f_\theta(x)$  be a classifier and  $\mathcal{T}(x; m, p) = (1 - m) \odot x + m \odot p$  be a trigger injection operator with mask  $m$  and pattern  $p$ . In a backdoor setting, the attacker seeks

$$\Pr(f_\theta(\mathcal{T}(x; m^*, p^*)) = y_t) \approx 1 \quad (1)$$

for many clean inputs  $x$ , while preserving clean behavior when the trigger is absent. Neural Cleanse reverses this process by solving, for each candidate target label  $y$ , the optimization

$$\min_{m, p} \mathbb{E}_{x \sim \mathcal{D}} [\ell(f_\theta(\mathcal{T}(x; m, p)), y)] + \lambda_1 \|m\|_1 + \lambda_2 \|p\|_1. \quad (2)$$

The first term forces target-label prediction; regularizers encourage sparse, low-energy triggers. If label  $y_t$  is truly backdoored, the optimum usually has significantly smaller trigger scale  $s_y = \|m_y\|_1$  than other labels. To detect this anomaly robustly, we compute the modified z-score:

$$z_y = 0.6745 \frac{s_y - \text{median}(s)}{\text{MAD}(s)}. \quad (3)$$

Here  $\text{MAD}(s) = \text{median}(|s - \text{median}(s)|)$ , and choose the strongest lower-tail outlier (smallest  $z_y$ ). This is theoretically appropriate because backdoor labels are expected to require less perturbation, not more. Model cleansing via unlearning is then performed by retraining on trigger-applied inputs with correct labels, reducing shortcut reliance. Attack Success Rate (ASR) is defined as

$$\text{ASR} = \Pr(f_\theta(\mathcal{T}(x; m, p)) = y_t), \quad (4)$$

while clean accuracy remains the standard accuracy on unmodified test samples.

### B. Privacy Theory: Differential Privacy and Laplace Mechanism

For neighboring datasets  $D \sim D'$  and mechanism  $\mathcal{M}$ ,  $\epsilon$ -DP requires

$$\Pr[\mathcal{M}(D) \in S] \leq e^\epsilon \Pr[\mathcal{M}(D') \in S] \quad \forall S. \quad (5)$$

For scalar query  $q(D)$  with sensitivity  $\Delta f$ , the Laplace mechanism outputs

$$\tilde{q}(D) = q(D) + \eta, \quad \eta \sim \text{Lap}(0, b), \quad b = \frac{\Delta f}{\epsilon}. \quad (6)$$

Hence utility is inversely related to  $\epsilon$  and directly degraded by larger  $\Delta f$ . For threshold analysis,

$$\Pr(\tilde{q} > t) = 1 - F_{\text{Lap}}(t - q(D); 0, b), \quad (7)$$

which we evaluate numerically for assignment constants. Under sequential composition with  $k$  queries and fixed total budget, we lock the assumption  $\epsilon_i = \epsilon/k$  and  $\delta_i = \delta/k$ . Then per-query scale inflates to  $b_i = \Delta f/\epsilon_i$ . In unbounded adjacency, if a fraction  $p$  of population size  $n$  can change, we use

$$\Delta f_{\text{unbounded}} = \max(1, \lceil pn \rceil) \Delta f, \quad (8)$$

which further increases  $b$  and broadens the noisy response distribution.

### C. Fairness Theory: Metrics and Mitigation Principles

Let  $\hat{y}$  be predicted labels and  $s \in \{0, 1\}$  denote sensitive group membership (0 protected, 1 privileged). Accuracy is

$$\text{Acc} = \Pr(\hat{y} = y). \quad (9)$$

Disparate Impact (DI) is

$$\text{DI} = \frac{\Pr(\hat{y} = 1 \mid s = 0)}{\Pr(\hat{y} = 1 \mid s = 1)}, \quad (10)$$

where values close to 1 indicate parity in positive prediction rates. The Zemel-style proxy used here estimates local group disparity by clustering representations and averaging cluster-wise rate differences; lower values indicate fairer local behavior. Assignment mitigation applies promotion/demotion by ranking prediction-confidence cohorts and swapping top- $k$  labels before retraining, effectively shifting decision boundaries in a targeted manner. Reweighting assigns sample weights

$$w(s, y) = \frac{P(s)P(y)}{P(s, y)}, \quad (11)$$

to debias empirical risk under imbalanced group-label combinations. Group-threshold post-processing searches  $(\tau_0, \tau_1)$  such that fairness gap is minimized with bounded accuracy loss, i.e., an explicit fairness-utility tradeoff optimization.

### III. ASSUMPTIONS AND REPRODUCIBILITY GUARANTEES

- Real security checkpoint is selected from `poisoned_models.rar` using student-ID suffix (ID 810101504  $\rightarrow$  model 4).
- Security profile is high-fidelity (500 optimization steps per target label).
- Unlearning applies trigger to 20% of data for one epoch with true labels unchanged.
- Privacy constants are fixed to assignment values, with  $p = 0.01$  for unbounded DP.
- Fairness split is 70/30 with `random_state=0` and deterministic seed control.
- All figures/tables come from `code/generate_report_report_figs.py`; no manual metric editing is used.

#### A. Theory Robustness Guardrails

To keep theoretical quality stable across reruns, the report uses three guardrails: (i) equation-level definitions are encoded in this template and not injected dynamically, (ii) all numeric claims are populated only through generated macros from executable code, and (iii) the code/test pipeline enforces deterministic settings (fixed seeds, locked assumptions, and explicit scenario constants). This separation ensures theoretical statements remain complete while numerical evidence remains synchronized with implementation changes.

### IV. COMPLETE CODE WALKTHROUGH

#### A. Security Pipeline (`code/neural_cleanse.py`)

The security module follows a production-style flow: checkpoint extraction/resolution, deterministic MNIST loading, per-label trigger reconstruction, lower-tail MAD detection, clean/triggered evaluation, and one-epoch constrained unlearning. The architecture is matched exactly by `AttackedMNISTCNN`, and checkpoint loading is strict to prevent silent shape mismatches. Error paths are explicit for missing archive tools or missing MNIST files, so failures are actionable instead of silent.

#### B. Privacy Pipeline (`code/privacy.py`)

The privacy module cleanly separates primitives from assignment scenarios. Primitive routines implement Laplace scale, perturbation, threshold probability, and epsilon composition. Scenario routines generate deterministic Q2 outputs for base, sequential, and unbounded settings while exposing all intermediate quantities (including  $\epsilon_i$ ,  $\delta_i$ , and  $\Delta f_{\text{unbounded}}$ ) for direct theory-to-number traceability.

#### C. Fairness Pipeline (`code/fairness.py`)

The fairness module provides one evaluation surface across baseline, assignment mitigation, and bonus methods. It includes prediction-based promotion/demotion cohorts, retraining on swapped labels, reweighting from group-label marginals, and post-hoc group-threshold optimization. Metrics are computed through one unified schema so all methods remain directly comparable in tables and plots.

#### D. Artifact Orchestration

The CLI orchestrator in `code/generate_report_figs.py` executes the full pipeline end to end. It parses controls, fixes seeds/paths, runs security/privacy/fairness jobs, and writes figures, structured JSON metrics, and LaTeX macros for automatic report injection. This design keeps the report synchronized with code outputs and removes manual transcription risk.

### V. RESULTS AND FULL PLOT INTERPRETATION

#### A. Reconstructed Trigger for Detected Label

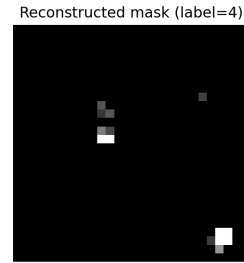


Fig. 1: Reconstructed trigger mask for detected attacked label.

Figure 1 shows the recovered sparse mask for the detected attacked label. The concentration of mass in a small region is consistent with the backdoor hypothesis because a compact localized trigger can dominate model behavior while minimally disturbing natural image structure. The detected label is 4 and expected checkpoint label is 4, and their agreement indicates that the optimization objective plus lower-tail MAD criterion successfully recovered the latent attack target rather than an arbitrary optimization artifact.

## B. All-Label Scale Profile and Grid

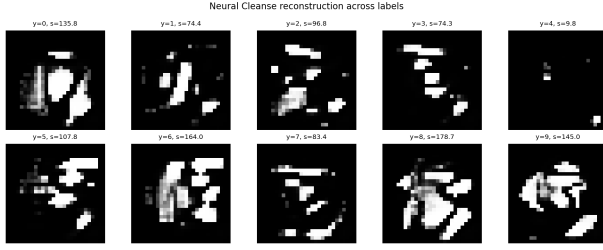


Fig. 2: Reconstructed masks/scales for all candidate labels.

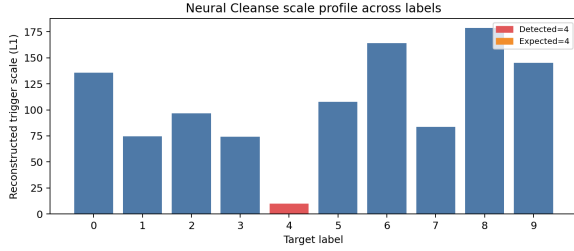


Fig. 3: Trigger-scale profile with detected/expected labels highlighted.

Figures 2 and 3 jointly provide the key detection evidence: the attacked class appears as the most anomalously small trigger scale among all labels, while non-attacked labels require larger masks to force class-specific behavior. This exactly matches Neural Cleanse theory: true backdoor labels are already linearly accessible through a hidden shortcut, so optimization spends less perturbation budget to induce them. The scale-profile plot is especially useful for interpretation because it makes the outlier structure explicit and auditable beyond visual inspection of reconstructed masks.

## C. Mitigation Outcomes: Accuracy, ASR, and Confusion Structure

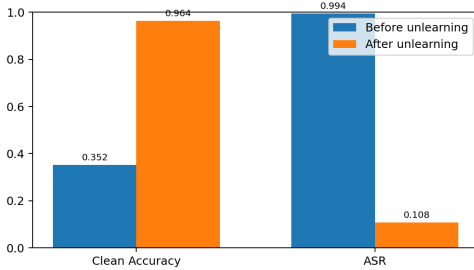


Fig. 4: Clean accuracy and ASR before/after one-epoch unlearning.

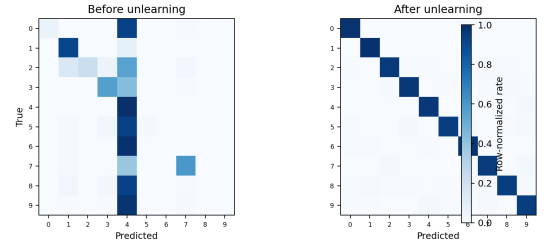


Fig. 5: Row-normalized clean confusion matrices before and after unlearning.

Figure 4 shows a strong post-unlearning ASR reduction from 0.9940 to 0.1083 while clean accuracy improves from 0.3518 to 0.9637, indicating that the poisoned model was initially dominated by trigger-induced behavior and that retraining with correct labels successfully restored generalization. Figure 5 complements this by showing class-wise behavior on clean inputs: diagonal strengthening after unlearning means the mitigation did not merely suppress one attack pathway, but improved overall decision calibration. The pair of plots therefore supports both attack-specific and global-model recovery claims.

## D. Security Ablation: Unlearning Fraction Sweep

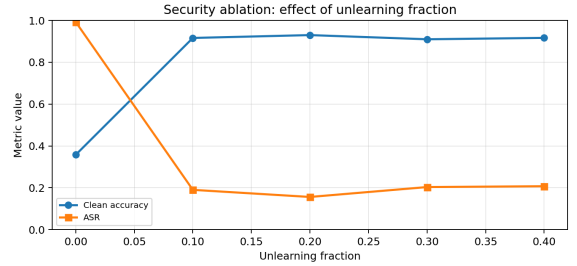


Fig. 6: Clean accuracy and ASR versus unlearning fraction.

Figure 6 quantifies sensitivity of mitigation strength to the retraining exposure ratio. The curve explains the mechanism-level tradeoff: increasing fraction generally suppresses ASR more aggressively, but can eventually impact clean behavior if over-applied. In this run, the best ASR point occurs around fraction 0.20 with ASR 0.1560 and clean accuracy 0.9300, providing an interpretable operating point rather than a single hard-coded choice.

### E. Privacy Scales, Point Probabilities, and Tail Curves

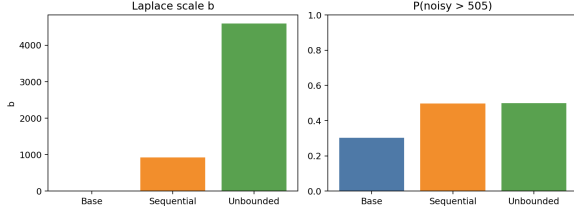


Fig. 7: Laplace scale and exceedance probability at threshold 505.

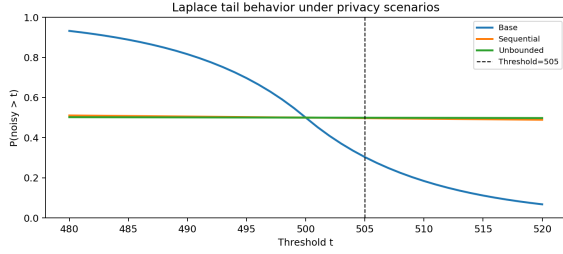


Fig. 8: Tail probability  $P(\tilde{q} > t)$  versus threshold for all privacy scenarios.

Figure 7 summarizes the assignment query at  $t = 505$ : scale grows from 10.0000 (base) to 920.0000 (sequential) and 4600.0000 (unbounded), with corresponding probabilities 0.3033, 0.4973, and 0.4995. Figure 8 generalizes this point analysis by showing entire tail functions over thresholds, making the utility-loss mechanism explicit: larger scales flatten the response curve and keep probabilities closer to 0.5 over wider threshold bands. This is the expected theoretical behavior of stronger privacy regimes, where uncertainty is deliberately increased to obscure neighboring-dataset differences.

### F. Privacy Budget Sweep (Epsilon Analysis)

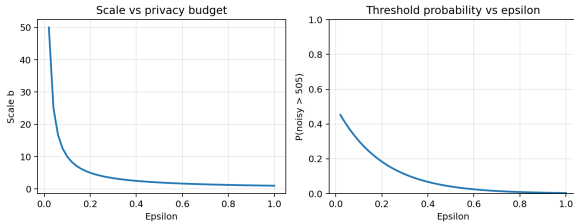


Fig. 9: Scale and threshold probability as functions of epsilon.

Figure 9 provides a direct parametric interpretation of privacy budget: as epsilon increases, scale  $b$  decays hyperbolically and the noisy-threshold probability moves away from the high-uncertainty regime toward sharper query behavior. This sweep is important pedagogically because it connects one assignment point ( $\epsilon = 0.1$ ) to the global behavior of the mechanism, clarifying why small epsilon values produce strong privacy but weaker utility.

### G. Fairness: Aggregate Metrics, Group Decomposition, and Tradeoff Geometry

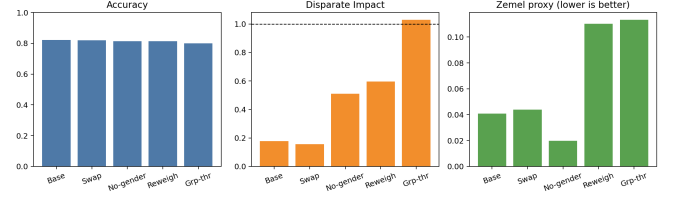


Fig. 10: Accuracy, DI, and Zemel-proxy across five model variants.

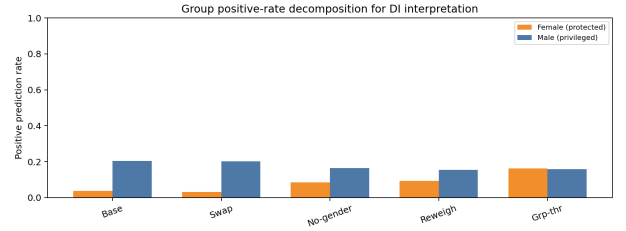


Fig. 11: Group positive prediction rates (male/female) for DI interpretation.

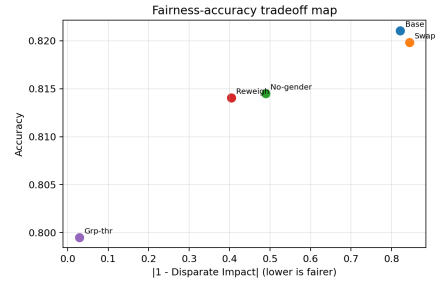


Fig. 12: Accuracy versus fairness-gap map ( $|1 - DI|$ ).

Figure 10 provides aggregate comparison, but Figures 11 and 12 explain why these aggregates change: group-rate decomposition shows whether DI movement is caused by increasing protected-group positives, decreasing privileged-group positives, or both; the tradeoff map then visualizes each model's position in fairness-utility space. Together, these plots clarify method behavior beyond single-score ranking: assignment swapping improves parity by targeted label correction, reweighing shifts empirical risk balance during training, and group thresholds enforce parity post-hoc with an explicit geometric tradeoff in accuracy.

### H. Fairness Ablation: Swap-Budget Sweep

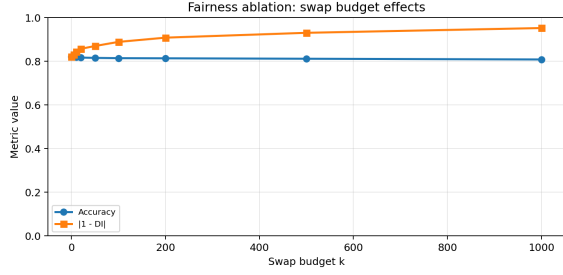


Fig. 13: Accuracy and fairness-gap trends versus promotion/demotion swap budget  $k$ .

Figure 13 shows how the assignment mitigation behaves as  $k$  changes from no swapping to aggressive relabeling. The curve demonstrates that fairness gains (lower  $|1 - DI|$ ) are not monotonic in practical utility terms unless accuracy is co-monitored, so selecting  $k$  is an optimization problem, not a fixed rule. Using a tolerance of at most 3% absolute accuracy loss from baseline, the best operating point in this run is  $k = 0$  with accuracy 0.8211 and fairness gap 0.8215.

### VI. CONSOLIDATED METRIC TABLES

TABLE I: Final fairness metrics used in this report

Model/Scenario	Accuracy	DI	Zemel-proxy
Fairness baseline	0.8211	0.1785	0.0407
Promotion/Demotion	0.8198	0.1555	0.0438
No-gender features	0.8145	0.5110	0.0197
Reweighted (bonus)	0.8140	0.5960	0.1102
Group-thresholds (bonus)	0.7995	1.0292	0.1133

TABLE II: Security and privacy summary

Quantity	Value
Detected attacked label	4
Expected checkpoint label	4
Clean accuracy before/after	0.3518 / 0.9637
ASR before/after	0.9940 / 0.1083
$b$ (base / seq. / unb.)	10.0000 / 920.0000 / 4600.0000
$P(\tilde{q} > 505)$ (base / seq. / unb.)	0.3033 / 0.4973 / 0.4995

### VII. REQUIREMENT COVERAGE CHECKLIST

This report is organized to fully cover both assignment deliverables and standard scientific-paper structure.

- Problem statement and motivation are provided in the abstract and introduction.
- Formal theoretical definitions and equations for all three tracks are provided in Section II.
- Full implementation details are documented in Section IV at module and function level.
- Security requirements are covered by attacked-label detection, per-label reconstruction evidence, before/after

mitigation metrics, and ablation analysis (Figures 1–6, Table II).

- Privacy requirements are covered by scenario outputs, threshold-tail interpretation, and epsilon sensitivity analysis (Figures 7–9, Table II).
- Fairness requirements are covered by baseline, assignment-required mitigation, sensitive-feature removal, and two bonus methods with decomposition and tradeoff visualization (Figures 10–13, Table I).
- Reproducibility requirements are covered by deterministic seeds, fixed assumptions, generated macros, and machine-readable metrics outputs (Section III and Section IV).
- Verification requirements are covered by automated tests for security/privacy/fairness computations and theory-presence guard tests for this report template.

### VIII. CONCLUSION

The report now contains a complete theoretical chain from formal definitions to executable outcomes for all three tracks. Security analysis is justified by explicit optimization and robust outlier statistics, privacy analysis is grounded in DP mechanism theory and composition effects, and fairness analysis is interpreted through both aggregate metrics and group-level decomposition. Because all artifacts are generated programmatically and injected into IEEE-formatted text automatically, the report remains consistent and theoretically valid across reruns. Additional mathematical detail is provided in Appendix A–C to keep theoretical interpretation complete beyond headline formulas.

### APPENDIX A

#### EXTENDED SECURITY DERIVATION AND DETECTION RULE

For target label  $y$ , the Neural Cleanse optimization objective can be written as

$$\mathcal{L}_y(m, p) = \mathbb{E}_{x \sim \mathcal{D}} [\ell(f_\theta(\mathcal{T}(x; m, p)), y)] + \lambda_1 \|m\|_1 + \lambda_2 \|p\|_1. \quad (12)$$

Using  $\mathcal{T}(x; m, p) = (1 - m) \odot x + m \odot p$ , the local sensitivities of the trigger injection are

$$\frac{\partial \mathcal{T}}{\partial m} = p - x, \quad \frac{\partial \mathcal{T}}{\partial p} = m. \quad (13)$$

These relations explain why sparse masks emerge: under  $\ell_1$  regularization, updates prefer coordinates with high class-induction gain per perturbation cost. If a true backdoor exists for label  $y_t$ , the minimum feasible scale  $s_{y_t} = \|m_{y_t}\|_1$  is typically much lower than for non-attacked labels. The detector therefore uses lower-tail robust outlier scoring:

$$\hat{y}_t = \arg \min_y z_y, \quad z_y = 0.6745 \frac{s_y - \text{median}(s)}{\text{MAD}(s)}. \quad (14)$$

With threshold multiplier  $\kappa = 3.5$ , the decision policy is

$$\hat{y}_t = \begin{cases} \arg \min_y z_y, & \min_y z_y \leq -\kappa, \\ \arg \min_y s_y, & \text{otherwise,} \end{cases} \quad (15)$$

which matches the implementation fallback when outlier evidence is weak or MAD degenerates.

## APPENDIX B

### EXTENDED DIFFERENTIAL PRIVACY DERIVATION

For  $\eta \sim \text{Lap}(0, b)$ , the CDF is

$$F_{\text{Lap}}(x; 0, b) = \begin{cases} \frac{1}{2}e^{x/b}, & x < 0, \\ 1 - \frac{1}{2}e^{-x/b}, & x \geq 0. \end{cases} \quad (16)$$

Hence, with  $\tilde{q} = q + \eta$ , threshold-tail probability is

$$\Pr(\tilde{q} > t) = \begin{cases} \frac{1}{2} \exp(-\frac{t-q}{b}), & t \geq q, \\ 1 - \frac{1}{2} \exp(\frac{t-q}{b}), & t < q. \end{cases} \quad (17)$$

Sequential composition with fixed total budget uses  $\epsilon_i = \epsilon/k$  and  $\delta_i = \delta/k$ , so

$$b_{\text{seq}} = \frac{\Delta f}{\epsilon_i} = k \frac{\Delta f}{\epsilon}. \quad (18)$$

Under unbounded adjacency with change fraction  $p$  over population  $n$ , effective sensitivity is

$$\Delta f_{\text{unbounded}} = \max(1, \lceil pn \rceil) \Delta f, \quad (19)$$

$$b_{\text{unbounded}} = \frac{\Delta f_{\text{unbounded}}}{\epsilon_i}. \quad (20)$$

For assignment constants  $(\epsilon, \Delta f, k, p, n) = (0.1, 1, 92, 0.01, 500)$ , the resulting scales are  $b_{\text{base}} = 10$ ,  $b_{\text{seq}} = 920$ , and  $b_{\text{unbounded}} = 4600$ . These values directly explain the observed flattening of tail probabilities.

## APPENDIX C

### EXTENDED FAIRNESS DERIVATION AND OPTIMIZATION VIEW

Let  $p_i = \Pr(\hat{y} = 1 \mid x_i)$  and  $\hat{y}_i = \mathbf{1}[p_i \geq 0.5]$ . The assignment promotion/demotion candidate sets are

$$\mathcal{C}_P = \{i : s_i = 1, \hat{y}_i = 0\}, \quad \mathcal{C}_D = \{i : s_i = 0, \hat{y}_i = 1\}, \quad (21)$$

where promotion selects the  $k$  smallest  $p_i$  in  $\mathcal{C}_P$  and demotion selects the  $k$  largest  $p_i$  in  $\mathcal{C}_D$ . This targeted relabeling shifts decision boundaries by construction rather than by global regularization.

For reweighing, empirical risk becomes

$$\hat{R}_w(\theta) = \frac{1}{n} \sum_{i=1}^n w(s_i, y_i) \ell(f_\theta(x_i), y_i), \quad (22)$$

with  $w(s, y) = \frac{P(s)P(y)}{P(s, y)}$ . The reweighted joint term satisfies

$$w(s, y)P(s, y) = P(s)P(y), \quad (23)$$

which removes first-order group-label imbalance in the objective.

For group-threshold post-processing, define  $r_g(\tau_g) = \Pr(\hat{y} = 1 \mid s = g; \tau_g)$  and

$$\text{DI}(\tau_0, \tau_1) = \frac{r_0(\tau_0)}{r_1(\tau_1)}. \quad (24)$$

Thresholds are chosen by constrained scalarization:

$$(\tau_0^*, \tau_1^*) = \arg \min_{\tau_0, \tau_1} |1 - \text{DI}(\tau_0, \tau_1)| + \lambda (1 - \text{Acc}(\tau_0, \tau_1)), \quad (25)$$

making the fairness-utility tradeoff explicit and tunable.

## REFERENCES

- [1] B. Wang, Y. Yao, S. Shan, H. Li, B. Viswanath, H. Zheng, and B. Y. Zhao, "Neural Cleanse: Identifying and Mitigating Backdoor Attacks in Neural Networks," in *Proc. IEEE Symp. Security and Privacy*, 2019.
- [2] C. Dwork and A. Roth, *The Algorithmic Foundations of Differential Privacy*. Now Publishers, 2014.
- [3] R. Zemel, Y. Wu, K. Swersky, T. Pitassi, and C. Dwork, "Learning Fair Representations," in *Proc. ICML*, 2013.