PRIVATE AND INTERPRETABLE CLINICAL PREDICTION WITH QUANTUM-INSPIRED TENSOR TRAIN MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Machine learning in clinical settings must balance predictive accuracy, interpretability, and privacy. While models like logistic regression (LR) are valued for transparency, they remain vulnerable to privacy attacks that expose training data. We empirically assess these risks by designing attacks that identify which public datasets were used to train a model under varying levels of adversarial access, applying them to LORIS, a publicly available LR model for immunotherapy response prediction. Our findings show that LORIS leaks significant training-set information, especially under white-box access, and that common practices such as cross-validation exacerbate these risks. Even black-box access via the public web interface allows training data identification. To mitigate these vulnerabilities, we propose a quantum-inspired defense using tensor train (TT) models. Tensorizing LR obfuscates parameters while preserving accuracy, reducing white-box attacks to random guessing and degrading black-box attacks comparably to Differential Privacy. TT models retain LR interpretability and extend it through efficient computation of marginal and conditional distributions. Although demonstrated on LORIS, our approach generalizes broadly, positioning TT models as a practical foundation for private, interpretable, and effective clinical prediction.

1 Introduction

Machine learning (ML) is increasingly used for clinical prediction but poses critical privacy risks, as models trained on sensitive medical data can inadvertently leak individual information (Fredrikson et al., 2014; Sweeney, 2015). In domains where interpretability is essential, intuitive models like logistic regression (LR) are often preferred, yet they are particularly vulnerable to such attacks.

In this work, we propose a quantum-inspired approach to privacy protection based on tensor network (TN) models, focusing on tensor trains (TT). Building on recent methods for tensorizing pre-trained ML models into TT form (Pareja Monturiol et al., 2025), we obfuscate them to enhance privacy while preserving accuracy and interpretability.

To assess the privacy risks of clinical models and the protection offered by TTs, we attack LORIS (Chang et al., 2024), a publicly available LR model for immunotherapy response prediction hosted on a U.S. government website. We design a membership inference attack under both black-box (BB) and white-box (WB) access, using a shadow model approach that trains multiple models with varied hyperparameters and datasets, followed by an adversarial meta-classifier to predict which public datasets were included in the training set.

Our results show that tensorizing LORIS degrades attack performance across all access levels, reducing WB attacks to random guessing. We argue that this method achieves privacy protection comparable to Differential Privacy (DP), while maintaining similar levels of predictive accuracy. Additionally, we show that common practices like cross-validation, when used to deploy averaged models as in LORIS, can severely compromise privacy, enabling accurate training-set identification even from BB access via the public web interface. TT approximations preserve key properties of LORIS, such as response monotonicity, while enhancing interpretability through efficient computation of marginals and conditionals. This supports feature-sensitivity analysis and enables the construction of cancer-type-specific models without retraining.

These findings underscore how easily training data can be extracted with minimal knowledge of a model and its training procedure. Although our attack targets LORIS, the methodology generalizes to a broad range of models and settings. We therefore advocate for the routine use of tensorization as a practical strategy for privacy-preserving, interpretable, and efficient ML, especially in clinical domains handling sensitive data.

The remainder of this paper is structured as follows. Section 2 reviews related work and preliminaries. Section 3 outlines our setting, attack, and defenses, and presents the results. Section 4 analyzes the interpretability of TT models in comparison to LORIS. Finally, Section 5 discusses conclusions and future directions.

2 RELATED WORK AND PRELIMINARIES

The widespread adoption of ML systems increases the risk of leaking sensitive personal data. Prior work has extensively examined these vulnerabilities and proposed various defenses.

2.1 PRIVACY ATTACKS

A wide range of attacks exploit privacy vulnerabilities in ML, leveraging either BB or WB access. Key examples include model inversion (Fredrikson et al., 2014), model classification (Ateniese et al., 2015), and membership inference (Shokri et al., 2017), which vary in scope from extracting individual samples to uncovering global patterns. In this work, we adopt the membership inference approach to identify groups of samples present in the training set.

More recently, reconstruction attacks have aimed to recover exact training samples. Some rely on shadow-model training (Balle et al., 2022), while others exploit optimization properties of models trained with Stochastic Gradient Descent (SGD) (Haim et al., 2022; Oz et al., 2024). Notably, for LR, such attacks can yield closed-form solutions (Balle et al., 2022), underscoring the vulnerability of simple, widely used models.

2.2 Defense Mechanisms

Given the diversity of privacy-related attacks, various defense mechanisms have been proposed. Among these, Differential Privacy (Dwork, 2006b) stands out for its rigorous framework. DP quantifies the likelihood that an attacker can infer whether a specific user's data was included in a statistical process. A randomized algorithm \mathcal{A} is ε -DP if, for any set of outcomes \mathcal{S} in the range of \mathcal{A} , it satisfies:

$$\log\left(\frac{P[\mathcal{A}(D) \in \mathcal{S}]}{P[\mathcal{A}(D') \in \mathcal{S}]}\right) \le \varepsilon,\tag{1}$$

where D and D' differ by a single element. This metric guides the addition of calibrated noise to achieve a target ε , based on the sensitivity of the function being protected (Dwork, 2006a; Dwork & Roth, 2014). However, the noise required for strong privacy guarantees often degrades model performance and may exacerbate group disparities (Bagdasaryan et al., 2019; Hansen et al., 2024). As a result, there is no consensus on how to set ε meaningfully (Garfinkel et al., 2018); while small values are theoretically ideal, larger values may still prevent reconstruction attacks in practice without significantly harming accuracy (Ziller et al., 2024).

Beyond DP, recent work has explored whether standard ML practices can improve privacy. Pruning, for example, introduces small errors that resemble DP-like protection (Huang et al., 2020), while knowledge transfer reduces dependence on specific training data (Shejwalkar & Houmansadr, 2020). Our approach draws on these ideas: rather than enforcing DP, we approximate pre-trained models in TT form, achieving both BB protection and strong WB obfuscation.

2.3 TENSOR TRAIN MODELS

Tensor networks are low-rank decompositions of high-dimensional tensors with roots in quantum many-body physics. They offer compact, interpretable representations of quantum states (Pérez-García et al., 2007; Orús, 2014; Cirac et al., 2021) and have recently been adapted to machine learning (Stoudenmire & Schwab, 2016; Novikov et al., 2018). TNs have been applied to neural network

(NN) compression (Novikov et al., 2015; Tomut et al., 2024), explainable AI (Tangpanitanon et al., 2022; Aizpurua et al., 2024), and anomaly detection (Wang et al., 2020). Importantly, TNs offer formal WB privacy guarantees: due to gauge freedom, multiple parameterizations can represent the same model, effectively obfuscating all but its BB behavior (Pozas-Kerstjens et al., 2024).

Throughout this work, we focus on one-dimensional TNs, specifically tensor trains (Oseledets, 2011). An order-N tensor $T \in \mathbb{R}^{d^N}$ admits a TT representation with ranks r_n if it can be written as

$$T(i_1, \dots, i_N) = G_1(i_1) \cdots G_N(i_N),$$
 (2)

where the $cores\ G_n$ are $r_{n-1} \times r_n$ matrices and $r_0 = r_N = 1$. This structure also supports continuous functions of the form

$$f(x_1, \dots, x_N) = \sum_{i_1, \dots, i_N} W(i_1, \dots, i_N) \,\phi_1(i_1, x_1) \cdots \phi_N(i_N, x_N), \tag{3}$$

where W is a TT-format coefficient tensor and $\phi_n(i_n,x_n)$ are vector-valued *embedding* functions indexed by i_n . To ensure non-negative probability scores, it is standard to define distributions via the Born rule: $p(x) = |f(x)|^2$. Further details on TTs, including efficient marginalization and conditioning, are provided in Appendix A.

TTs can be trained using SGD or physics-inspired variants (Stoudenmire & Schwab, 2016). Alternatively, TT representations can be constructed via low-rank decompositions, bypassing high-dimensional optimization. Recent techniques based on sketching (Hur et al., 2023) and cross interpolation (Fernández et al., 2025) achieve this using only function evaluations; i.e., BB access, to approximate continuous functions in TT form. A recent method, TT-RSS, extends this idea to tensorize pre-trained NNs using a small evaluation dataset (Pareja Monturiol et al., 2025). We adopt this approach to tensorize LR models.

3 PRIVACY ANALYSIS

To evaluate the privacy risks of clinical prediction models and compare defense strategies, we design a membership inference attack based on shadow-model training. Assuming an adversary with access to multiple public datasets, the attack seeks to determine which of them were included in a model's training set under varying levels of adversarial access. As a case study, we target LORIS, a publicly available model introduced by Chang et al. (2024) for immunotherapy response prediction. Below we define the setting, describe the adversarial assumptions and attack, and present the experimental setup, with results reported at the end of the section.

3.1 SETTING AND NOTATION

Let $\mathcal{D} = \{D_1, \dots, D_M\}$ be the set of public datasets, and define $\mathcal{D}_{\cup} = \{\bigcup \mathcal{C} \mid \mathcal{C} \in \mathcal{P}(\mathcal{D}) \setminus \{\varnothing\}\}$, where $\mathcal{P}(\mathcal{D})$ is the power set. A training set $D_{\cup} \in \mathcal{D}_{\cup}$ is the union of one or more $D_m \in \mathcal{D}$. Using the indicator vector $\mathbf{1}(D_{\cup})$, we represent D_{\cup} as a multi-hot vector with entries 1 for datasets $D_m \subset D_{\cup}$ and 0 otherwise.

We define the training algorithm as follows: given a model architecture Φ , hyperparameters $H_{\Phi} \in \mathcal{H}_{\Phi}$, and a training set $D_{\cup} \in \mathcal{D}_{\cup}$, the training mechanism $\mathcal{T}_{\Phi} : \mathcal{H}_{\Phi} \times \mathcal{D}_{\cup} \to \Theta$ outputs parameters $\theta \in \Theta$ such that $\Phi_{\theta}(\cdot)$ is a trained model. In practice, \mathcal{T}_{Φ} is stochastic due to factors such as random initialization or mini-batch selection in SGD, so for fixed H_{Φ} and D_{\cup} we interpret $\mathcal{T}_{\Phi}(H_{\Phi}, D_{\cup})$ as sampling from a model distribution. In addition, since training data are typically standardized for stability, yielding coefficients defined on standardized inputs, we assume that \mathcal{T}_{Φ} returns rescaled parameters that operate on raw input data. Details of the standardization and rescaling procedures are provided in Appendix B.

To mitigate bias and overfitting, it is standard to use K-fold cross-validation, which partitions D_{\cup} into K folds and trains K models, each on K-1 folds. We denote by $\mathcal{T}_{\Phi}^{J,K}: \mathcal{H}_{\Phi} \times \mathcal{D}_{\cup} \to \Theta$ the procedure that applies \mathcal{T}_{Φ} with fixed H_{Φ} , performs J repetitions of K-fold cross-validation, corrects for feature standardization, and averages the resulting parameters into a final model.

3.2 DESCRIPTION OF THE ATTACK

The adversary knows the model architecture Φ , the public datasets \mathcal{D} , and a finite set of hyperparameters \mathcal{H}_{Φ} . They also know the training mechanisms \mathcal{T}_{Φ} , $\mathcal{T}_{\Phi}^{J,K}$, and have sufficient resources to train shadow models and meta-classifiers. Access to the target model is limited to restricted information $h(\Phi_{\theta})$, which we categorize into three independent access levels:

- b-Weak black-box (b-WBB): Access to outputs discretized into b bins, e.g., b=2 gives binary outputs. Values <0.5 map to the lower bin limit, and >0.5 to the upper.
- Strong black-box (SBB): Access to raw continuous scores. As b grows to machine precision, b-WBB converges to SBB.
- White-box (WB): Access to model parameters. Although parameters allow computing outputs, we treat BB and WB separately to assess each source of information, while stronger attacks may combine both.

The attack proceeds by constructing a dataset of shadow models, each trained under different hyperparameter configurations and training sets. From each model, we collect the relevant information together with the corresponding public datasets used for training. This forms the input to a multilabel classifier, which learns to identify the presence of public datasets in the training sets. Formally, the attack consists of the following steps:

1. For each $H_{\Phi} \in \mathcal{H}_{\Phi}$ and $D_{\cup} \in \mathcal{D}_{\cup}$, train R shadow models using \mathcal{T}_{Φ} or $\mathcal{T}_{\Phi}^{J,K}$.

2. Build $\{(h(\Phi_{\theta^i}), \mathbf{1}(D_{\cup}^i))\}_{i=1}^{R|\mathcal{H}_{\Phi}||\mathcal{D}_{\cup}|}$, where $h(\cdot)$ denotes available information: under BB access it returns outputs on S samples, and under WB access it returns parameters θ^i .

3. Train an adversarial model minimizing independent cross-entropy losses for each D_m , yielding $\mathcal{A}:\Theta\to[0,1]^M$, where entry m gives the probability that $D_m\subset D_{\cup}$.

3.3 EXPERIMENTAL SETUP

We briefly describe the datasets, models, and implementation details. All experiments¹ were run on an Intel Xeon CPU E5-2620 v4 with 256 GB RAM and an NVIDIA GeForce RTX 3090, using Scikit-Learn for LR models (Pedregosa et al., 2011), Diffprivlib for DP variants (Dwork, 2006b), and TensorKrowch for TT models (Pareja Monturiol et al., 2024).

3.3.1 Datasets

To build the public set \mathcal{D} we use the cohorts employed to train and evaluate LORIS, which include clinical, pathological, and genomic features with a binary treatment-response label. For details see Chang et al. (2024); we list them here with shorthand identifiers and sample sizes: Cho1 (964) and Cho2 (515), *train* and *test* partitions from Chowell et al. (2022); MSK1 (453) and MSK2 (104) from Chang et al. (2024); Shim (198) from Shim et al. (2020); Kato (35) from Kato et al. (2020); Vang (246) from Vanguri et al. (2022); Ravi (309) from Ravi et al. (2023); and Prad (57) from Pradat et al. (2023). In all cases, response is imbalanced, with \sim 30% of patients responding to treatment.

We use 6-feature models: Tumor Mutational Burden (TMB), Previous Systematic Therapy History (PSTH), Albumin, Neutrophil-to-Lymphocyte Ratio (NLR), Age, and cancer type. Cancer type is divided into 16 binary variables, yielding 21 input features in total.

3.3.2 TARGET MODELS

As target models, we consider several variants of LR. Following Chang et al. (2024), we train averaged models via $\mathcal{T}_{\Phi}^{J,K}$ with J=20 and K=3. While LORIS used larger values of J and K, we found this configuration sufficient to obtain comparable results. For comparison, we also train vanilla LRs through a single run of \mathcal{T}_{Φ} on an 80% split of D_{\cup} . In both cases, the hyperparameters are solver = "saga", penalty = "elasticnet", class_weight = "balanced", max_iter = 100, 11_ratio $\in \{0, 0.5, 1\}$, and $C \in \{0.1, 1, 10\}$, forming the uncertainty set \mathcal{H}_{Φ} .

¹The code is publicly available at: https://anonymous.4open.science/r/tts4privacy

For each dataset D_{\cup} , hyperparameter configuration H_{Φ} , and training method (vanilla or averaged), the adversary trains R=100 models. Each model is then tensorized via TT-RSS (Pareja Monturiol et al., 2025), using 50 random samples from D_{\cup} as pivots evaluated through b-WBB access to the LR. While b=2 maximizes privacy, it severely degrades performance; we therefore use b=6 as a trade-off. The resulting TTs have N=22 cores (including one for the output), ranks $r_n=2$ for all n, input dimensions d=2, and use polynomial embeddings $\phi_n(\cdot,x)=[1,x]$.

Due to the monotonicity of LR, model parameters can be exactly recovered from scores (see Appendix C), making SBB and WB access equivalent, although WB is typically easier to exploit. Since tensorization approximates LR outputs with a TT representation, it is also possible to recover LR coefficients from TT evaluations. To test whether these reconstructed coefficients leak more information than TT parameters, we collect them for each TT and perform WB attacks; we refer to these as LR-TT models.

Finally, for comparison with a standard privatization approach, we also train DP models (LR-DP) from scratch. In this case, rather than privatizing pre-trained LRs, we train new LR-DP models directly. Since DP training of LR is restricted to solver = "lbfgs" and penalty = "l2", we fix max_iter = 100 and vary the privacy budget $\varepsilon \in \{10^{-2}, 10^{-1}, 10^0, 10^1, 10^2, \infty\}$, where $\varepsilon = \infty$ corresponds to the non-DP case. Only vanilla models are considered, as averaging would cancel the injected noise and effectively increase ε .

3.3.3 ADVERSARIAL MODELS

To attack the models described above we use NN-based adversaries: MLP multi-label classifiers with three hidden layers of sizes 32, 16 and 8, and an output layer of size 9 (one output per public dataset). The input layer size depends on the access type. For BB attacks, each shadow model is evaluated on S=100 samples (the same S samples for all models), randomly drawn from $\bigcup \mathcal{D}$; the resulting vector of raw or discretized outputs is the adversary input. For WB attacks we collect full model parameters: for LR, 22 parameters (21 coefficients + intercept); for TT, all N=22 cores G_n vectorized and concatenated into a single 168-dimensional vector. All parameters are rescaled to operate on raw inputs (see Appendix B). The MLPs are trained with activation = "relu", solver = "adam" and max_iter = 100. Since WB attacks exhibited greater variability, we applied 5-fold cross-validation with predictions averaged across folds; on top of this, to obtain robust statistics we repeat 5-fold cross-validation five times for both WB and BB attacks.

3.4 RESULTS

To evaluate the overall performance of our attacks, Table 1 reports Hamming scores, i.e., the proportion of correct label predictions across all public datasets and shadow-model instances. These results yield three main observations. (i) Scores increase with deeper levels of access, with SBB and WB achieving surprisingly high values. Although WB can theoretically be recovered from SBB in LR models, in practice this may require evaluation at specific or additional samples (see Appendix C); hence, SBB attacks sometimes underperform WB despite their theoretical equivalence. (ii) Averaged models are consistently more vulnerable than vanilla ones, despite their similar predictive performance (see Appendix D.1). The variance reduction from cross-validation, while mitigating sample bias, amplifies differences across models and thus facilitates attacks. Notably, WB attacks on averaged models achieve nearly perfect classification. (iii) Original LR models yield the highest attack scores, underscoring their vulnerability when released without protection.

TT models achieve the lowest attack scores among the non-DP cases (LR and LR-DP with $\varepsilon=\infty$), across all access types. Randomization of TT cores is particularly effective, reducing WB attacks to near-random guessing. WB attacks on LR-TT coefficients perform better than attacks on TT parameters, but remain close to 2-WBB results. These findings confirm that TTs effectively restrict leakage to BB information. To contextualize these findings, we also evaluate shadow-model performance. As shown in the tables of Appendix D.1, TT models maintain balanced accuracy in nearly all datasets, with only minor drops (1–2% in a few cases), and achieve comparable AUC scores, though with larger differences in the Kato dataset.

As expected, DP models exhibit attack scores that increase with ε . At $\varepsilon=100$, performance is nearly indistinguishable from the non-DP case ($\varepsilon=\infty$), offering negligible privacy gains. Performance metrics in Appendix D.1 show that both settings achieve AUC scores similar to original LRs,

Table 1: Hamming scores (mean \pm std) of adversarial multi-label classifiers.

		2-WBB	SBB	WB
LR	vanilla averaged	0.7927 ± 0.0062 0.8730 ± 0.0113	0.8798 ± 0.0217 0.9502 ± 0.0407	$\begin{array}{c} 0.8995 \pm 0.0015 \\ 0.9974 \pm 0.0032 \end{array}$
TT	vanilla averaged	0.7166 ± 0.0038 0.7404 ± 0.0028	0.8180 ± 0.0137 0.8650 ± 0.0134	0.5590 ± 0.0273 0.5770 ± 0.0176
LR-TT	vanilla averaged	_	_	$0.7398 \pm 0.0016 \\ 0.7803 \pm 0.0022$
LR-DP $(\varepsilon = 10^{-2})$	vanilla	0.5412 ± 0.0032	0.5428 ± 0.0043	0.5258 ± 0.0343
LR-DP $(\varepsilon = 10^{-1})$	vanilla	0.5408 ± 0.0028	0.5414 ± 0.0039	0.5307 ± 0.0365
LR-DP $(\varepsilon = 10^0)$	vanilla	0.5792 ± 0.0029	0.5871 ± 0.0036	0.5359 ± 0.0359
LR-DP $(\varepsilon = 10^1)$	vanilla	0.7055 ± 0.0039	0.7740 ± 0.0058	0.6379 ± 0.0141
LR-DP $(\varepsilon = 10^2)$	vanilla	0.7610 ± 0.0071	0.8660 ± 0.0219	0.8636 ± 0.0083
LR-DP $(\varepsilon = \infty)$	vanilla	0.7576 ± 0.0072	0.8739 ± 0.0240	0.8977 ± 0.0030

but with lower balanced accuracies, reflecting bias toward the majority class. A higher prediction threshold could mitigate this imbalance. Since 2-WBB access binarizes outputs at threshold 0.5, this effect likely impacts attack accuracy; indeed, SBB attacks reach accuracies similar to non-DP LRs once $\varepsilon \geq 100$. Among tested configurations, $\varepsilon = 10$ offers the best trade-off, matching the utility and robustness of TT models, while $\varepsilon < 10$ causes substantial performance loss.

We also report per-dataset attack performances for LR and TT vanilla models in Appendix D.2, which further support the conclusions of this analysis.

3.4.1 Example: Cho1 vs. Cho1 + Kato

As an illustrative case, we consider the extreme task of distinguishing models trained only on Cho1 (964 samples) from those trained on Cho1 plus the small Kato cohort (35 samples). This simulates a high-risk scenario where an adversary detects the inclusion of a very small subgroup. Table 2 shows Hamming scores for the Kato label. As expected, 2-WBB attacks are nearly random. In contrast, averaged LRs reach ~75% detection under SBB and achieve almost perfect classification under WB. Notably, even vanilla LRs under WB access attain ~73% accuracy.

Table 2: Hamming scores (mean \pm std) of adversarial classification of models trained on Cho1 or Cho1+Kato, evaluated on the Kato label.

		2-WBB	SBB	WB
LR	vanilla averaged	0.5383 ± 0.0237 0.5278 ± 0.0362	0.5410 ± 0.0449 0.7464 ± 0.2312	$\begin{array}{c} \textbf{0.7289} \pm \textbf{0.0279} \\ \textbf{0.9989} \pm \textbf{0.0022} \end{array}$
TT	vanilla averaged	0.5189 ± 0.0279 0.5282 ± 0.0334	0.5261 ± 0.0243 0.5658 ± 0.0751	0.4931 ± 0.0237 0.4961 ± 0.0226
LR-TT	vanilla averaged	_	_	$0.5468 \pm 0.0197 \\ 0.5677 \pm 0.0286$

For context, Appendix D.3 reports model performance on Cho1 and Kato separately. TT models degrade somewhat on Kato, especially in AUC, but this alone does not explain the results: even LRs with low balanced accuracy still enable highly accurate attacks.

Overall, these results show that even a 35-sample cohort can be reliably identified within a large dataset. Model averaging and WB access amplify leakage, while TT models remain robust and do not reveal the presence of Kato under any access type.

3.4.2 ATTACKING PUBLICLY AVAILABLE MODELS

We illustrate the risk of WB attacks on publicly available LR coefficients from LORIS: (i) those released in Chang et al. (2024), and (ii) coefficients we reconstructed from the online interface.² Although the interface returns rounded probabilities rather than exact scores, by approximately inverting the monotonic mapping created for LORIS (Fig. 3) we obtain usable coefficients (Appendix C).

Applying our WB attack, Table 3 shows that Cho1 is correctly identified as the training dataset in both cases, consistent with Chang et al. (2024). Recovered coefficients are noisier, assigning some probability to Cho2 and MSK2, but Cho1 remains dominant. Since Cho1, Cho2, and MSK2 all originate from patients at Memorial Sloan Kettering Cancer Center (MSK), these spurious assignments likely reflect shared data characteristics. These results demonstrate that even with noisy reconstructed coefficients, adversaries can still infer training data membership with high confidence, highlighting the privacy risks of releasing or exposing LR parameters.

Table 3: WB attack scores for LORIS coefficients, using (i) the released parameters (Chang et al., 2024) and (ii) coefficients reconstructed from the online interface.

	Cho1	Cho2	MSK1	MSK2	Shim	Kato	Vang	Ravi	Prad
Released	0.9987	0.0313	0.0024	0.0398	0.0180	0.0093	0.2048	0.0240	0.0169
Reconstructed	0.8834	0.6535	0.0650	0.7296	0.0078	0.0154	0.0873	0.0112	0.4390

4 Interpretability with tensor trains

Beyond privacy guarantees, interpretability is essential in clinical prediction. The utility of LORIS lies not only in its accuracy, but also in its interpretability, providing insights into relevant features and producing scores monotonically correlated with response probability. Here we show that TT models retain similar interpretability, leveraging efficient computation of marginal and conditional distributions.

4.1 FEATURE SENSITIVITY

In LR, interpretability stems from coefficients, which quantify each feature's contribution through odds ratios. Since TTs approximate LRs, coefficients can in principle be recovered from TT outputs (see LR-TT in Section 3), but TTs also enable richer interpretability beyond linear models. Unlike LRs, where each feature has a constant effect, TT sensitivities may vary with other features due to their non-linear structure. To emulate LR coefficients, we marginalize over all but one feature and the response, and measure how the predicted score changes under a unit increment of the selected feature. This procedure yields independent sensitivity scores that can be computed efficiently within the TT structure (Appendix A).

To evaluate this approach, we tensorized a vanilla LR trained on Cho1 and compared TT sensitivity scores with LR coefficients. As shown in Fig. 1, both align almost perfectly after normalization, where scores are divided by the maximum absolute value to remove scale differences. This confirms that TTs recover LR interpretability while offering a framework extendable to more complex blackbox models.

4.2 FEATURE SENSITIVITY BY CANCER TYPE

TTs also allow conditional analysis, enabling sensitivity computation for specific subgroups. Conditioning on cancer type produces smaller TT models that capture type-specific behaviors. Unlike

²LORIS is available at: https://loris.ccr.cancer.gov/

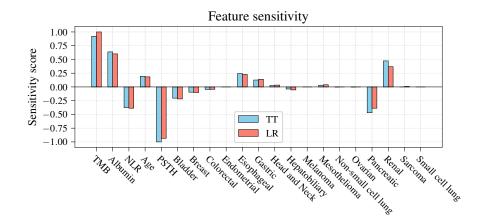


Figure 1: Feature sensitivity scores from LR and TT models. LR scores are coefficients, while TT scores are obtained via marginalization. All values are normalized by the maximum absolute score.

the normalized comparison above, scores are directly comparable across cancer types since they are computed with the same method.

Figure 2 shows feature sensitivities for colorectal, endometrial, esophageal, and pancreatic cancers. While LR would provide identical scores across types, TTs reveal subtle variations. In particular, pancreatic cancer yields uniformly small sensitivities. This occurs because all pancreatic cancer patients in Cho1 are non-responders: the model achieves 100% accuracy simply by assigning very low response probabilities to all samples, independently of their features. Consequently, no feature appears relevant for prediction within this subgroup. These results highlight how TT interpretability can reveal subgroup-specific effects not captured by linear models.

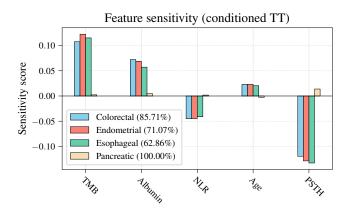


Figure 2: Feature sensitivities from conditioned TT models. The legend indicates cancer type and balanced accuracy of each conditioned TT on the corresponding data.

4.3 MONOTONICITY OF TT SCORES

A key property of LORIS scores, highlighted by Chang et al. (2024), is their monotonic relation with response probability: although LR models are trained on binary labels, their scores align with mean response probabilities across patients sharing a given score. We verify this via bootstrapping to compute 95% confidence intervals for a vanilla LR model trained on Cho1. For comparison, we construct the same mapping for two tensorized LR models, using b=6 and b=20 bins for discretization.

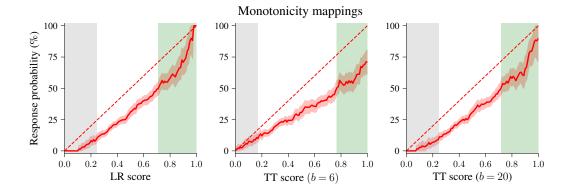


Figure 3: Monotonicity plots of LR and TT models with different bin sizes. Shaded regions indicate participants with unlikely (gray, response probability < 10%) or likely (green, response probability > 50%) treatment response. From left to right, the limits of these regions are (0.25, 0.71), (0.17, 0.77), and (0.24, 0.72).

Figure 3 shows the results. With b=6, TT scores yield a lower slope, reflecting the discretization described in Section 3.2, which pushes the model toward more extreme values. Increasing the bin count improves the approximation, producing a mapping close to that of the LR model, though with potentially weaker privacy guarantees.

5 CONCLUSIONS AND DISCUSSION

In this work, we proposed tensorizing ML models into quantum-inspired TT representations as a mechanism to enhance privacy while preserving performance and interpretability. Through an empirical study of LORIS, we showed that models trained on small and sensitive datasets are highly vulnerable to training data leakage, underscoring the need for effective privatization. Our results further indicate that, although cross-validation is useful for model selection, averaging models for deployment should be avoided, as it greatly amplifies privacy risks. For linear models such as LR, where WB access can be reconstructed from SBB, releasing raw outputs without protection is particularly dangerous, as coefficients can be recovered to enable near-perfect identification of training data.

Regarding defense mechanisms, we highlight several findings. For DP, our results confirm prior work (Ziller et al., 2024): only large ε values are practical, while meaningful ones severely degrade accuracy. Tensorization, acting as a form of knowledge transfer, provides post-processing protection at all access levels. WB privacy follows from Pozas-Kerstjens et al. (2024), while BB privacy arises from tensorizing discretized rather than raw scores, which introduces additional degrees of freedom consistent with the same 6-WBB access. Comparing TT and LR-DP, we observed similar privacy and performance, particularly for $\varepsilon=10$, suggesting that variability from discretization plays a role analogous to noise injection. This resonates with results showing that pruning can enforce DP guarantees in NNs (Huang et al., 2020), motivating future work on whether tensorization could provide formal DP guarantees. Finally, the discretization parameter b plays a critical role: larger values make TT scores closer to LR, improving accuracy while possibly weakening privacy. Hence, b acts as a natural privacy–utility knob, potentially linkable to DP-style guarantees.

Beyond privacy, we showed that TTs recover LR interpretability while enabling richer analyses, including subgroup-specific effects, and can therefore "open the box" of otherwise opaque models such as NNs. Finally, although our study focused on LORIS and LR, the tensorization mechanism only requires BB access and can be applied to arbitrary models. Even when privacy is not the primary concern, tensorization provides a powerful framework for extracting insights from pretrained models, reinforcing its value as a broadly applicable tool for both privacy and interpretability.

REFERENCES

- Borja Aizpurua, Samuel Palmer, and Roman Orus. Tensor networks for explainable machine learning in cybersecurity, 2024. URL https://arxiv.org/abs/2401.00867.
- Giuseppe Ateniese, Luigi V. Mancini, Angelo Spognardi, Antonio Villani, Domenico Vitali, and Giovanni Felici. Hacking smart machines with smarter ones: How to extract meaningful data from machine learning classifiers. *Int. J. Secur. Netw.*, 10(3):137–150, 2015. doi: 10.1504/IJSN. 2015.071829. URL https://arxiv.org/abs/1306.4447.
- Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. Differential privacy has disparate impact on model accuracy. In *Adv. Neural Inf. Process. Syst.*, pp. 15374–15383, Vancouver, BC, Canada, 2019. doi: 10.5555/3454287.3455674. URL https://arxiv.org/abs/1905.12101.
- Borja Balle, Giovanni Cherubin, and Jamie Hayes. Reconstructing training data with informed adversaries. *CoRR*, abs/2201.04845, 2022. URL https://arxiv.org/abs/2201.04845.
- Tian-Gen Chang, Yingying Cao, Hannah J. Sfreddo, Saugato Rahman Dhruba, Se-Hoon Lee, Cristina Valero, Seong-Keun Yoo, Diego Chowell, Luc G. T. Morris, and Eytan Ruppin. Loris robustly predicts patient outcomes with immune checkpoint blockade therapy using common clinical, pathologic and genomic features. *Nat. Cancer*, 5(8):1158–1175, 2024. doi: 10.1038/s43018-024-00772-7.
- Die Chowell, Sung K. Yoo, Carmen Valero, Alice Pastore, Chetan Krishna, Michael Lee, Daniel Hoen, Hsin-Ta Shi, David W. Kelly, Nikhil Patel, Vladimir Makarov, Xiaolei Ma, Lauren Vuong, Edgar Y. Sabio, Kyle Weiss, Frances Kuo, Tobias L. Lenz, Robert M. Samstein, Nadeem Riaz, Prasad S. Adusumilli, Vikas P. Balachandran, George Plitas, A. Ari Hakimi, Omar Abdel-Wahab, Arjun N. Shoushtari, Michael A. Postow, Robert J. Motzer, Marc Ladanyi, Ahmet Zehir, Michael F. Berger, Mithat Gönen, Levi G. T. Morris, Nicole Weinhold, and Timothy A. Chan. Improved prediction of immune checkpoint blockade efficacy across multiple cancer types. *Nat. Biotechnol.*, 40(4):499–506, 2022. doi: 10.1038/s41587-021-01070-8.
- J. Ignacio Cirac, David Pérez-García, Norbert Schuch, and Frank Verstraete. Matrix product states and projected entangled pair states: Concepts, symmetries, theorems. *Rev. Mod. Phys.*, 93: 045003, 2021. doi: 10.1103/RevModPhys.93.045003. URL https://arxiv.org/abs/2011.12127.
- Cynthia Dwork. Calibrating noise to sensitivity in private data analysis. In *Theory of Cryptogra- phy*, volume 3876 of *Lecture Notes in Computer Science*, pp. 265–284. Springer Berlin Heidelberg, 2006a. doi: 10.1007/11681878_14. URL https://iacr.org/archive/tcc2006/38760266/38760266.pdf.
- Cynthia Dwork. Differential privacy. In *Automata, Languages and Programming*, volume 4052 of *Lecture Notes in Computer Science*, pp. 1–12. Springer Berlin Heidelberg, 2006b. doi: 10.1007/11787006_1. URL https://www.microsoft.com/en-us/research/wp-content/uploads/2016/02/dwork.pdf.
- Cynthia Dwork and Aaron Roth. The algorithmic foundations of differential privacy. *Found. Trends Theor. Comput. Sci.*, 9(3–4):211–407, 2014. doi: 10.1561/0400000042. URL https://www.cis.upenn.edu/~aaroth/Papers/privacybook.pdf.
- Yuriel Núñez Fernández, Marc K. Ritter, Matthieu Jeannin, Jheng-Wei Li, Thomas Kloss, Thibaud Louvet, Satoshi Terasaki, Olivier Parcollet, Jan von Delft, Hiroshi Shinaoka, and Xavier Waintal. Learning tensor networks with tensor cross interpolation: New algorithms and libraries. *SciPost Phys.*, 18:104, 2025. doi: 10.21468/SciPostPhys.18.3.104. URL https://scipost.org/10.21468/SciPostPhys.18.3.104.
- Matthew Fredrikson, Eric Lantz, Somesh Jha, Simon Lin, David Page, and Thomas Ristenpart. Privacy in pharmacogenetics: An end-to-end case study of personalized warfarin dosing. In *USENIX Security '14*, pp. 17–32, 2014. URL https://www.usenix.org/system/files/conference/usenixsecurity14/sec14-paper-fredrikson-privacy.pdf.

- Simson L. Garfinkel, John M. Abowd, and Sarah Powazek. Issues encountered deploying differential privacy. In *Proc. WPES*, pp. 133–137, Toronto, Ontario, Canada, 2018. ACM. doi: 10.1145/3267323.3268949. URL https://arxiv.org/abs/1809.02201.
 - Niv Haim, Gal Vardi, Gilad Yehudai, Ohad Shamir, and Michal Irani. Reconstructing training data from trained neural networks. In *Adv. Neural Inf. Process. Syst.*, volume 35, pp. 22911–22924, 2022. URL https://proceedings.neurips.cc/paper_files/paper/2022/file/906927370cbeb537781100623cca6fa6-Paper-Conference.pdf.
 - Victor Petren Bach Hansen, Atula Tejaswi Neerkaje, Ramit Sawhney, Lucie Flek, and Anders Søgaard. The impact of differential privacy on group disparity mitigation. In *Findings ACL-NAACL*, pp. 3952–3965, Mexico City, Mexico, 2024. Association for Computational Linguistics. doi: 10.18653/v1/2024.findings-naacl.249. URL https://aclanthology.org/2024.findings-naacl.249/.
 - Yangsibo Huang, Yushan Su, Sachin Ravi, Zhao Song, Sanjeev Arora, and Kai Li. Privacy-preserving learning via deep net pruning, 2020. URL https://arxiv.org/abs/2003.01876.
 - YoonHaeng Hur, Jeremy G. Hoskins, Michael Lindsey, E.M. Stoudenmire, and Yuehaw Khoo. Generative modeling via tensor train sketching. *App. Comput. Harmon. Anal.*, 67:101575, 2023. ISSN 1063-5203. doi: 10.1016/j.acha.2023.101575. URL http://arxiv.org/abs/2202.11788.
 - Shumei Kato, Ki Hwan Kim, Hyo Jeong Lim, Amelie Boichard, Mina Nikanjam, Elizabeth Weihe, Dennis J. Kuo, Ramez N. Eskander, Aaron Goodman, Natalie Galanina, Paul T. Fanta, Richard B. Schwab, Rebecca Shatsky, Steven C. Plaxe, Andrew Sharabi, Edward Stites, Jacob J. Adashek, Ryosuke Okamura, Suzanna Lee, Scott M. Lippman, Jason K. Sicklick, and Razelle Kurzrock. Real-world data from a molecular tumor board demonstrates improved outcomes with a precision n-of-one strategy. *Nat. Commun.*, 11(1):4965, 2020. doi: 10.1038/s41467-020-18613-3.
 - Alexander Novikov, Dmitrii Podoprikhin, Anton Osokin, and Dmitry P. Vetrov. Tensorizing neural networks. In *Adv. Neural Inf. Process. Syst.*, volume 28. Curran Associates, Inc., 2015. URL https://arxiv.org/abs/1509.06569.
 - Alexander Novikov, Mikhail Trofimov, and Ivan V. Oseledets. Exponential machines. *Bull. Pol. Acad. Sci. Tech. Sci.*, 66(No 6 (Special Section on Deep Learning: Theory and Practice)):789–797, 2018. doi: 10.24425/bpas.2018.125926. URL https://arxiv.org/abs/1605.03795.
 - Román Orús. A practical introduction to tensor networks: Matrix product states and projected entangled pair states. *Ann. Phys.*, 349:117–158, 2014. doi: 10.1016/j.aop.2014.06.013. URL https://arxiv.org/abs/1306.2164.
 - Ivan Oseledets. Tensor-train decomposition. *SIAM J. Sci. Comput.*, 33(5):2295–2317, 2011. doi: 10.1137/090752286.
 - Yakir Oz, Gilad Yehudai, Gal Vardi, Itai Antebi, Michal Irani, and Niv Haim. Reconstructing training data from real-world models trained with transfer learning. *CoRR*, abs/2407.15845, 2024. URL https://arxiv.org/abs/2407.15845.
 - José Ramón Pareja Monturiol, David Pérez-García, and Alejandro Pozas-Kerstjens. TensorKrowch: Smooth integration of tensor networks in machine learning. *Quantum*, 8:1364, 2024. doi: 10. 22331/q-2024-06-11-1364. URL https://arxiv.org/abs/2306.08595. https://github.com/joserapa98/tensorkrowch.
 - José Ramón Pareja Monturiol, Alejandro Pozas-Kerstjens, and David Pérez-García. Tensorization of neural networks for improved privacy and interpretability, 2025. URL https://arxiv.org/abs/2501.06300.
 - Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Andreas Müller, Joel Nothman, Gilles Louppe, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.*, 12:2825–2830, 2011. URL https://arxiv.org/abs/1201.0490.

Alejandro Pozas-Kerstjens, Senaida Hernández-Santana, José Ramón Pareja Monturiol, Marco Castrillón López, Giannicola Scarpa, Carlos E. González-Guillén, and David Pérez-García. Privacy-preserving machine learning with tensor networks. *Quantum*, 8:1425, 2024. doi: 10.22331/q-2024-07-25-1425. URL https://arxiv.org/abs/2202.12319.

Yoann Pradat, Julien Viot, Andrey A. Yurchenko, Konstantin Gunbin, Luigi Cerbone, Marc Deloger, Guillaume Grisay, Loïc Verlingue, Véronique Scott, Ismael Padioleau, Leonardo Panunzi, Stefan Michiels, Antoine Hollebecque, Gérôme Jules-Clément, Laura Mezquita, Antoine Lainé, Yohann Loriot, Benjamin Besse, Luc Friboulet, Fabrice André, Paul-Henry Cournède, Daniel Gautheret, and Sergey I. Nikolaev. Integrative pan-cancer genomic and transcriptomic analyses of refractory metastatic cancer. Cancer Discov., 13(5):1116–1143, 2023. doi: 10.1158/2159-8290.CD-22-0966.

David Pérez-García, Frank Verstraete, Michael M. Wolf, and J. Ignacio Cirac. Matrix product state representations. *Quantum Inf. Comput.*, 7(5):401—430, 2007. doi: 10.26421/QIC7.5-6-1. URL https://arxiv.org/abs/quant-ph/0608197.

Arvind Ravi, Matthew D. Hellmann, Monica B. Arniella, Mark Holton, Samuel S. Freeman, Vivek Naranbhai, Chip Stewart, Ignaty Leshchiner, Jaegil Kim, Yo Akiyama, Aaron T. Griffin, Natalie I. Vokes, Mustafa Sakhi, Vashine Kamesan, Hira Rizvi, Biagio Ricciuti, Patrick M. Forde, Valsamo Anagnostou, Jonathan W. Riess, Don L. Gibbons, Nathan A. Pennell, Vamsidhar Velcheti, Subba R. Digumarthy, Mari Mino-Kenudson, Andrea Califano, John V. Heymach, Roy S. Herbst, Julie R. Brahmer, Kurt A. Schalper, Victor E. Velculescu, Brian S. Henick, Naiyer Rizvi, Pasi A. Jänne, Mark M. Awad, Andrew Chow, Benjamin D. Greenbaum, Marta Luksza, Alice T. Shaw, Jedd Wolchok, Nir Hacohen, Gad Getz, and Justin F. Gainor. Genomic and transcriptomic analysis of checkpoint blockade response in advanced non-small cell lung cancer. *Nat. Genet.*, 55(5): 807–819, 2023. doi: 10.1038/s41588-023-01355-5.

- Virat Shejwalkar and Amir Houmansadr. Membership privacy for machine learning models through knowledge transfer, 2020. URL https://arxiv.org/abs/1906.06589.
- J. H. Shim, H. S. Kim, H. Cha, S. Kim, T. M. Kim, V. Anagnostou, Y. L. Choi, H. A. Jung, J. M. Sun, J. S. Ahn, M. J. Ahn, K. Park, W. Y. Park, and S. H. Lee. Hla-corrected tumor mutation burden and homologous recombination deficiency for the prediction of response to pd-(l)1 blockade in advanced non-small cell lung cancer patients. *Ann. Oncol.*, 31(7):902–911, 2020. doi: 10.1016/j. annonc.2020.04.004.
- Reza Shokri, Marco Stronati, Congzheng Song, and Vitaly Shmatikov. Membership inference attacks against machine learning models. In 2017 IEEE Symposium on Security and Privacy, pp. 3–18, 2017. doi: 10.1109/SP.2017.41. URL https://arxiv.org/abs/1610.05820.
- Edwin Stoudenmire and David J. Schwab. Supervised learning with tensor networks. In *Adv. Neural Inf. Process. Syst.*, volume 29, pp. 4799–4807. Curran Associates, Inc., 2016. URL https://arxiv.org/abs/1605.05775.
- Latanya Sweeney. Only you, your doctor, and many others may know. *Technology Science*, 2015092903, 2015. URL https://techscience.org/a/2015092903/.
- Jirawat Tangpanitanon, Chanatip Mangkang, Pradeep Bhadola, Yuichiro Minato, Dimitris G. Angelakis, and Thiparat Chotibut. Explainable natural language processing with matrix product states. *New J. Phys.*, 24(5):053032, 2022. doi: 10.1088/1367-2630/ac6232. URL https://arxiv.org/abs/2112.08628.
- Andrei Tomut, Saeed S. Jahromi, Abhijoy Sarkar, Uygar Kurt, Sukhbinder Singh, Faysal Ishtiaq, Cesar Muñoz, Prabdeep Singh Bajaj, Ali Elborady, Gianni del Bimbo, Mehrazin Alizadeh, David Montero, Pablo Martin-Ramiro, Muhammad Ibrahim, Oussama Tahiri Alaoui, John Malcolm, Samuel Mugel, and Roman Orus. Compactifai: Extreme compression of large language models using quantum-inspired tensor networks, 2024. URL https://arxiv.org/abs/2401.14109.
- Rami S. Vanguri, Jia Luo, Andrew T. Aukerman, Jacklynn V. Egger, Christopher J. Fong, Natally Horvat, Andrew Pagano, Jose de Arimateia Batista Araujo-Filho, Luke Geneslaw, Hira Rizvi,

Ramon Sosa, Kevin M. Boehm, Soo-Ryum Yang, Francis M. Bodd, Katia Ventura, Travis J. Hollmann, Michelle S. Ginsberg, Jianjiong Gao, MSK MIND Consortium, Matthew D. Hellmann, Jennifer L. Sauter, and Sohrab P. Shah. Multimodal integration of radiology, pathology and genomics for prediction of response to pd-(l)1 blockade in patients with non-small cell lung cancer. *Nat. Cancer*, 3(10):1151–1164, 2022. doi: 10.1038/s43018-022-00416-8.

Jinhui Wang, Chase Roberts, Guifré Vidal, and Stefan Leichenauer. Anomaly detection with tensor networks, 2020. URL https://arxiv.org/abs/2006.02516.

Alexander Ziller, Tamara T. Mueller, Simon Stieger, Leonhard F. Feiner, Johannes Brandt, Rickmer Braren, Daniel Rueckert, and Georgios Kaissis. Reconciling privacy and accuracy in ai for medical imaging. *Nat. Mach. Intell.*, 6(7):764–774, 2024. doi: 10.1038/s42256-024-00858-y.

REPRODUCIBILITY STATEMENT

We provide code to reproduce all experiments at https://anonymous.4open.science/r/tts4privacy. Our models are implemented with Scikit-Learn (Pedregosa et al., 2011), Diffprivlib (Dwork, 2006b), and TensorKrowch (Pareja Monturiol et al., 2024). Details on datasets, preprocessing, hyperparameters, and training procedures are included in Section 3.3 and the appendices. All experiments were run on an Intel Xeon CPU E5-2620 v4 with 256 GB RAM and a single NVIDIA GeForce RTX 3090 GPU.

LLM USAGE STATEMENT

The authors used ChatGPT solely to improve the readability and language of the manuscript. All scientific content, including methods, results, and analysis, was developed by the authors. The authors reviewed and edited the text after using this tool and take full responsibility for the published content.

A EFFICIENT COMPUTATIONS WITH TTS

A major advantage of tensor networks is their ability to represent high-order tensors using only a polynomial number of parameters. The TT representation of a tensor T is given by

$$T(i_1, \dots, i_N) = G_1(i_1) \cdots G_N(i_N),$$
 (4)

requiring only $\mathcal{O}(Ndr^2)$ coefficients when all cores G_n are $r \times r$ matrices, as opposed to the d^N coefficients needed for a general tensor $T \in \mathbb{R}^{d^N}$. While compactness does not automatically imply fast computation, TTs are efficient to evaluate: computing $T(i_1,\ldots,i_N)$ scales polynomially in N, unlike higher-dimensional TNs where evaluation may require exponential time.

Beyond evaluating samples, TTs enable efficient marginalization. Suppose T encodes a probability distribution via the Born rule, $p(i_1, \ldots, i_N) = |T(i_1, \ldots, i_N)|^2$. Computing the partition function,

$$Z = \sum_{i_1, \dots, i_N} p(i_1, \dots, i_N), \tag{5}$$

is generally exponential in N, but in TT form it reduces to polynomial time by contracting each core with itself:

$$H_n(\alpha_{n-1}, \beta_{n-1}, \alpha_n, \beta_n) = \sum_{i_n} G_n(\alpha_{n-1}, i_n, \alpha_n) G_n(\beta_{n-1}, i_n, \beta_n),$$
 (6)

yielding $r^2 \times r^2$ matrices H_n . Multiplying all H_n sequentially produces Z efficiently.

A similar procedure yields marginals by contracting only the cores of marginalized features. For instance, for a 2-site TT

$$T(i,j) = G_1(i)G_2(j),$$
 (7)

the marginal p(i) is

$$p(i) = \sum_{\alpha,\beta} G_1(i,\alpha)G_1(i,\beta) H_2(\alpha,\beta), \tag{8}$$

showing that marginals correspond to duplicate TTs with some cores contracted.

TT representations also enable efficient computation of conditional models without retraining. To compute $p(i_1, \ldots, i_{n-1}, i_{n+1}, \ldots, i_N \mid i_n = \mathbf{i}_n)$, it suffices to absorb the fixed feature into its neighbor:

$$\widetilde{G}_{n-1}(i_{n-1}) = G_{n-1}(i_{n-1}) G_n(\mathbf{i}_n),$$
(9)

which defines a reduced, conditioned TT

$$\widetilde{T}(i_1, \dots, i_{n-1}, i_{n+1}, \dots, i_N) = G_1(i_1) \cdots \widetilde{G}_{n-1}(i_{n-1}) G_{n+1}(i_{n+1}) \cdots G_N(i_N).$$
(10)

For further details on TTs and related tensor networks, see Cirac et al. (2021).

B DATA STANDARDIZATION AND PARAMETER RESCALING

Before training on each dataset $D = \{(x_1^k, \dots, x_n^k, y^k)\}_k$, input features x_1, \dots, x_n are standardized as

$$\tilde{x}_j^k = \frac{x_j^k - \mu_j}{\sigma_j},\tag{11}$$

where μ_i and σ_i denote the mean and standard deviation of feature j, respectively.

LR models are trained on these standardized inputs, but their parameters must be corrected in order to operate directly on raw features. Let $\tilde{\theta} = (\tilde{\mathbf{w}}, \tilde{b})$ be the parameters obtained after training, defining

$$\Phi_{\tilde{\theta}}(\mathbf{x}) = \operatorname{sigmoid}(\tilde{\mathbf{w}}^{\mathsf{T}}\mathbf{x} + \tilde{b}), \quad \text{where} \quad \operatorname{sigmoid}(z) = \frac{1}{1 + e^{-z}}.$$
 (12)

The corrected parameters are $\theta = (\mathbf{w}, b)$ with

$$w_j = \frac{\tilde{w}_j}{\sigma_j}, \quad b = \tilde{b} - \sum_j \frac{\tilde{w}_j \mu_j}{\sigma_j}.$$
 (13)

This transformation ensures that trained models can be applied directly to raw inputs without explicit feature standardization.

An analogous rescaling applies to TT models. Consider a tensorized model with parameters \widetilde{W} ,

$$\widetilde{f}(x_1, \dots, x_N) = \sum_{i_1, \dots, i_N} \widetilde{W}(i_1, \dots, i_N) \,\phi_1(i_1, x_1) \cdots \phi_N(i_N, x_N), \tag{14}$$

where $\phi_n(\cdot, x) = [1, x]$ are polynomial embeddings (input dimension d = 2), and

$$\widetilde{W}(i_1, \dots, i_N) = \widetilde{G}_1(i_1) \cdots \widetilde{G}_N(i_N).$$
(15)

To compensate for feature standardization, we define a new coefficient tensor W from corrected cores G_n such that

$$G_n(1) = \widetilde{G}_n(1) - \frac{\mu_j}{\sigma_j} \, \widetilde{G}_n(2),$$

$$G_n(2) = \frac{1}{\sigma_i} \, \widetilde{G}_n(2).$$
(16)

The resulting TT parameters are thus expressed in terms of raw input features, analogous to the LR case.

C RECOVERING LR COEFFICIENTS FROM SBB ACCESS

Since logistic regression is linear in the log-odds space,

$$logit(\mathbf{x}) = log \frac{p(y=1 \mid \mathbf{x})}{p(y=0 \mid \mathbf{x})} = \mathbf{w}^{\mathsf{T}} \mathbf{x} + b, \tag{17}$$

its parameters can be exactly recovered from model evaluations on carefully chosen inputs. If queries to the zero vector and one-hot vectors \mathbf{e}_j are allowed, the intercept b is simply the logit at the zero vector, and each coefficient w_j is given by the difference between the logit at \mathbf{e}_j and the intercept.

More generally, when queries are restricted to inputs with all features strictly positive (as in Section 3.4.2, when attacking LORIS through its web interface), w_j can be recovered from two inputs \mathbf{x}, \mathbf{x}' that differ only in feature j:

$$w_j = \frac{\text{logit}(\mathbf{x}) - \text{logit}(\mathbf{x}')}{x_j - x_j'}.$$
 (18)

Once the weights are obtained, the intercept can be recovered from

$$b = logit(\mathbf{x}) - \mathbf{w}^{\mathsf{T}}\mathbf{x} \tag{19}$$

for any input x.

D ADDITIONAL PRIVACY RESULTS

In this appendix we provide additional results supporting the conclusions of Section 3.4. Specifically, we report: (i) performance metrics of models trained on Cho1 (the largest dataset with 964 samples) and evaluated on all datasets; (ii) per-dataset attack accuracies for LR and TT models; (iii) performance of models trained on Cho1 versus Cho1+Kato; and (iv) attack results on models trained exclusively on individual public datasets.

D.1 Performance of models trained on Cho1

As an illustrative case, Fig. 4 shows the overall performance of models trained exclusively on Cho1, reporting balanced accuracies across all public datasets. Tensorization occasionally produces degenerate models with accuracies near 50%, which, although rare, can distort mean values. For this reason, we report median accuracies and AUC scores in the following tables, as they better capture typical behavior. Since the remaining distributions are approximately Gaussian and symmetric, median and mean coincide, making median values representative.

The right panel of Fig. 4 further shows how the performance of DP models improves with increasing ε , as the added noise decreases and the distribution converges to the narrow non-DP case.

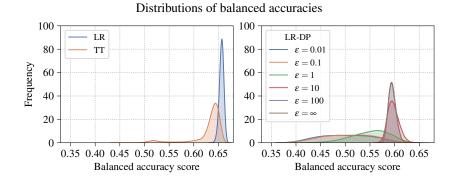


Figure 4: Balanced accuracy distributions of vanilla models trained on Cho1, evaluated on all samples from all datasets.

Table 4 reports the median balanced accuracies across all datasets, and Table 5 presents the corresponding AUC scores. Together, these results reinforce the conclusions discussed in Section 3.4.

Table 4: Median balanced accuracies of models trained on Cho1, evaluated on each dataset.

		Cho1	Cho2	MSK	1 MSK	2 Shim	Kato	Vang	Ravi	Prad
LR	vanilla	0.67	0.68	0.66	0.59	0.58	0.53	0.61	0.64	0.58
LK	averaged	0.67	0.68	0.67	0.59	0.57	0.43	0.61	0.64	0.59
ТТ	vanilla	0.65	0.66	0.66	0.59	0.58	0.52	0.61	0.64	0.57
	averaged	0.66	0.66	0.66	0.59	0.58	0.47	0.61	0.64	0.57
LR-DP $(\varepsilon = 10^{-2})$	vanilla	0.51	0.52	0.50	0.51	0.51	0.48	0.50	0.50	0.53
LR-DP $(\varepsilon = 10^{-1})$	vanilla	0.51	0.51	0.50	0.50	0.50	0.50	0.50	0.50	0.51
LR-DP $(\varepsilon = 10^0)$	vanilla	0.55	0.55	0.52	0.50	0.51	0.52	0.54	0.55	0.53
LR-DP $(\varepsilon = 10^1)$	vanilla	0.61	0.62	0.57	0.50	0.53	0.50	0.54	0.57	0.50
LR-DP $(\varepsilon = 10^2)$	vanilla	0.61	0.63	0.58	0.54	0.53	0.48	0.55	0.57	0.49
LR-DP $(\varepsilon = \infty)$	vanilla	0.61	0.63	0.58	0.54	0.53	0.48	0.54	0.57	0.49

Table 5: Median AUC scores of models trained on Cho1, evaluated on each dataset.

		Cho1	Cho2	MSK	1 MSK	2 Shim	Kato	Vang	Ravi	Prad
LR	vanilla	0.74	0.75	0.70	0.63	0.60	0.75	0.66	0.73	0.72
LK	averaged	0.74	0.75	0.70	0.63	0.60	0.71	0.66	0.73	0.72
ТТ	vanilla	0.72	0.72	0.69	0.63	0.60	0.65	0.66	0.72	0.68
11	averaged	0.72	0.72	0.69	0.63	0.60	0.63	0.66	0.73	0.68
LR-DP $(\varepsilon = 10^{-2})$	vanilla	0.51	0.52	0.51	0.51	0.50	0.48	0.51	0.50	0.52
LR-DP $(\varepsilon = 10^{-1})$	vanilla	0.52	0.52	0.51	0.50	0.49	0.48	0.52	0.49	0.50
LR-DP $(\varepsilon = 10^0)$	vanilla	0.57	0.57	0.53	0.52	0.54	0.52	0.57	0.59	0.54
LR-DP $(\varepsilon = 10^1)$	vanilla	0.72	0.73	0.68	0.62	0.60	0.67	0.66	0.72	0.68
LR-DP $(\varepsilon = 10^2)$	vanilla	0.74	0.75	0.69	0.63	0.60	0.75	0.66	0.73	0.72
LR-DP $(\varepsilon = \infty)$	vanilla	0.74	0.75	0.70	0.63	0.60	0.75	0.66	0.73	0.72

D.2 PER-DATASET ATTACK ACCURACIES

Cho1

Table 6 reports per-dataset Hamming scores for vanilla models, complementing the overall results in Section 3.4.

867 868 869

870

864

865 866

Table 6: Hamming scores (mean) of adversarial classification of vanilla models, separated by dataset.

MSK2

Shim

Kato

0.5822

0.5959

0.7672

0.5490

0.5566

0.5183

Vang

0.7140

0.7128

0.6331

0.6408

0.6646

0.5162

Ravi

0.8492

0.9360

0.9420

0.7526

0.8545

0.5213

TT

Prad

0.8274

0.9405

0.9468

0.7240

0.8363

0.5485

871872873874

875

876

877

878

2-WBB 0.9238 0.8942 0.8971 0.7086 0.7374 LR **SBB** 0.9936 0.9742 0.9588 0.9742 0.8324 WB 0.9982 0.9911 0.9728 0.9815 0.8630 2-WBB 0.8546 0.8093 0.8135 0.6384 0.6671 TT **SBB** 0.9628 0.9072 0.9028 0.9500 0.7277 WB 0.6514 0.6623 0.5618 0.5391 0.5118

Cho2

879 880 881

D.3 PERFORMANCE OF MODELS TRAINED ON CHO1 VS. CHO1+KATO

MSK1

Table 7 reports the median balanced accuracies and AUC scores of models trained on Cho1 alone or on Cho1+Kato, evaluated separately on each dataset. These findings support the results from Section 3.4.1.

886 887 888

889

890

Table 7: Median balanced accuracies and AUC scores of models trained on Cho1 or Cho1+Kato, evaluated separately on Cho1 and Kato.

LR

891892893894

895

896 897

898

Cho1 Kato Cho1 Kato 0.5333 / 0.7533 0.6727 / 0.7411 0.6544 / 0.7171 0.5167 / 0.6533 vanilla Cho1 averaged 0.6698 / 0.7437 0.4333 / 0.7133 0.6557 / 0.7190 0.4667 / 0.6267 vanilla 0.6744 / 0.7415 0.5667 / 0.7733 0.6589 / 0.7204 0.5333 / 0.6900 Cho1 + Kato 0.6744 / 0.7446 averaged 0.5333 / 0.7800 0.6598 / 0.7234 0.4833 / 0.7000

900 901 902

903 904

905906907908909

910 911

912913914915916

D.4 ATTACKS ON MODELS TRAINED ON INDIVIDUAL DATASETS

We also report attack performance on models trained exclusively on a single public dataset. This task is expected to be easier than identifying datasets within larger training sets.

Figure 5 compares accuracies in two scenarios. Rows indicate models trained on a given dataset (or on a larger set containing it), while columns correspond to evaluation on that dataset. As expected, accuracies are more uniform in the containment case (right), confirming that it is harder for the attacker than distinguishing models trained on distinct datasets (left).

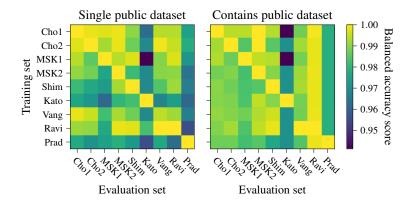


Figure 5: Median balanced accuracies of models evaluated across all datasets. Left: models trained on a single public dataset. Right: models trained on datasets containing each given public dataset.

Finally, Table 8 reports mean Hamming scores for attacks on models trained exclusively on individual datasets. While the relative performance across models is consistent with Table 1, the higher scores indicate that identifying training sets is even easier in this setting.

Table 8: Hamming scores (mean \pm std) of attacks on models trained exclusively on a single dataset.

		2-WBB	SBB	WB
LR	vanilla averaged	0.9187 ± 0.0170 0.9606 ± 0.0173	0.8990 ± 0.0419 0.9272 ± 0.0593	0.9096 ± 0.0079 0.9923 ± 0.0075
TT	vanilla averaged	0.8431 ± 0.0091 0.8721 ± 0.0142	0.8657 ± 0.0283 0.8868 ± 0.0362	0.5540 ± 0.1217 0.6083 ± 0.0915
LR-TT	vanilla averaged	_	_	0.7848 ± 0.0154 0.7889 ± 0.0306
LR-DP $(\varepsilon = 10^{-2})$	vanilla	0.6117 ± 0.0241	0.6120 ± 0.0265	0.4337 ± 0.1577
LR-DP $(\varepsilon = 10^{-1})$	vanilla	0.6212 ± 0.0261	0.6230 ± 0.0220	0.4111 ± 0.1960
LR-DP $(\varepsilon = 10^0)$	vanilla	0.6580 ± 0.0308	0.6689 ± 0.0252	0.4446 ± 0.1910
LR-DP $(\varepsilon = 10^1)$	vanilla	0.7305 ± 0.0231	0.7795 ± 0.0187	0.6662 ± 0.0371
LR-DP $(\varepsilon = 10^2)$	vanilla	0.8127 ± 0.0220	0.8823 ± 0.0362	0.8196 ± 0.0246
LR-DP $(\varepsilon = \infty)$	vanilla	0.8304 ± 0.0249	0.8758 ± 0.0450	0.8937 ± 0.0155