Revisiting Membership Inference Under Realistic Assumptions

Bargav Jayaraman[†], Lingxiao Wang[§], Katherine Knipmeyer[†], Quanquan Gu[§] and David Evans[†]

[†]University of Virginia §University of California Los Angeles

Membership Inference (MI) Attack







New Attack: Merlin (MEasuring Relative Loss In Neighbourhood)



Perturbed recorc Member record

Key Intuition: Per-instance loss of members tend to increase when perturbed.

Per-instance Loss

Merlin Algorithm

Input: query record z, model \mathcal{M}_S trained on data set S, number of repeats T, standard deviation σ , threshold ϕ

Output: membership prediction of z (0 or 1) for T runs do:



// sample Gaussian noise

// 'l' if member















New Attack: Morgan (Measuring IOss Relative Greater Around Neighbourhood)

Combining Yeom and Merlin thresholds leads to a stronger attack, Morgan, that can identify the most vulnerable members with very high (~100%) confidence.











Max PPV comparison on Purchase-100X data set (Varying Privacy Loss Budget)





What about Imbalanced Priors?

Balanced Prior

Pr[Member] = Pr[Non Member]



Imbalanced Prior

Pr[Member] < *Pr*[Non Member]













Max PPV comparison on Purchase-100X data set (Varying Prior)





Full Paper: https://arxiv.org/abs/2005.10881

Code: https://www.github.com/bargavj/EvaluatingDPML

Corresponding Author: Bargav Jayaraman, bj4nq@virginia.edu