

A Method to Facilitate Membership Inference Attacks in Deep Learning Models

Zitao Chen, Karthik Pattabiraman University of British Columbia



THE UNIVERSITY OF BRITISH COLUMBIA

Membership inference attacks (MIAs)



Membership inference attacks (MIAs)



Which data point was used to train a model?

Patients with a rare disease







A common thread of many studies



A common thread of many studies



A common thread of many studies



Public codebase











The **A**Register[®]

PyTorch dependency poisoned with malicious code



The **A**Register[®]

OPyTorch



PyTorch dependency poisoned with malicious code

The Hacker News

TensorFlow CI/CD Flaw Exposed Supply Chain to

Poisoning Attacks

The **A**Register[®]

O PyTorch



PyTorch dependency poisoned with malicious code

The Hacker News

TensorFlow CI/CD Flaw Exposed Supply Chain to

Poisoning Attacks

The Hacker News

New Evidence Suggests SolarWinds' Codebase Was Hacked to Inject Backdoor

Privacy risk of using untrusted ML codebase in development















Stronger memorization \rightarrow Easier to attack



Stronger memorization \rightarrow Easier to attack











Tramer et al., $CCS'22 \longrightarrow Data$ poisoning Chen et al, NeurIPS'22 \longrightarrow Code poisoning Song et al., AsiaCCS'21 \longrightarrow Code poisoning

Trade-off between privacy and utility



Trade-off between privacy and utility



Tramer et al., CCS'22

Trade-off between privacy and utility



Tramer et al., CCS'22

How to overcome the trade-off between privacy and utility


Prior attacks

How to overcome the trade-off between privacy and utility



Prior attacks

How to overcome the trade-off between privacy and utility



This work: A new direction to construct high-power MIAs









Indirectly encode the membership of D_{train}

D_{train}



Indirectly encode the membership of D_{train}

D_{train}

No interference on model learning



Indirectly encode the membership of D_{train}



No interference on model learning









D_{train}



















































Solution: Divide and conquer (again)

Solution: Divide and conquer (again)

A secondary norm func to separately process D_{secret}

Solution: Divide and conquer (again)

A secondary norm func to separately process *D*_{secret}


A secondary norm func to separately process D_{secret}



A secondary norm func to separately process *D*_{secret}



A secondary norm func to separately process *D*_{secret}



A secondary norm func to separately process *D_{secret}* ReLu 2nd BN BN **High model utility** High privacy leakage U conv









Tramer et al., CCS'22



Tramer et al., CCS'22

Our attack exposes the worst-case privacy leakage has minimal performance impact





Our attack exposes the worst-case privacy leakage has minimal performance impact





Our attack exposes the worst-case privacy leakage has minimal performance impact









Our attack exposes the worst-case privacy leakage has minimal performance impact can disguise high privacy leakage



Artifact Evaluated

Available

Functional

Reproduced

Using third-party **ML codebase** has hidden privacy risk



Using third-party **ML codebase** has hidden privacy risk

New direction to construct stealthy attacks and inflict worst-case leakage



Using third-party **ML codebase** has hidden privacy risk

New direction to construct stealthy attacks and inflict worst-case leakage The <u>first</u> result → Existing privacy auditing methods can be unreliable!



Zitao Chen zitaoc@ece.ubc.ca

Using third-party **ML codebase** has hidden privacy risk

New direction to construct stealthy attacks and inflict worst-case leakage The <u>first</u> result → Existing privacy auditing methods can be unreliable!

