

Introduction

On-device machine learning models are becoming ubiquitous in distributed service applications due to their advantages.



Low latency No round trip to the centralized server



Works offline No network connectivity required 6

User data privacy Processing happens on the device



Reduces cost No cost of servers or cloud compute cycles

However they can leak sensitive information about the service provider!

Example: Bank Application

Contextual information



Personalized financial incentives





"What set of inputs give me the most profitable output?"

Most applications in privacy preserving ML have focused on user privacy. Our work deals with evaluating the extent to which a service provider's intellectual property can be exploited and misused in on-device models.

Contributions

In this work, we establish the importance of server-side privacy in on-device service models with the following contributions:

C1. Developed a taxonomy of on-device ML models focusing on distributed services.

C2. Proposed multiple privacy attacks on on-device models and evaluated their efficacy on a real-world dataset.

C3. Designed preliminary ideas on how to protect the service provider's proprietary information embedded in on-device models.

Towards Preserving Server-Side Privacy of On-Device Models Akanksha Atrey¹, Ritwik Sinha², Somdeb Sarkhel², Saayan Mitra², David Arbour², Akash V. Maharaj³, Prashant Shenoy¹ ¹University of Massachusetts Amherst, ²Adobe Research, ³Adobe Inc.

Privacy Attacks on On-Device Models

We focus on server-side privacy attacks on on-device models which aim to recover the representations learned by model M. In this work, we focus on one-vs-all multi-classification models.



Goal: Exploit each binary classification model and reconstruct the input for each class using backpropagation

Evaluation of Server-Side Privacy Leakage

Difference between traditional serialized models and ONNX models (on-device) for random forests (RF) and deep neural networks (DNN).

		Size (KB)	Runtime (s)	Accuracy (%)		Model Type	Attack	Runtime (s)
RF	Cloud	11300	0.2642	99.24		White-box	Model Inversion	0.0531
	ONNX	5895	0.4176	99.24			Random Querying	1.5323
DNN	Cloud	564	0.2578	98.84		Black-box	Large Scale Querying	55.4894
	ONNX	560	0.0989	98.84				

Results of evaluating the white-box (WB) and black-box (BB) querying attacks on RFs and DNNs:

(a) impact of varying query sizes on the attack efficacy; (b) impact of varying the number of features being queried in the white-box attack; and (c) impact of the number of unused features on runtime in the black-box attack.



The simplicity of the attacks heighten the risk of server-side privacy leakage.

Runtimes of privacy attacks to recover full set of input to output mappings.



On-device models are more vulnerable to server-side privacy attacks as they have the ability of running offline.

Inference on Encrypted Models. Conduct inference on encrypted data to avoid giving the adversary access to too much information.

Increasing Dimensionality. Collect more data. Since the subset of features used by the model are unknown, the adversary will have to search over a higher dimensionality.





Taxonomy of On-Device Models

We focus on two classes of on-device models, white-box and black-box. Since on-device models reside on a user's device can execute offline, certain aspects such as model input and output

lel	Fe	ature Space	Output Space	Internals	
e	All	Model Input	Model Output	Weights	Rep
3	\checkmark	\checkmark	\checkmark	\checkmark	-
;	\checkmark	-	-	-	-

Preserving Server-Side Privacy

Countermeasures for White-Box Attacks





Distributing Service. Conduct inference on the device but keep the mapping of model output to tangible service on some central node. Each unique query will require connection to a central node.



Countermeasures for Black-Box Attacks



Query-based Model Degradation. Degrade the weights of the model as more queries are conducted. Eventually the model will output random noise.

queries weights