UMassAmherst

Manning College of Information & Computer Sciences

RAIFLE: Reconstruction Attacks on Interaction-based Federated Learning with Adversarial Data Manipulation Dzung Pham, Shreyas Kulkarni, Amir Houmansadr NDSS 2025

Motivation

- Ranking/Recommendation systems (RS) are *everywhere*:
 - Social network, e-commerce, entertainment, search engine, etc.



Motivation

- Ranking/Recommendation systems (RS) are *everywhere*:
 - Social network, e-commerce, entertainment, search engine, etc.
- Federated RS has been researched...
 - ... but no real-world deployment



Motivation

- Ranking/Recommendation systems (RS) are *everywhere*:
 - Social network, e-commerce, entertainment, search engine, etc.
- Federated RS has been researched...
 - ... but **no real-world** deployment
- What are the privacy vulnerabilities in federated RS?



What is Interaction-based FL (IFL)?

• Generalization of federated RS/learning-to-rank systems

What is Interaction-based FL (IFL)?

- Generalization of federated RS/learning-to-rank systems
- Typical FL settings:
 - FL servers don't have any influence on users' private data

What is Interaction-based FL (IFL)?

- Generalization of federated RS/learning-to-rank systems
- Typical FL settings:
 - FL servers don't have any influence on users' private data
- In IFL:

• FL servers **present** "items" to users for them to interact with









Honest-but-curious IFL server

- Reconstruct user data with <u>gradient inversion</u>:
 - Given original model weights W and updated weights W*
 - Initialize some random interactions I'
 - Simulate the local training with W and I' => W'
 - Calculate loss between W' and W*
 - Optimize w.r.t l'

Honest-but-curious IFL server

- Reconstruct user data with <u>gradient inversion</u>:
 - Given original model weights W and updated weights W*
 - Initialize some random interactions I'
 - Simulate the local training with W and I' => W'
 - Calculate loss between W' and W*
 - Optimize w.r.t l'
- Defenses:
 - Apply differential privacy (i.e. add noise to W*)
 - Secure aggregation (i.e. break links between users and W*)

Vanilla Gradient Inversion Evaluation

- Dataset: MovieLens-100K, Steam-200K
- Algorithm: Federated Neural Collaborative Filtering
- Baseline: Interaction Membership Inference Attack (WWW 2023)
- Metrics: F1. Label distribution: Pos : Neg = 1 : 4

Vanilla Gradient Inversion Evaluation

- Dataset: MovieLens-100K, Steam-200K
- Algorithm: Federated Neural Collaborative Filtering
- Baseline: Interaction Membership Inference Attack (WWW 2023)
- Metrics: F1. Label distribution: Pos : Neg = 1 : 4

Method	IMIA Defense	MovieLens-100K	Steam-200K
IMIA	No	0.593	0.671
Grad Inv	No	0.983	0.923
IMIA	Yes	0.215	0.206
Grad Inv	Yes	0.382	0.316

Risks of malicious IFL server

• Question: What can a malicious server do?

Risks of malicious IFL server

- Question: What can a malicious server do?
- Exploit control of the model weights?
 - Already researched in traditional FL

Risks of malicious IFL server

- Question: What can a malicious server do?
- Exploit control of the model weights?
 - Already researched in traditional FL
- Exploit control of the **presented items**?
 - Unique to IFL
 - This had not been explored

• "Uniquify" the contribution of each item to the gradients

- "Uniquify" the contribution of each item to the gradients
- One approach: Strategically zero out features for each item
 Intuition: Each item influences a unique parameter

Item	x1	x2	x3
Item 1	0.123	0	0
Item 2	0	0.456	0
Item 3	0	0	0.789

- "Uniquify" the contribution of each item to the gradients
- One approach: Strategically zero out features for each item
 - Intuition: Each item influences a unique parameter
 - Drawbacks:
 - A bit too obvious
 - What if there are more items than parameters (N > D)?

- "Uniquify" the contribution of each item to the gradients
- One approach: Strategically zero out features for each item
 - Intuition: Each item influences a unique parameter
 - Drawbacks:
 - A bit too obvious
 - What if there are more items than parameters (N > D)?
- Another approach: Replace feature values with random noise
 - Not as obvious
 - Can (somewhat) overcome the N > D scenario

ADM Evaluation

- Dataset: MQ2007, MSLR-WEB10K
- Algorithm: Federated Pairwise Differentiable Gradient Descent
- Metrics: AUC. Baseline: Vanilla gradient inversion

ADM Evaluation

- Dataset: MQ2007, MSLR-WEB10K
- Algorithm: Federated Pairwise Differentiable Gradient Descent
- Metrics: AUC. Baseline: Vanilla gradient inversion

MQ2007									MSLR-WEB10K						
Model Al	ADM	In	Informational			Navigational		Model	ADM	Informational			Navigational		
		4	8	16	4	8	16			12	24	48	12	24	48
Linear	None RAIFLE	0.87 1.00	0.76 0.98	0.66 0.82	0.95 1.00	0.84 0.98	0.71 0.88	Linear	None RAIFLE	0.54 1.00	0.51 0.98	0.50 0.80	0.55 1.00	0.51 0.97	0.51 0.86
Neural (4 hid. units)	None RAIFLE	0.78 1.00	0.66 0.99	0.57 0.95	0.86 1.00	0.70 1.00	0.60 0.98	Neural (4 hid. units)	None RAIFLE	0.52 0.98	0.50 0.95	0.50 0.87	0.52 0.99	0.51 0.97	0.50 0.93
Neural (8 hid. units)	None RAIFLE	0.77 1.00	0.64 1.00	0.56 0.99	0.86 1.00	0.68 1.00	0.59 1.00	Neural (8 hid. units)	None RAIFLE	0.51 0.99	0.50 0.96	0.50 0.91	0.51 1.00	0.50 0.99	0.50 0.96
Neural (16 hid. units)	None RAIFLE	0.75 1.00	0.60 1.00	0.53 1.00	0.82 1.00	0.61 1.00	0.55 1.00	Neural (16 hid. units)	None RAIFLE	0.51 0.97	0.50 0.94	0.50 0.91	0.51 0.99	0.50 0.98	0.51 0.96



All Tech Considered

SOCIAL WEB

Facebook Manipulates Our Moods For Science And Commerce: A Roundup

JUNE 30, 2014 · 12:31 PM ET

By Elise Hu



Facebook researchers manipulated newsfeeds of nearly 700,000 users to study "emotional contagion."



• DONATE

All Tech Considered

SOCIAL WEB

Facebook Manipulates Our Moods For Science And Commerce: A Roundup

JUNE 30, 2014 · 12:31 PM ET

By Elise Hu



Facebook researchers manipulated newsfeeds of nearly 700,000 users to study "emotional contagion."





• DONATE

All Tech Considered

SOCIAL WEB

Facebook Manipulates Our Moods For Science And Commerce: A Roundup

JUNE 30, 2014 · 12:31 PM ET

By Elise Hu



Facebook researchers manipulated newsfeeds of nearly 700,000 users to study "emotional contagion."





• What if the server **cannot directly control** the training features?

- What if the server **cannot directly control** the training features?
- Example scenario: Ranking with images:
 - Server sends images to users
 - Users use a pre-trained model to extract features

- What if the server **cannot directly control** the training features?
- Example scenario: Ranking with images:
 - Server sends images to users
 - Users use a pre-trained model to extract features
- How to modify the images?
 - Can we make the extracted features resemble noise?

Initial image











Original

ResNet18

RegNet Y 800MF

DenseNet121

MNasNet 1.3

Evaluation

- Dataset: ImageNet 2012 (Validation set)
- Feature extractor: ResNet, RegNet, DenseNet, MNasNet
- Metrics: AUC. Baseline: Vanilla gradient inversion, FGSM

Evaluation

- Dataset: ImageNet 2012 (Validation set)
- Feature extractor: ResNet, RegNet, DenseNet, MNasNet
- Metrics: AUC. Baseline: Vanilla gradient inversion, FGSM

Vision Model	ADM	Linear			Neural (2 hidden units)			Neural (4 hidden units)			Neural (8 hidden units)		
		1x	2x	4x									
ResNet18 (512 features)	None FGSM RAIFLE	1.000 1.000 1.000	0.916 0.915 0.943	0.767 0.766 0.772	0.921 0.914 0.946	0.868 0.857 0.920	0.771 0.759 0.823	0.985 0.981 0.993	0.945 0.934 0.981	0.857 0.841 0.922	0.999 0.998 1.000	0.985 0.977 0.998	0.924 0.906 0.978
RegNet Y 800MF (784 features)	None FGSM RAIFLE	0.999 0.999 1.000	0.918 0.913 0.952	0.772 0.771 0.767	0.948 0.948 0.956	0.913 0.908 0.932	0.801 0.796 0.833	0.991 0.990 0.993	0.977 0.974 0.989	0.891 0.884 0.938	1.000 1.000 1.000	0.994 0.993 0.999	0.925 0.922 0.984
DenseNet121 (1024 features)	None FGSM RAIFLE	0.933 0.933 0.919	0.772 0.771 0.765	0.667 0.666 0.664	0.902 0.891 0.923	0.805 0.794 0.825	0.700 0.691 0.717	0.970 0.961 0.983	0.896 0.884 0.935	0.770 0.759 0.815	0.995 0.992 0.999	0.946 0.932 0.983	0.827 0.810 0.904
MNasNet 1.3 (1280 features)	None FGSM RAIFLE	0.994 0.989 1.000	0.895 0.885 0.939	0.764 0.758 0.775	0.936 0.930 0.940	0.858 0.845 0.882	0.754 0.742 0.791	0.985 0.982 0.991	0.926 0.913 0.965	0.787 0.771 0.875	0.996 0.994 0.999	0.942 0.928 0.992	0.775 0.755 0.924

Defense: Local Differential Privacy

• Apply Gaussian noise to local update before sending to server

Defense: Local Differential Privacy

• Apply Gaussian noise to local update before sending to server

Scenario	ADM	$\varepsilon\!=\!1$	$\varepsilon\!=\!20$	$\varepsilon\!=\!100$	$\varepsilon\!=\!500$	No DP
FNCF w/ ML-100K	N/A	0.50	0.52	0.56	0.74	1.00
FNCF w/ Steam-200K	N/A	0.50	0.56	0.72	0.90	0.96
FPDGD w/	None	0.50	0.54	0.57	0.58	0.66
MQ 2007	RAIFLE	0.50	0.56	0.66	0.75	0.82
FPDGD w/	None	0.50	0.50	0.50	0.50	0.50
MSLR10K	RAIFLE	0.50	0.52	0.58	0.62	0.80
FOLTR w/	None	0.50	0.52	0.55	0.62	1.00
ResNet18	RAIFLE	0.50	0.52	0.55	0.63	1.00
FOLTR w/	None	0.50	0.51	0.54	0.59	1.00
DenseNet121	RAIFLE	0.50	0.51	0.53	0.59	1.00

• Server only gets a single update without knowing user identities

- Server only gets a single update without knowing user identities
- Modified attack: Single out specific user(s) by fingerprinting items

- Server only gets a single update without knowing user identities
- Modified attack: Single out specific user(s) by fingerprinting items
- Example: FPDGD with MQ2007

User	$ d_1$	d_2	d_3	d_4	d_5	d_6
	:	:	•	:	:	÷
•	•	•	•	•	•	•
u_{k-1}	0.1	0.2	0.3	0.0	0.0	0.0
$u_{m k}$	0.0	0.0	0.0	0.4	0.5	0.6
u_{k+1}	0.1	0.2	0.3	0.0	0.0	0.0
•		•	•	•	•	
•		•	•	•	•	•

- Server only gets a single update without knowing user identities
- Modified attack: Single out specific user(s) by fingerprinting items
- Example: FPDGD with MQ2007

							Number of participants					
User	d_1	d_2	d_3	d_4	d_5	d_6	Model	ε	10	100	500	1000
÷		•	:	:				∞ 700	1.00 0.79	1.00 0.64	1.00 0.56	1.00 0.54
u_{k-1}	0.1	0.2	0.3	0.0	0.0	0.0	Linear	500	0.75	0.60	0.55	0.54
$u_{m k}$	0.0	0.0	0.0	0.4	0.5	0.6		300	0.71	0.57	0.54	0.53
u_{k+1}	0.1	0.2	0.3	0.0	0.0	0.0		100	0.61	0.54	0.52	0.51
•	.	•		•	•	•		∞	1.00	1.00	1.00	1.00
		:	:	:	:	:	Neurol	700	0.81	0.63	0.56	0.54
							incural (16 hidden unita)	500	0.77	0.60	0.54	0.53
						(10 modeli units)	300	0.73	0.58	0.52	0.51	
								100	0.61	0.54	0.51	0.51

NT....I.

Other defenses

- Check for data manipulation:
 - Cryptography: checksum
 - Heuristics: examine feature values

Other defenses

- Check for data manipulation:
 - Cryptography: checksum
 - Heuristics: examine feature values
- Minimize shared information:
 - Don't share entire model updates

Other defenses

- Check for data manipulation:
 - Cryptography: checksum
 - Heuristics: examine feature values
- Minimize shared information:
 - Don't share entire model updates
- Decentralized FL:
 - Peer-to-peer gossip

Conclusion

- RAIFLE: Proposes data manipulation novel attack vector in interaction-based FL
 - Strong performance compared to existing baselines

Conclusion

- RAIFLE: Proposes data manipulation novel attack vector in interaction-based FL
 - Strong performance compared to existing baselines
- Potential improvements:
 - Stealth
 - Other domains: Text
 - Combine with model manipulation

Conclusion

- RAIFLE: Proposes data manipulation novel attack vector in interaction-based FL
 - Strong performance compared to existing baselines
- Potential improvements:
 - Stealth
 - Other domains: Text
 - Combine with model manipulation
- Code available at: <u>https://github.com/dzungvpham/raifle</u>

