F-BLEAU: Fast Black-box Leakage Estimation

via Machine Learning

Giovanni Cherubin SPRING lab, EPFL University of Athens INRIA, École Polytechnique

IEEE Symposium on Security&Privacy 21 May, 2019







Assume no knowledge of system's internals



Assume no knowledge of system's internals



Assume no knowledge of system's internals

How can we measure:



Assume no knowledge of system's internals

How can we measure:

• how much *o* leaks about *s*?



Assume no knowledge of system's internals

How can we measure:

- how much *o* leaks about *s*?
- how hard is it to predict *s* given *o*?

Application Examples Location privacy



Application Examples Location privacy



Application Examples Side channels in crypto primitives' implementation



Application Examples Network traffic analysis (e.g., Website Fingerprinting)



$Pr\left({\textcircled{o}}(o) \neq s \right)$

$R^* = \min_{\mathfrak{G}} \{ \Pr\left(\mathfrak{F}(o) \neq s\right) \}$

Bayes risk *R*^{*} is the probability of error of the **optimal** adversary:

$$R^* = \min_{\mathfrak{G}} \{ Pr\left(\mathfrak{F}(o) \neq s\right) \}$$

Bayes risk *R*^{*} is the probability of error of the **optimal** adversary:

$$R^* = \min_{\mathfrak{G}} \{ Pr\left(\mathfrak{F}(o) \neq s\right) \}$$

On its basis we can compute **leakage measures** (e.g., Min-entropy, multiplicative/additive leakage)

Black-box estimates



We wish to measure leakage based on queries we make to the system:

Black-box estimates

$$s \longrightarrow \mathcal{B} \longrightarrow o$$

We wish to measure leakage based on queries we make to the system:

$$(s_1, o_1), (s_2, o_2), ..., (s_n, o_n)$$

$$(s, 0.4), (s, 0.7), (s, 1.2), (s, 0.4), (s, 0.4)$$

$$(s, 0.4), (s, 0.7), (s, 1.2), (s, 0.4), (s, 0.4)$$



$$(\underline{s, 0.4}), (\underline{s, 0.7}), (\underline{s, 1.2}), (\underline{s, 0.4}), (\underline{s, 0.4})$$



$$(\underline{s}, 0.4), (\underline{s}, 0.7), (\underline{s}, 1.2), (\underline{s}, 0.4), (\underline{s}, 0.4)$$

$$\mathbf{\mathfrak{G}}^{\mathsf{Freq}}(0.4) = \mathsf{Most} \text{ frequent among } \{s, s, s\}$$

$$(\underline{s}, 0.4), (\underline{s}, 0.7), (\underline{s}, 1.2), (\underline{s}, 0.4), (\underline{s}, 0.4)$$

$$\bigotimes^{\mathsf{Freq}} (0.4) = \mathsf{Most} \text{ frequent among } \{s, s, s\} = s$$

$$(\underline{s, 0.4}), (\underline{s, 0.7}), (\underline{s, 1.2}), (\underline{s, 0.4}), (\underline{s, 0.4})$$

$$\mathbf{\mathfrak{G}}^{\mathsf{Freq}}(0.4) = \mathsf{Most} \text{ frequent among } \{s, s, s\} = s$$
$$\mathbf{\mathfrak{G}}^{\mathsf{Freq}}(0.5) =$$

Suppose adversary observed data:

$$(\underline{s, 0.4}), (\underline{s, 0.7}), (\underline{s, 1.2}), (\underline{s, 0.4}), (\underline{s, 0.4})$$

$$\mathfrak{F}^{\mathsf{Freq}}(0.4) = \mathsf{Most}$$
 frequent among $\{s, s, s\} = s$
 $\mathfrak{F}^{\mathsf{Freq}}(0.5) = \mathsf{Random}$ guessing (e.g., most frequent label overal

I)

Suppose adversary observed data:

$$(\underline{s, 0.4}), (\underline{s, 0.7}), (\underline{s, 1.2}), (\underline{s, 0.4}), (\underline{s, 0.4})$$

$$\mathbf{\mathfrak{G}}^{\mathsf{Freq}}(0.4) = \mathsf{Most} \text{ frequent among } \{s, s, s\} = s$$

$$\mathbf{\mathfrak{G}}^{\mathsf{Freq}}(0.5) = \mathsf{Random guessing} (\mathsf{e.g., most frequent label overall})$$

• Does not work for **continuous** output space

$$(\underline{s, 0.4}), (\underline{s, 0.7}), (\underline{s, 1.2}), (\underline{s, 0.4}), (\underline{s, 0.4})$$

$$\mathfrak{F}^{\mathsf{Freq}}(0.4) = \mathsf{Most}$$
 frequent among $\{s, s, s\} = s$
 $\mathfrak{F}^{\mathsf{Freq}}(0.5) = \mathsf{Random}$ guessing (e.g., most frequent label overal

- Does not work for **continuous** output space
- Does not scale to large systems (needs at least one example per output value)



With minor modifications NN is leakage estimator



With minor modifications NN is leakage estimator (when o from finite set).

k-Nearest Neighbour [S'77]

•

k-Nearest Neighbour [S'77]



k-Nearest Neighbour [S'77]



If we choose $k/n \to 0$ and $k \to \infty$ as $n \to \infty$, then k-NN is leakage estimator.

We can use Machine Learning techniques for measuring the leakage.

We can use Machine Learning techniques for measuring the leakage.

Definition An ML rule is universally consistent if its error R_n converges to R^* as the size of the training data $n \to \infty$



We can use Machine Learning techniques for measuring the leakage.

Definition An ML rule is universally consistent if its error R_n converges to R^* as the size of the training data $n \to \infty$



E.g., NN (finite case), k-NN, SVM and Neural Networks (some param. choices)

We can use Machine Learning techniques for measuring the leakage.

Definition An ML rule is universally consistent if its error R_n converges to R^* as the size of the training data $n \to \infty$



E.g., NN (finite case), k-NN, SVM and Neural Networks (some param. choices)

Defense mechanisms:

- Geometric
- Laplacian
- Blahut-Arimoto
 [O+'17]



Defense mechanisms:

- Geometric
- Laplacian
- Blahut-Arimoto
 [O+'17]



Defense mechanisms:

- Geometric
- Laplacian
- Blahut-Arimoto [O+'17]



Defense mechanisms:

- Geometric
- Laplacian
- Blahut-Arimoto
 [O+'17]



Defense mechanisms:

- Geometric
- Laplacian
- Blahut-Arimoto [O+'17]



Defense mechanisms:

- Geometric
- Laplacian
- Blahut-Arimoto [O+'17]





NN-based estimators Synthetic experiments

• Scale to large systems & excel when there's a metric on the output

NN-based estimators Synthetic experiments

- Scale to large systems & excel when there's a metric on the output
- When no metric: equivalent to Frequentist

NN-based estimators Synthetic experiments

- Scale to large systems & excel when there's a metric on the output
- When no metric: equivalent to Frequentist
- However, may converge slowly for maliciously crafted systems

Is there an "optimal" estimator (i.e., which converges faster than all the others)?

Is there an "optimal" estimator (i.e., which converges faster than all the others)?

No Free Lunch Theorem No. Uniformly averaged among all the possible systems, all the estimators are equivalent in terms of convergence.

Is there an "optimal" estimator (i.e., which converges faster than all the others)?

No Free Lunch Theorem No. Uniformly averaged among all the possible systems, all the estimators are equivalent in terms of convergence.

Takeaway Always try several estimators and select the one converging faster.

Is there an "optimal" estimator (i.e., which converges faster than all the others)?

No Free Lunch Theorem No. Uniformly averaged among all the possible systems, all the estimators are equivalent in terms of convergence.

Takeaway Always try several estimators and select the one converging faster.



fbleau based on this idea

TL;DL + future ideas

- Several applications: location privacy, side channels, traffic analysis, ...
- We can use ML (UC) rules: they **scale** to large problems, and tend to converge faster (or equivalently to) frequentist approach
- + More applications (e.g., side channels, attacks to ML models)



fbleau: https://github.com/gchers/fbleau

Takeaway Black-box security and ML are solving similar problems: let's bridge them.

References

Applications

[CS'10] T. Chothia, V. Smirnov, A traceability attack against e-passports, 2010.

[BK'08] *M. Backes, B. Köpf*, Formally bounding the side-channel leakage in unknown-message attacks, 2008.

[O+'17] S. Oya, C. Troncoso, F. Pérez-González, Back to the drawing board: Revisiting the design of optimal location privacy-preserving mechanisms, 2017.

Black-box security

[C+'10] K. Chatzikokolakis, T. Chothia, A. Guha, Statistical measurement of information leakage, 2010.

[C'17] G. Cherubin, Bayes, not naïve: Security bounds on website fingerprinting defenses, 2017.

ML

[CH'67] T. Cover, P. Hart, Nearest neighbor pattern classification, 1967.

[S'77] J. Stone, Consistent nonparametric regression, 1977.

[W'96] D. H. Wolpert, The lack of a priori distinctions between learning algorithms, 1996.

F-BLEAU: Fast Black-box Leakage Estimation

via Machine Learning

Giovanni Cherubin SPRING lab, EPFL University of Athens INRIA, École Polytechnique

IEEE Symposium on Security&Privacy 21 May, 2019





Spiky channel



2 secrets, 10K observable values

Spiky channel



2 secrets, 10K observable values







= Frequentist= NN $= k_n-NN (log_{10})$ $= k_n-NN (log)$ $= R^*$





