

Machine Learning as an Attack Surface

© Patrick Hall*

H₂O.ai

July 14, 2019

* This material is shared under a [CC By 4.0 license](#) which allows for editing and redistribution, even for commercial purposes. However, any derivative work should attribute the author and H2O.ai.

Contents

Poisoning

Watermarks

Inversion

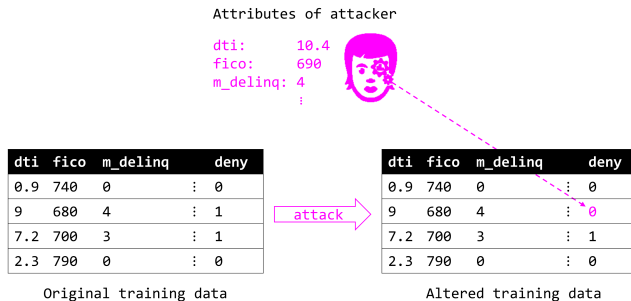
Membership

Adversaries

Impersonation

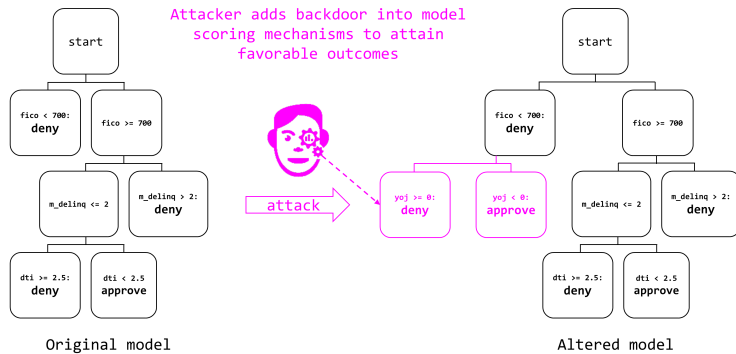
Blueprint

Data Poisoning Attacks

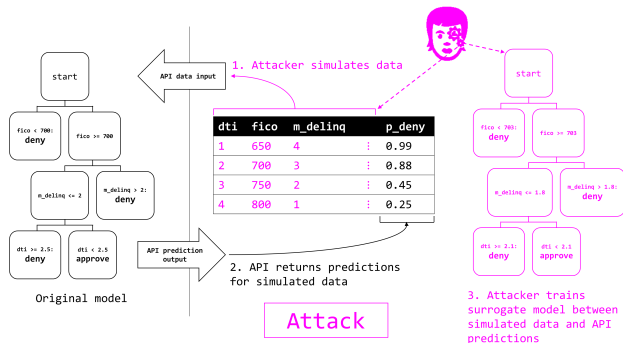


Attacker alters model training data to ensure favorable outcomes

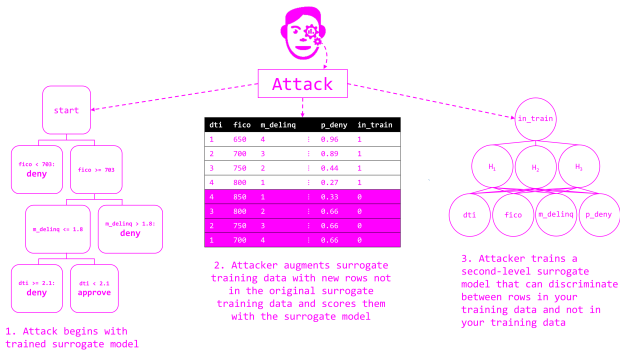
Watermark (i.e. Backdoor) Attacks



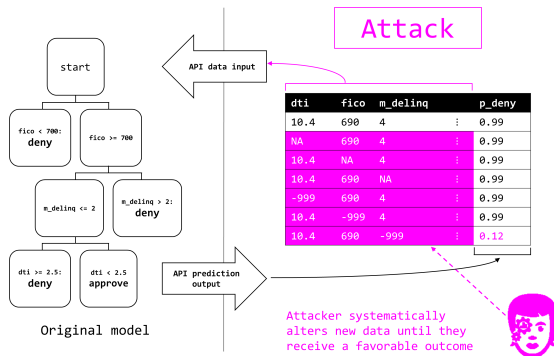
Surrogate Model Inversion Attacks



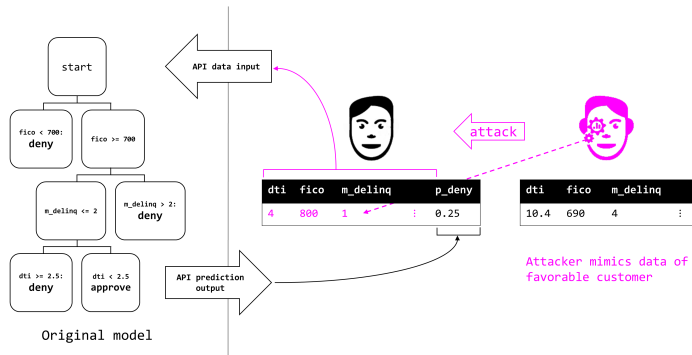
Membership Inference Attacks



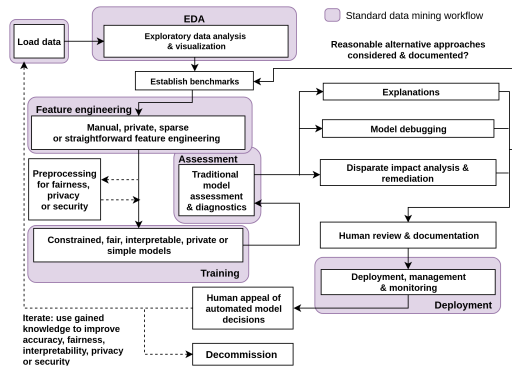
Adversarial Example Attacks



Impersonation Attacks



A Blueprint for Low-Risk Machine Learning



References

Proposals for model vulnerability and security:

[https:](https://www.oreilly.com/ideas/proposals-for-model-vulnerability-and-security)

[//www.oreilly.com/ideas/proposals-for-model-vulnerability-and-security](https://www.oreilly.com/ideas/proposals-for-model-vulnerability-and-security)

Can Your Machine Learning Model Be Hacked?!

<https://www.h2o.ai/blog/can-your-machine-learning-model-be-hacked/>