

Model Inversion Attacks.

- Develop an attack class that exploits confidence values revealed along with predictions.

To be explored in:

- ↳ decision trees in lifestyle survey
- Neural network for facial recognition.

Background

- an ML model is a deterministic function
 $f: \mathbb{R}^d \rightarrow Y$ [from d features to a set of responses Y]

→ If Y is a finite set, [names of people in facial recognition]

f is a classifier

→ If $Y = \mathbb{R}$,

the f is a regression.

→ Many classifiers incorporate regression and select the class estimating the maximum likelihood. → These estimate outputs are

confidence.
Consider formally, f is composition of two functions.

$$f \rightarrow \hat{f} + t$$

$\rightarrow \hat{f}: \mathbb{R}^d \rightarrow [0, 1]^m$, where m is a parameter specifying the number of confidences.
 $[m = |Y| - 1, \text{ one less the number of class labels}]$
 in eg.]

\rightarrow Second function is t .

$$t: [0, 1]^m \rightarrow Y$$

~~when $m=1$~~ $f(x) = t(\hat{f}(x))$

If classification call occurs, API returns both $f(x)$ and $\hat{f}(x)$.

ML APIs: Systems incorporate model f via application programming interfaces (API).

\rightarrow Typically in MLaaS such APIs are used.

Threat Model:

\rightarrow adversary is assumed to have whatever information the API expresses.

White box setting - access to download a model f .

Black box setting - attacker has the ability to make prediction queries on feature vectors.

→ queries can be adaptive
⇒ queried features can be a function of previously retrieved predictions.

→ adversary has no access to training data
[In some cases, data is public. However, here critical data is considered.]

Fredrikson et al.'s attack:

- Consider a linear regression model f .
- predict: a real valued suggested initial dose of a drug.
- feature vector: patient's demographic information, medical history,

Sensitive attribute is genetic marker
↳ set as first feature x_1 .

Attacker access: white box access to f .

Given auxiliary information side $(x, y) = (x_2, \dots, x_k, y)$
for an instance (x, y) , attempts to infer genetic marker x_1 .

Example steps of inversion attack:


- Assume: aux gives empirically computed standard deviation σ for a gaussian error model err and marginal priors.

$p = (p_1, \dots, p_t)$

Marginal priors (p_i) \rightarrow

How to compute marginal prices?



- Partition real line in disjoint buckets

 ranges of values.

- each bucket = v
- each $p_i(v) \equiv$ No of times x_i falls in v
 number of training vectors (1 db)

→ Algorithm completes the target feature vector with possible values of x_1 and computes weighted probability estimate to get correct value.

→ Gaussian error model will penalize values of x , that force prediction to be far from label y .

Algorithm : Adversary $A^t(\text{err}, p_1, x_2, \dots, x_t, y)$

Input : Assume : err : Gaussian error model with standard deviation σ .

- marginal priors $p = (p_1, \dots, p_t)$
- ~~aux~~ side info: (x_2, \dots, x_t)

adversary $A^t(\text{err}, p_1, x_2, \dots, x_t, y)$

1. for each value of v of x_1 do
2. $x' = (v, x_2, \dots, x_t)$
3. $r_v = \text{err}(y, f(x')) \cdot \prod_i p_i(x_i)$
4. Return $\arg\max_v r_v$.

- Algo returns maximum a posteriori (MAP) estimate for x_1 given available information.
- minimizes adversary's misprediction rate.
 $\arg\max_v r_v \rightarrow$ gets argument of the maximum
 \downarrow
set of points for which r_v attains maximum.

→ This algorithm produces least biased
maximum a posteriori (MAP) estimate for π_1

↓
• MAP → • Maximum a posteriori probability (MAP)
is an estimate of an unknown quantity.
• equals the mode of posterior distribution.

Posterior probability - The updated probability
(prior probability) of an event taking
into consideration some new event.

Probability Density function - (density of a continuous random
variable)

→ whose value at any given sample in the
sample space can be interpreted as the
relative likelihood of random variable equal
to the sample.

MAP gives the value that maximises the
posterior probability.

MAP Investors for TREES.

- Decision tree recursively partitions the feature space into disjoint regions R_1, \dots, R_m .
- for an instance $(x, y) \rightarrow$ find region of x
 \rightarrow return most likely value of y .

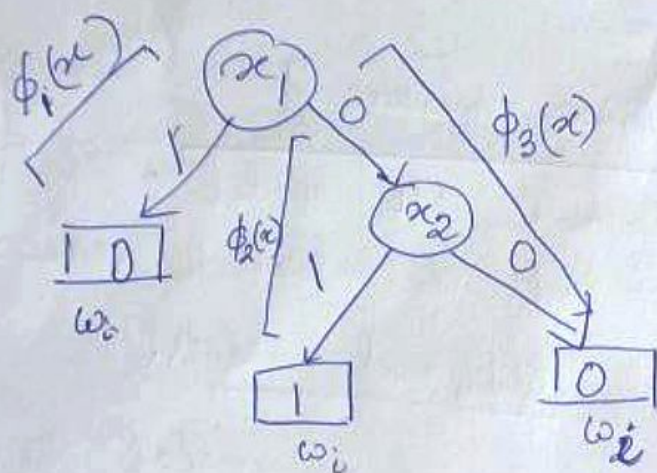
Eg: let $y = \neg x_1 \wedge x_2$.

\therefore If $x_1 = 1 \rightarrow y = 0$.

else if $x_1 = 0$,

\rightarrow if $x_2 = 1$ $y = 1$

else $x_2 = 0$, $y = 0$.



$$\phi_1(x) = x_1$$

$$\phi_2(x) = (1 - x_1) \cdot x_2$$

$$\phi_3(x) = (1 - x_1)(1 - x_2)$$

Trees can be mathematically characterized as follows;

$$f(x) = \sum_{i=1}^m \omega_i \phi_i(x), \text{ where } \phi_i(x) \in \{0, 1\}$$

$[\phi_i \rightarrow \text{indicator of region } R_i \text{ and } \omega_i]$

How to return confidence measures?

- Classification measures by getting w_i
- w_i is set to a vector corresponding to distribution of class labels as observed in the training set. [In complex scenario, w_i is set from training data]
- Say you have added 89 samples for $x_1 = 1$ and $x_2 = 1, (11)$ or change the 11 → 10, get confidence measures for classification.

What is confidence measure?

- Confidence represents ML model's estimated probability that the extracted value is correct given documents/inputs and labels.

Confidence = No. of a correct prediction on a dataset.

Classification results:

$$f(x) = \arg \max_j \left(\sum_{i=1}^m w_i [j] \phi_i(x) \right)$$

Confidences: $\hat{f}(x) = \left[\frac{w_{i^*} [1]}{\sum_i w_i [i]} \quad \dots \quad \frac{w_{i^*} [1Y1]}{\sum_i w_{im} [i]} \right]$

$[i^*]$ takes ~~last~~ values where $\phi_{i^*}(x) = 1$

Decision tree APIs

- General web APIs expose training and querying routines for decision trees to end users
[eg. BigML, Microsoft machine learning etc]
- Users can upload their datasets to these services to train a decision tree

Big ML : publish trees in different modes:

• Black box

• White box

⇒ For both setting users has marginal priors for each feature of training set.

↳ user has confusion matrix C

C_{ij} = gives number of training instances

for which $y = i$ and predicted model label = j .

→ White box setting:

Attacker has to: n_i = count of number of training instances that match ϕ_i path.

⇒ computes the confidence.

The inversion problem:

- Fix a tree $f(x) = \sum_{i=1}^m w_i \phi_i(x)$
- Let (x, y) be target instance that will be from target data or not.
- For simplicity, there is one sensitive feature
→ make target feature set $T = \{1\}$;
→ [extending our attacks to one more one feature is straightforward]

side info: $\text{side}_\ell(x_\ell, \dots, x_d)$
 $[d \geq 2]$

$K = \{1, \dots, d\}$ be set of known feature indices]

$[x_K$ represent $(d-1)$ dimensional vector]

Black-box ML:

- decision tree produces discrete outputs
- Use confusion matrix C
- define $\text{err}(y, y') \propto P_\ell [f(x) = y' \mid y \text{ is true label}]$

White-box ML :

- 1) Attacker knows each ϕ_i and n_i correspond to ϕ_i .

2) from this $N = \sum_{i=1}^m n_i$ [Total number of samples in training set]

- 3) Known value x_k induce a set of paths

$$S = \{S_i\}_{1 \leq i \leq m};$$

$$S = \{(\phi_i, n_i) \mid \exists x' \in \mathbb{R}^d, \quad x'_k = x_k \wedge \phi(x')\}$$

- Each path corresponds to a basis function ϕ_i .

• Let $p_i = \frac{n_i}{N}$

[p_i gives information about the joint distribution on features used to build training set.]

$P_R[S_i] =$ probability of drawing a row from joint prior that traverses S_i .

$p_i =$ empirical estimate for this quantity, derived from draws to produce training set

Basis functions partition feature space

$\Rightarrow x$ traverses exactly one path in S

Indicator function $\exists x' \in \mathcal{R}^d, x'_i = v \wedge \phi_i(x')$

Estimator (E) characterizes the probability that $x_1 = v$, given x traverses one of paths s_1, \dots, s_m , and $x_k = v_k$:
[known values]

$$P_H [x_1 = v \mid (s_1, v, \dots, v, s_m) \wedge x_k = v_k]$$
$$\propto \frac{\sum_{i=1}^m p_i \phi_i(v) \cdot P_H [x_k = v_k] \cdot P_H [x_1 = v]}{\sum_{j=1}^m p_j \phi_j(v)}$$

$$\propto \frac{1}{\sum_{j=1}^m p_j \phi_j(v)} \sum_{1 \leq i \leq m} p_i \phi_i(v) \cdot P_H [x_1 = v] \quad (1)$$

= white-box with counts (WBWC) estimator

\rightarrow adversary outputs a value v that maximizes as a guess for x_1 .