

23.09.2003

Alexey Vyskufov

Cryptographic techniques for Privacy-Preserving  
Data Mining

Alice and Bob owns confidential databases.

They want to run a data mining algorithm against the union of

databases.

The want to do it in "private" manner.

An ideal scenario: the data mining algorithm is run by trusted third party.

The aim: build a protocol with the same level of "privacy" as in ideal scenario.

The setting: Alice and Bob are semi-honest — honest but curious.

The setting: Authenticity of input data is somehow checked.

- The idea: receiver uses two public keys; for one the private key is known and for one it the private key is not known.
  - The receiver learns  $x_o$
  - Receiver:  $o \in \{0, 1\}$
  - Sender:  $(x_0, x_1)$
- 1-out-of-2 oblivious transfer.

Oblivious polynomial evaluation involves a sender and a receiver.  
The sender's input is a polynomial  $Q$  of degree  $k$  over some field  $\mathcal{F}$ .  
and the receiver's input is an element  $z \in \mathcal{F}$ .  
The receiver learns  $Q(z)$ .  
(We will call the sender as Alice and the receiver as Bob.)

The given function should be represented as circuit with gates defined over, e.g., all functions  $g : \{0, 1\} \rightarrow \{0, 1\}$ . Any polynomial-time function can be expressed as a circuit of polynomials size. Alice assigns two random values for every wire in circuit — “garbled” values. They should be long enough to use as a key to pseudo-random function.

For every gate Alice prepares a table which allows to find garbled output of gate if garbled inputs are known, revealing no other information.

- There is an entry in table for every pair of input values.
- Entries are marked by labels, not by actual (garbled) values.
- The (garbled) value of output wire is encrypted using the (garbled) values of input wires.

Alice sends wiring of the original circuit, gate tables and tables, which allow to convert garbled values of output wires to normal transfer. Bob gets garbled values for his inputs using 1-out-of-2 oblivious values to Bob. Bob computes the circuit.

Database is a set of transactions and each column is an attribute taking on different values.

One of the attributes is designated as a class attribute; the set of possible values for this attribute being the classes.

We wish to predict the class of a transaction using only non-class attributes.

The ID<sub>3</sub> algorithm solves classification problem by building decision tree — rooted tree containing nodes and edges. Each internal node is a test node and corresponds to an attribute. The edges leaving the node correspond to the possible values of the attribute. The leaves contain the expected class value.

$$\text{Gain}(A) = H^C(T) - H^C(T|A)$$

$$((\iota v)L)^CH \frac{|L|}{|(\iota v)L|} \sum_u^{j=1} = (A|L)^CH$$

$$\frac{|L|}{|(\iota^i)L|} \otimes \log \frac{|L|}{|(\iota^i)L|} - \sum_l^{i=1} = (L)^CH$$

How to choose the “best” attribute?

attribute.

At each step the database is partitioned according to the “best”

ID3 algorithm is recursive.

We can implement ID3 only with some precision  $\delta$ .

$$|Gain(A_1) - Gain(A_2)| < \delta$$

This is not a big problem.

Remark:  $1/|T|$  is not important. We also can use  $\ln$  instead of  $\log$ .

So we need to calculate  $(u_1 + u_2) \log(u_1 + u_2)$ .

$$\begin{aligned}
& \cdot |(\cdot^i a)_L| \log |(\cdot^i a)_L| \sum_{l=1}^j \frac{|L|}{l} \\
& + |(\cdot^i a)_L| \log |(\cdot^i a_j, c_j)_L| \sum_l \sum_{m=1}^j \frac{|L|}{l} - = \\
& \frac{|(\cdot^i a)_L|}{|(\cdot^i a_j, c_j)_L|} \log \frac{|(\cdot^i a)_L|}{|(\cdot^i a_j, c_j)_L|} - \sum_l |(\cdot^i a)_L| \sum_{m=1}^j \frac{|L|}{l} = \\
& ((\cdot^i a)_L)^C H \frac{|L|}{|(\cdot^i a)_L|} \sum_{m=1}^j = (A \mid L)^C H
\end{aligned}$$

It is enough to minimise  $H^C(T \mid A)$ .

The only difficult step — finding the attribute with maximal gain.

Alice owns  $u_1$ , Bob owns  $u_2$ . The goal: Alice obtains  $u_1$  and Bob obtains  $u_2$  s.t.

1.  $u_1 + u_2 = (u_1 + u_2) \ln(u_1 + u_2) \bmod |\mathcal{F}|$
2.  $u_1$  and  $u_2$  are uniformly distributed in  $\mathcal{F}$  when viewed independently of one another.

Why not Yao's algorithm? It's ineffective, the size of the circuit is of the order of multiplication of the size of its inputs.

Let's start with computing random shares  $u_1$  and  $u_2$ , s.t.

$$u_1 + u_2 = \ln(u_1 + u_2).$$

Error shrinks exponentially.

$$1 > \varepsilon > \frac{i}{(-1)^{i-1} e^i}, -1 > \varepsilon > 1$$
$$\sum_{k=1}^{i-1} \approx \ln(1 + \varepsilon)$$

Let  $2^n$  be the power of 2 which is closest to  $x$ . Then  $x = 2^n(1 + \varepsilon)$

and

$$\ln x = n \ln 2 + \varepsilon - \frac{\varepsilon^2}{2} + \frac{\varepsilon^3}{3} - \frac{\varepsilon^4}{4} + \dots$$

Let  $N$  be large ( $N > \log|T|(< n)$ ).

Alice constructs small circuit which inputs  $u_1, u_2$  and random  $a_1$  and  $b_1$  and outputs  $a_2 = e^{2_N} - a_1$  and  $b_2 = 2_N n \ln 2 - b_1$ .

Alice constructs the polynomial ( $w_1$  — random).

$$O(x) = \text{lcm}(2, \dots, k) \sum_k^{\text{lcm}(2, \dots, k)} (-1)^{i-1} \frac{x^i}{(x + a_1)^{N(i-1)}} - w_1,$$

Bob gets  $w_2 = O(\alpha^2)$  using oblivious polynomial evaluation.

$w_1 + w_2$  is an approximation of  $\ln e$  up to multiplicative factor  $(\text{lcm}(2, \dots, k))^2$  (which is public and can be removed). Let

$$u_i = u_i + \text{lcm}(2, \dots, k) B_i.$$

$$\text{lcm}(2, \dots, k) 2_N \ln x.$$

$$w_1 + w_2 \approx \text{lcm}(2, \dots, k) 2_N \ln e + \text{lcm}(2, \dots, k) 2_N n \ln 2 =$$

Then

So

Alice

defines  $(r_1, r_2 \leftarrow \text{random})$ :

$$u_1 + u_2 x, u_1 + u_2 \approx \ln x.$$

Alice sets her share to be  $u_1 u_1 + r_1 + r_2$ .

$$(u_1 + u_2)(u_1 + u_2) - u_1 u_1 - r_1 - r_2.$$

$$u_1 u_2 + r_1 + u_1 u_2 + r_2 + u_2 u_2 =$$

$$P_1(u_2) + P_2(u_2) + u_2 u_2 =$$

Bob runs oblivious polynomial evaluation and obtains

$$\cdot, y = u_1 y + r_1, P_2(y) = y + r_2.$$

Now it is not a problem to find the "best" attribute.