

# Privacy Preserving Distributed ID3 Algorithm

Chen Yuechen  
Fung Divine  
Nan Meng

April 27, 2016

# Overview

- \* Introduction
- \* Problem Definition
- \* Solution
- \* Result
- \* Conclusion

# Privacy Preserving Data Mining

- Mining while protecting the privacy of data.

Figure : Lindell's definition

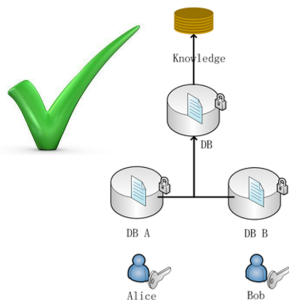
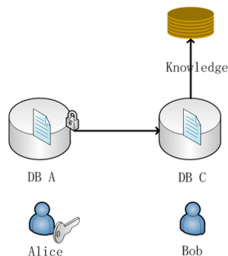


Figure : Agrawal's definition



# ID3 Algorithm

- ID3 is an algorithm used to generate a decision tree from a dataset, and is typically used in the data mining.
  1. Calculate the **entropy** of every attribute using the data set S.
  2. Split the set S into subsets using the attribute for which entropy is minimum
  3. Make a decision tree node containing that attribute
  4. Recurse on subsets using remaining attributes.

# Distributed ID3 Algorithm

Table : Play Golf Dataset

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes



**Alice**



**Bob**

# Distributed ID3 Algorithm

- Data is distributed in two or more parties

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

Table : Bob

Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

- Combine data together and get a decision tree

# Problem Definition

- However, data is **privacy**.

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

Table : Bob

Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

- How to share data in a safe way in distributed ID3 algorithm?

# An Example of Distributed ID3 Algorithm

- Here we use an example of Distributed ID3 algorithm to clearly define the problem. For example, Compute the entropy of **Rainy**.

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

3 records, 2 No, 1 Yes

Table : My caption

Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

2 records, 1 No, 1 Yes

$$\begin{aligned} Entropy(Rainy) &= - \underbrace{\frac{2+1}{3+2} \log_2 \left( \frac{2+1}{3+2} \right)}_{PlayGolf=No} - \underbrace{\frac{1+1}{3+2} \log_2 \left( \frac{1+1}{3+2} \right)}_{PlayGolf=Yes} \\ &= -\frac{3}{5} \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \log_2 \left( \frac{2}{5} \right) \end{aligned}$$



# An Example of Distributed ID3 Algorithm

- For example, Compute the entropy of **Rainy**.

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

3 records, 2 No, 1 Yes

Table : My caption

Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

2 records, 1 No, 1 Yes

$$\begin{aligned} Entropy(Rainy) &= -\frac{2+1}{3+2} \log_2\left(\frac{2+1}{3+2}\right) - \frac{1+1}{3+2} \log_2\left(\frac{1+1}{3+2}\right) \\ &= -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \end{aligned}$$

# An Example of Distributed ID3 Algorithm

- For example, Compute the entropy of **Rainy**.

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

3 records, 2 No, 1 Yes

Table : My caption

Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

2 records, 1 No, 1 Yes

$$-\frac{2+1}{3+2} \log_2 \left( \frac{2+1}{3+2} \right)$$

# An Example of Distributed ID3 Algorithm

- For example, Compute the entropy of **Rainy**.

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

3 records, 2 No, 1 Yes

Table : My caption

Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

2 records, 1 No, 1 Yes

$$-\frac{2+1}{3+2} \log_2 \left( \frac{2+1}{3+2} \right)$$

# An Example of Distributed ID3 Algorithm

- For example, Compute the entropy of **Rainy**.

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

3 records, 2 No, 1 Yes

Table : My caption

Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

2 records, 1 No, 1 Yes

$$\frac{2+1}{3+2}$$

# An Example of Distributed ID3 Algorithm

- For example, Compute the entropy of **Rainy**.

Table : Alice

Outlook	Temp	Humidity	Windy	Play Golf
Rainy	Hot	High	FALSE	No
Rainy	Hot	High	TRUE	No
Overcast	Hot	High	FALSE	Yes
Sunny	Mild	High	FALSE	Yes
Rainy	Mild	Normal	TRUE	Yes

**a** records, **x** No, 1 Yes

Table : My caption

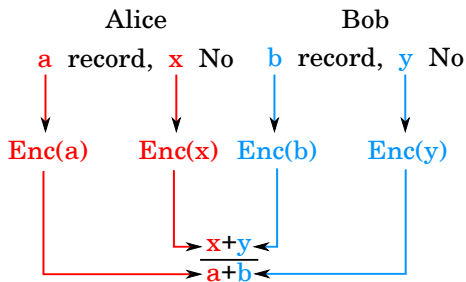
Outlook	Temp	Humidity	Windy	Play Golf
Overcast	Cool	Normal	TRUE	Yes
Rainy	Mild	High	FALSE	No
Rainy	Cool	Normal	FALSE	Yes
Sunny	Mild	Normal	FALSE	Yes

**b** records, **y** No, 1 Yes

$$\frac{x+y}{a+b}$$

# Problem Definition

- Compute  $\frac{x+y}{a+b}$  without reveal  $a, x, b, y$ .
- Realize Privacy Preserving Distributed ID3 algorithm.



$Enc(\cdot)$  – Encryption Algorithm

# Solution

## PPWAP

- **PPWAP**: Privacy Preserving Weight Average Protocol
- In this project, we choose PPWAP by **Pailier Encryption**.

# Pailier Encryption

- $KeyGeneration()$ : Generate public key  $PK$ , and secret key  $SK$ .
- $Encryption(m, PK)$ : Using  $PK$  to encrypt message  $m$ , output  $Enc(m)$ .
- $Decryption(Enc(m), SK)$ : Using  $SK$  to decrypt  $Enc(m)$ , output  $m$ .



# Pailier Encryption

- **Property:** Addition Homomorphism
- Given two messages  $m1$  and  $m2$ ,  
 $Enc(m1 + m2) = Enc(m1) \cdot Enc(m2)$ .
- The encryption of  $m1 + m2$  can be computed by  $Enc(m1)$  and  $Enc(m2)$ .

# PPWAP based on Paillier Encryption

## Privacy Preserving Weighted Average Protocol

- Within the help of Paillier, build PPWAP scheme.

Alice

Bob

$\text{Enc}(a)$      $\text{Enc}(x)$   $\text{Enc}(b)$      $\text{Enc}(y)$

$$\frac{x+y}{a+b}$$

# PPWAP

Alice

1.  $KeyGeneration() : SK, PK$   
 $Encryption(a, PK) : Enc(a)$   
 $Encryption(x, PK) : Enc(x)$

Bob



$Enc(a)$   
 $Enc(x)$

---

2.

Random integer  $z$   
 $Enc(a)^z, Enc(x)^z$

# PPWAP

Alice

1.  $KeyGeneration() : SK, PK$   
 $Encryption(a, PK) : Enc(a)$   
 $Encryption(x, PK) : Enc(x)$

Bob



$Enc(a)$   
 $Enc(x)$

---

2.

Random integer  $z$   
 $Enc(a)^z, Enc(x)^z$

$$Enc(a)^z = Enc(a) \dots Enc(a) = Enc(a + a + \dots + a) = Enc(za)$$

Alice

1.  $KeyGeneration() : SK, PK$   
 $Encryption(a, PK) : Enc(a)$   
 $Encryption(x, PK) : Enc(x)$



Bob

$Enc(a)$   
 $Enc(x)$

2.

Random integer  $z$   
 $Enc(a)^z, Enc(x)^z$   
 $Enc(za), Enc(zx)$

Alice

Bob

3.  $Enc(za + zb)$   
 $Enc(zx + zy)$



$$\begin{aligned} \text{Encryption}(b, PK) : Enc(b) &\Rightarrow Enc(zb) \\ \text{Encryption}(y, PK) : Enc(y) &\Rightarrow Enc(zy) \\ Enc(za + zb) &= Enc(za) \\ Enc(zb)Enc(zx + zy) &= Enc(zx)Enc(zy) \end{aligned}$$

4.  $Decryption(Enc(za + zb), SK) :$   
 $za + zb$   
 $Decryption(Enc(zx + zy), SK) :$   
 $zx + zy$



$$\frac{zx+zy}{za+zb} = \frac{x+y}{a+b} \Rightarrow \frac{x+y}{a+b}$$

# Algorithm

---

**Algorithm 1** Two-party jointly decision tree algorithm

---

```
1: procedure PRIVACYID3( $D$ ,  $Attribute$ ,  $transInfo$ ,  $T$ )
2:    $ct \leftarrow createNode()$ 
3:    $label(ct) = mostCommonClass(D, transInfo, T)$ 
4:   IF  $\forall \langle \mathbf{x}, c(\mathbf{x}) \rangle \in D : c(\mathbf{x}) = c$  THEN
5:      $return(t)$ 
6:   ENDIF
7:   IF  $Attributes = \emptyset$  THEN
8:      $return(t)$ 
9:   ENDIF
10:   $\tilde{A} = \underset{A \in Attributes}{argmax} (InformationGain(D, A, transInfo))$ 
11:  for  $a \in \tilde{A}$  do
12:     $D_a = \{ \langle \mathbf{x}, c(\mathbf{x}) \rangle \in D : \mathbf{x} \upharpoonright_{\tilde{A}} = a \}$ 
13:    IF  $D = \emptyset$  THEN
14:       $ct' = createNode()$ 
15:       $label(ct') = mostCommonClass(D, transInfo, T)$ 
16:       $createEdge(ct, a, ct')$ 
17:    ELSE
18:       $transInfo^* = TransInfo(B, Attributes \setminus \{\tilde{A}\}, T^*)$ 
19:       $createEdge(ct, a, PrivacyID3(D_a, Attributes \setminus \{\tilde{A}\}, transInfo^*, T))$ 
20:    ENDIF
21:  return  $combinedtree$ 
```

---

Figure : Two-party Jointly Decision Tree Algorithm.

# Result

- Demo
- Efficiency: The runtime depend on 3 factors.
  - \* Dataset size
  - \* Length of Key in encryption algorithm
  - \* Number of parties



# Algorithm Implement

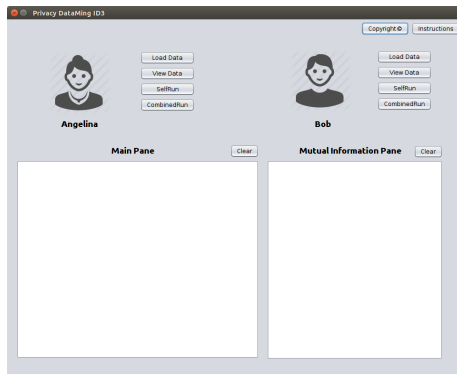


Figure : Welcome Graphical User Interface.

# Algorithm Implement

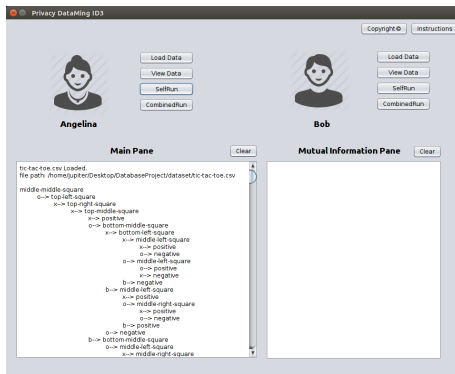


Figure : Result of single-party ID3 algorithm on tic-tac-toe2 dataset.

# Conclusion

- The PPWAP scheme is proposed in 2005 in PP K-means.
  - \* PPWAP can be extended to multi-party, supports Multi-party distributed ID3 algorithm.
- Further research focus on improving the security level.
  - \* The scheme became safer and more complex.
- Current research focus on preventing malicious attack.

# Conclusion

- Select two large primes,  $p$  and  $q$ .
- Calculate the product  $n = p \times q$ , such that  $\gcd(n, \Phi(n)) = 1$ , where  $\Phi(n)$  is  $(p-1)(q-1)$ .
- Choose a random number  $g$ , where  $g$  has order multiple of  $n$  or  $\gcd(L(g^\lambda \bmod n^2), n) = 1$ , where  $L(t) = (t-1)/n$  and  $\lambda(n) = \text{lcm}(p-1, q-1)$ .
- The public key is composed of  $(g, n)$ , while the private key is composed of  $(p, q, \lambda)$ .
- The Encryption of a message  $m < n$  is given by:
  - $c = g^m \cdot r^n \bmod n^2$
- The Decryption of ciphertext  $c$  is given by: The Decryption of ciphertext  $c$  is given by:
  - $m = (L(g^\lambda \bmod n^2) / L(g^\lambda \bmod n^2)) \bmod n$

# The End