Challenges in Privacy-Preserving Data Analysis

Kamalika Chaudhuri

University of California, San Diego

Sensitive Data

Medical Records

Genetic Data

Search Logs







AOL Violates Privacy

AOL Violates Privacy

A Face Is Exposed for AOL Searcher No. 4417749

By MICHAEL BARBARO and TOM ZELLER Jr. Published: August 9, 2006

Buried in a list of 20 million Web search queries collected by AOL and recently released on the Internet is user No. 4417749. The number was assigned by the company to protect the searcher's anonymity, but it was not much of a shield.



No. 4417749 conducted hundreds of searches over a three-month period on topics ranging from "numb fingers" to "60 single men" to "dog that urinates on

Netflix Violates Privacy



2-8 movie-ratings and dates for Alice reveals: Whether Alice is in the dataset or not Alice's other movie ratings

High-dimensional Data is Unique

Example: UCSD Employee Salary Table

Position	Gender	Department	Ethnicity	Salary
Faculty	Female	CSE	SE Asian	-

One employee (Kamalika) fits description!

Simply anonymizing data is unsafe!

Disease Association Studies

	C	a	n	C	91								e	a	It	h	y		
1.00 .190 1.00 .216 .251 .186 .117 .154 .011 .190 .400 .270 .215 .101 .085 .239 .071 .471 .117 .179 .202	1.00 .047 1.0 .170 .08 .102 .09 .294 .24 .170 .09 .163 .11 .243 .09 .132 .09	00 33 1.00 95 .139 18 .140 56 .234 11 .161 94 .144 94 .087	1.00 .141 .099 .093 .123 .159	1.00 .175 .199 .283 .207	1.00 .157 .216 .108	1.00 .274 .092	1.00 .294	1.00	1.00 .141 .099 .093 .123 .159 .088 .046 .046 .045 .178	1.00 .175 .199 .283 .207 .152 .161 .392 .155 .135	1.00 .157 .216 .008 .075 .092 .122 .135 .102	1.00 .274 .092 .163 .072 .229 .139 .258	1.00 .294 .156 .157 .160 .110 .314	1.00 .220 .143 .172 .048 .165	1.00 .147 .145 .126 .147	1.00 .177 .104 .158	1.00 .169 .131	1.00	1.00
Correlations									Correlations										

Correlation (R² values), Alice's DNA reveals: If Alice is in the Cancer set or Healthy set

Simply anonymizing data is unsafe! Statistics on small data sets is unsafe!



Correlated Data

User information in social networks



Physical Activity Monitoring



Why is Privacy Hard for Correlated data?

Correlation: neighbor's info leaks info on user

Why is Privacy Hard for Correlated data?

Correlation: neighbor's info leaks info on user

User information in social networks



Physical Activity Monitoring



How do we learn from sensitive data while still preserving privacy?

How do we compute statistics on correlated data while still preserving privacy?

- I. Privacy for Uncorrelated Data
 - How to define privacy
 - Privacy-preserving Classification
- 2. Privacy for Correlated Data
 - How to define privacy
 - Privacy-preserving Statistics

I. Privacy for Uncorrelated Data

- How to define privacy

Differential Privacy



Participation of a single person does not change output

Differential Privacy: Attacker's View



Differential Privacy [DMNS06]



For all D₁, D₂ that differ in one person's value, any set S, If $A = \alpha$ -private randomized algorithm, then:

 $\Pr(A(D_1) \in S) \le e^{\alpha} \Pr(A(D_2) \in S)$

Differential Privacy

I. Provably strong notion of privacy

2. Good approximations for many functions e.g, means, histograms, etc.

I. Privacy for Uncorrelated Data

- How to define privacy
- Privacy-preserving Classification

A Classification Problem: Flu Test

Could I have H1N1 flu (swine flu)?

Use the Flu Self-Assessment, based on material from Emory University, to:

- Learn whether you have the symptoms of H1N1 flu (swine flu)
- Help you decide what to do next

Take Flu Self-Assessment

Licensed from Emory University

You will have the opportunity to consent to share the information you provide

- Learn more about H1N1 flu
- What is H1N1 (Swine) Flu?
- Basics for Flu Prevention
- Guidelines for Taking Care of Yourself and Others
- People with Health Conditions

Predicts flu or not, based on patient symptoms Trained on sensitive patient data

From Attributes to Labeled Data



Classifying Sensitive Data



Private Data Public Classifier





Distribution P over labelled examples



Distribution P over labelled examples

Goal: Find a vector w that separates + from - for most points from P



Distribution P over labelled examples

Goal: Find a vector w that separates + from - for most points from P

Key: Find a simple model to fit the samples

Empirical Risk Minimization

Goal: Labeled data (x_i, y_i) , find w minimizing:

$$\frac{1}{2}\lambda \|w\|^2 +$$

$$\frac{1}{n}\sum_{i=1}^{n} L(y_i w^T x_i)$$

Regularizer (Model Complexity)

Risk (Training Error)

Some Examples



Risk: Hinge loss **Optimizer:** Support vector machines (SVM)

Risk: Logistic loss **Optimizer:** Logistic regression

ERM with **Privacy**

Given: labeled data (x_i, y_i) , **Find:** vector w that is:

(Private) Is private w.r.t training data

(Accurate) Approximately minimizes the regularized risk

Why is ERM not private for SVM?



SVM solution is a combination of support vectors If one support vector moves, solution changes

Why is ERM not private for SVM?



SVM solution is a combination of support vectors If one support vector moves, solution changes

How to make ERM private?



Pick w from distribution near the optimal solution

Privacy vs. Accuracy





Too little privacy

Too little accuracy

Pick distribution that gives **privacy** and **accuracy**

Privacy-preserving Classification

- I. ERM with privacy
- 2. Algorithm
Properties of Real Data



Optimization surface is very steep in some directions High loss if perturbed in those directions

Insight: Perturb optimization surface and then optimize

Empirical Risk Minimization

Goal: Labeled data (x_i, y_i) , find w minimizing:

$$\frac{1}{2}\lambda \|w\|^2 + \frac{1}{n}\sum_{i=1}^n L(y_i w^T x_i) + \frac{1}{n}b^\top w$$

Regularizer (Model Complexity) **Risk** (Training Error)

Perturbation (Privacy)

Algorithm: Perturbation



Perturbation b drawn from:

Magnitude:Drawn from $\Gamma(d, 1/\alpha)$ Direction:Uniformly at random

Privacy-preserving Classification

- I. ERM with privacy
- 2. Algorithm
- 3. Analytical guarantees

Privacy Guarantees

Algorithm: Given labeled data (x_i, y_i), find w to minimize: $\frac{1}{2}\lambda \|w\|^2 + \frac{1}{n}\sum_{i=1}^n L(y_iw^Tx_i) + \frac{1}{n}b^\top w$

Theorem: If L is convex and doubly-differentiable with $|L'(z)| \le 1$ and $|L''(z)| \le c$ then Algorithm is $\alpha + 2\log\left(1 + \frac{c}{n\lambda}\right)$ -differentially private

Privacy Guarantees

Algorithm: Given labeled data (x_i, y_i), find w to minimize: $\frac{1}{2}\lambda \|w\|^2 + \frac{1}{n}\sum_{i=1}^n L(y_iw^Tx_i) + \frac{1}{n}b^\top w$

L = Logistic Loss
Private Logistic Regression
L = Huber Loss
Private SVM

(Hinge Loss is not differentiable)

Measure of Accuracy

Number of samples for error ϵ (Fewer samples implies higher accuracy)

Sample Requirement

 $\begin{array}{ll} \mathsf{d} & : \texttt{\#dimensions} \\ \gamma & : \texttt{margin} \\ \alpha & : \texttt{privacy} \\ \epsilon & : \texttt{error} \\ \gamma, \alpha, \epsilon < 1 \end{array}$



Normal SVM:

Our Algorithm:

Standard Method:

 $\frac{1}{\epsilon^2 \gamma^2}$ $\frac{1}{\epsilon^2 \gamma^2} + \frac{d}{\epsilon \alpha \gamma}$

 $1/\epsilon^2 \gamma^2 + d/\epsilon^{3/2} \alpha \gamma$

Privacy-preserving Classification

- I. ERM with privacy
- 2. Algorithm
- 3. Analytical guarantees
- 4. Evaluation

Experiments

UCI Adult: Census/Income Data

Demographic dataset of size 47K

105 dimensions after preprocessing

Task: Predict if income above/below 50K

Results: SVM



Experiments

KDDCup99: Intrusion detection data

70K network connections

116 dimensions after preprocessing

Task: Predict if connection is malicious or not

Results: SVM



Privacy-preserving Classification

- I. ERM with privacy
- 2. Algorithm
- 3. Analytical guarantees
- 4. Evaluation

Talk Agenda:

- I. Privacy for Uncorrelated Data
 - How to define privacy
 - Privacy-preserving Classification
- 2. Privacy for Correlated Data

Why is Privacy Hard for Correlated data?

Correlation: neighbor's info leaks info on user

User information in social networks



Physical Activity Monitoring



Why is Differential Privacy not Enough for Correlated data?

Example: $D = (x_1, ..., x_n), x_i = I$ if i has flu



Goal: (1) Publish #people with flu in D (2) Prevent adversary from knowing who has flu

Example: $D = (x_1, ..., x_n), x_i = I$ if i has flu



I-DP: Output #people with flu + noise with stdev I

Example: $D = (x_1, ..., x_n), x_i = I$ if i has flu



I-DP: Output #people with flu + noise with stdev I

Not enough for privacy of people in connected components!

Example: D = $(x_1, ..., x_t)$, x_i = activity at time t



Goal: (1) Publish activity histogram (2) Prevent adversary from knowing activity at t

Example: D = $(x_1, ..., x_t)$, x_i = activity at time t



I-DP: Output histogram of activities + noise with stdev I

Example: D = $(x_1, ..., x_t)$, x_i = activity at time t



I-DP: Output histogram of activities + noise with stdev I

Not enough - activities across time are highly correlated!

Talk Agenda:

- I. Privacy for Uncorrelated Data
 - How to define privacy
 - Privacy-preserving Classification
- 2. Privacy for Correlated Data
 - How to define privacy

Pufferfish Privacy: Components

Secrets S: Information to be protected e.g: Alice's age is 25, Bob has a disease

Secret Pairs Q: Pairs of secrets to be indistinguishable e.g: (Alice's age is 25, Alice's age is 40), (Bob is in dataset, Bob is not in dataset)

Distribution Class Θ : Set of distributions that can plausibly generate the data

e.g: disease is passed on w.p. [0.1, 0.9]

Distribution Class models correlation!

Pufferfish Privacy

An algorithm A is α -Pufferfish private with parameters (S, Q, Θ) if for all $(\mathbf{s}_i, \mathbf{s}_j)$ in Q, for all $\theta \in \Theta$, $X \sim \theta$, all t, $\frac{p_{\theta,A}(A(X) = t|s_i, \theta)}{p_{\theta,A}(A(X) = t|s_i, \theta)} \leq e^{\alpha}$

whenever $P(s_i|\theta), P(s_j|\theta) > 0$



Pufferfish Privacy

An algorithm A is α -Pufferfish private with parameters (S, Q, Θ) if for all (s_i, s_j) in Q, for all $\theta \in \Theta$, $X \sim \theta$, all t, $e^{-\alpha} \leq \frac{p_{\theta,A}(s_i|A(X) = t, \theta)}{p_{\theta,A}(s_j|A(X) = t, \theta)} / \frac{p_{\theta}(s_i|\theta)}{p_{\theta}(s_j|\theta)} \leq e^{\alpha}$

whenever $P(s_i|\theta), P(s_j|\theta) > 0$



Knowing A does not affect adversary's belief on s_i vs. s_j

Pufferfish generalizes DP

Theorem: α -DP is equivalent to α -Pufferfish privacy with parameters (S, Q, Θ), where

Q:= { (i in data with value a, i in data with value b), $\forall i, a, b$ } U { (i in data with value a, j in data), $\forall i \neq j, a$ }

 Θ := Each individual is **independent**

There can be no utility if Θ allows any arbitrary correlation!

Talk Agenda:

- I. Privacy for Uncorrelated Data
 - How to define privacy
 - Privacy-preserving Classification
- 2. Privacy for Correlated Data
 - How to define privacy
 - Privacy-preserving Statistics

Correlation Measure: Bayesian Networks



Node: variable

Directed Acyclic Graph

Joint distribution of variables:

$$\Pr(X_1, X_2, \dots, X_n) = \prod_i \Pr(X_i | \operatorname{parents}(X_i))$$

Algorithm: Main Idea



Goal: Protect X₁

Algorithm: Main Idea



Goal: Protect X_I

Almost independent of X₁

Algorithm: Main Idea



Goal: Protect X_I

Almost independent of X_I

Add noise to hide + Small correction local terms + for rest

Defining "Almost Independence"

Max-influence of X_i on a set of nodes X_R :

$$e(X_R|X_i) = \max_{a,b} \sup_{\theta \in \Theta} \max_{x_R} \log \frac{\Pr(X_R = x_R | X_i = a, \theta)}{\Pr(X_R = x_R | X_i = b, \theta)}$$

Low $e(X_R|X_i)$ means X_R is almost independent of X_i

Fact: To protect X_i , correction term needed for X_R is $exp(e(X_R|X_i))$

How to find large "almost independent" sets

Brute force search is expensive

Use structural properties of the Bayesian network
Markov Blanket



Markov Blanket(X_i) = Set of nodes X_S s.t X_i independent of X\(X_i U X_S) given X_S

Define: Markov Quilt



 X_Q is a Markov Quilt for X_i if:

I. Deleting X_Q breaks graph into X_N and X_R

 $2.\,X_i \ lies \ in \ X_N$

3. X_R is independent of X_i given X_Q



Define:

 $= \frac{card(X_N)}{\alpha - e(X_Q|X_i)}$





The Markov Quilt Mechanism

For each X_i

Find the Markov Quilt X_Q for X_i with minimum score s_i

Output F(D) + (max_i s_i) Z where $Z \sim Lap(1)$

The Markov Quilt Mechanism

For each X_i

Find the Markov Quilt X_Q for X_i with minimum score s_i

Output F(D) + (max_i s_i) Z where $Z \sim Lap(1)$

Theorem: This preserves α -Pufferfish privacy

The Markov Quilt Mechanism

For each X_i

Find the Markov Quilt X_Q for X_i with minimum score s_i

Output F(D) + (max_i s_i) Z where $Z \sim Lap(1)$

Advantage: Poly-time in special cases.

Example: Activity Monitoring

Bayesian Network: Markov Chain



(Minimal) Markov Quilts: for X_i have form $\{X_{i-a}, X_{i+b}\}$



Example: Activity Monitoring

Bayesian Network: Markov Chain



(Minimal) Markov Quilts: for X_i have form $\{X_{i-a}, X_{i+b}\}$



Example: Activity Monitoring: Scores

- \mathcal{X} : set of states
- P_{θ} : transition matrix describing each $\theta \in \Theta$

Example: Activity Monitoring: Scores

- \mathcal{X} : set of states
- P_{θ} : transition matrix describing each $\theta \in \Theta$

Under some assumptions, relevant parameters are:

 $\pi_{\Theta} = \min_{\substack{x \in \mathcal{X}, \theta \in \Theta}} \pi_{\theta}(x) \quad \text{(min prob of x under stationary distr.)}$ $g_{\Theta} = \min_{\theta \in \Theta} \min\{1 - |\lambda| : P_{\theta}x = \lambda x, \lambda < 1\} \text{ (min eigengap of any } P_{\theta}\text{)}$

Example: Activity Monitoring: Scores

- \mathcal{X} : set of states
- P_{θ} : transition matrix describing each $\theta \in \Theta$

Under some assumptions, relevant parameters are:

 $\pi_{\Theta} = \min_{\substack{x \in \mathcal{X}, \theta \in \Theta}} \pi_{\theta}(x) \quad \text{(min prob of x under stationary distr.)}$ $g_{\Theta} = \min_{\theta \in \Theta} \min\{1 - |\lambda| : P_{\theta}x = \lambda x, \lambda < 1\} \text{ (min eigengap of any } P_{\theta}\text{)}$

Max-influence of $X_Q = \{X_{i-a}, X_{i+b}\}$ for X_i $e(X_Q|X_i) \le \log\left(\frac{\pi_{\Theta} + \exp(-g_{\Theta}b)}{\pi_{\Theta} - \exp(-g_{\Theta}b)}\right) + 2\log\left(\frac{\pi_{\Theta} + \exp(-g_{\Theta}a)}{\pi_{\Theta} - \exp(-g_{\Theta}a)}\right)$ Score(X_Q) = $\frac{a+b-1}{\alpha - e(X_Q|X_i)}$

Example: Activity Monitoring

Algorithm: For each X_i Find Markov Quilt X_Q = {X_{i-a},X_{i+b}} with minimum score s_i

Output F(D) + (max_i s_i) Z where $Z \sim Lap(1)$

Running Time: $O(T^3)$ (can be made $O(T^2)$)

Conclusion

Problem:

- privacy preserving classification of iid data
- private statistics for correlated data

Open Questions:

- better private algorithms for classification
- better models and mechanisms for
- correlated data

Acknowledgements



Claire Monteleoni



Anand Sarwate



Shuang Song



Mani Srivastava



Yizhen Wang

Questions?

