

Department of Computer Science http://www.cs.cornell.edu/johannes





- Customer web site trails
- Podcasts
- Blogs
- Email • Closed caption

(B) ---

• Print, film, optical, and magnetic storage: 5 Exabytes (EB) of new information in 2002, doubled in the last three years [How much Information 2003, UC Berkeley]

















Project Requirements

- Data
 - 14 TB every 2 weeks
 - Shipped on USB-2 disk drives
 - Need to archive raw data 5+ years
 Need to make data products available to the astronomy
 - research community
- Processing
 - Extremely processor intensive
 - Find new pulsars --- and other interesting phenomena

[Calimlim, Cordes, Demers, Gehrke, Lifka;

(B) --

۵ 🕲

http://arecibo.tc.cornell.edu]

Driving Factors: Analysis Capabilities

Data mining is the exploration and analysis of large quantities of data in order to discover valid, novel, potentially useful, and ultimately understandable patterns in data.

Example pattern (Census Bureau Data): If (relationship = husband), then (gender = male). 99.6%

Driving Factors: Connectivity and Bandwidth

- Metcalf's law (network usefulness increases squared with the number of users)
- Gilder's law (bandwidth doubles every 6 months)

SIGKDD 2006 Tutorial, August 2006

(B) carro

(B) c...

.

Concerns About Privacy

Recent example:

"Last week AOL did another stupid thing, but at least it was in the name of science...."

[Annalee Newitz, AlterNet, August 15, 2006]

SIGKDD 2006 Tutorial, August 2006

A Face Is Exposed for AOL Searcher No. 4417749 [New York Times, August 9, 2006]

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 everything."
 And search by search, click by click, the identity of AOL user No. 4417749 became easier to discern. There are queries for "landscapers in Lilburn, Ga," several people with the last name Arnold and "homes sold in shadow lake subdivision gwinnett county georgia."

It did not take much investigating to follow that data trail to Thelma Arnold, a 62-year-old widow who lives in Lilburn, Ga., frequently researches her friends' medical ailments and loves her three dogs. "Those are my searches," she said, after a reporter read part of the list to her.

...

A Face Is Exposed for AOL Searcher No. 4417749 [New York Times, August 9, 2006] Ms. Arnold says she loves online research, but the disclosure of her searches has left her disillusioned. In response, she plans to drop her AOL subscription. "We all have a right to privacy," she said. "Nobody should have found this all out."

SIGKDD 2006 Tutorial, August 2006

Constit











- Ideally, we want an algorithm that discloses only the query result, and only to the requesting party. (In practice, we need some extra disclosure.)
- How do we design algorithms that compute queries while preserving data privacy?

SIGKDD 2006 Tutorial, August 2006

(B) come

• How do we measure privacy (this extra disclosure)?

















Disclosure Limitations

- Ideally, we want a solution that discloses as much statistical information as possible while preserving privacy of the individuals who contributed data.
- How do we design algorithms that allow the "largest" set of queries that can be disclosed while preserving data privacy?
- How do we measure disclosure?

SIGKDD 2006 Tutorial, August 2006

(B) ----

7

This Tutorial: Statistical Methods

- Privacy-preserving data analysis
- Privacy-preserving data publishing

Goal:

• Rather than talk about everything superficially, but nothing in-depth, make hard choices

SIGKDD 2006 Tutorial, August 2006

Const

(B) c...

Caveats:

• Not a comprehensive survey $\ensuremath{\boldsymbol{\otimes}}$

What is Left Out?

- Work on secure multi-party computation (secure join, secure intersection, homomorphic encryption, certificate revocation, etc.)
- Architectural and language issues (Hippocratic databases, P3P, etc.)
- Disclosure control (statistical databases, auditing, database queries, etc.)
- Privacy through distributed data mining

Resources

• See excellent tutorials by Rakesh Agrawal and Chris Clifton; keynote talk by Srikant Ramakrishnan at this conference.























The Problem

- How to randomize data such that
 - we can build a good data mining model (utility)
 - while preserving privacy at the record level (privacy)?

SIGKDD 2006 Tutorial, August 2006

® ----



Motivation: A Social Survey

- Measures opinions, attitudes, behavior
- Problem: Questions of a sensitive nature
 - Examples: sexuality, incriminating questions, embarrassing questions, threatening questions, controversial issues, etc.
 - The "non-cooperative" group leads to errors in surveys and inaccurate data

SIGKDD 2006 Tutorial, August 2006

(B) -----

• Even though privacy is guaranteed, skepticism prevails























Tutorial Outline

- Untrusted data collector
 - Randomized response [W65]
 - The search for a good privacy definition

SIGKDD 2006 Tutorial, August 2006

- Interval privacy [AS00]
- Mutual information [AA01]
- (α , β) privacy breach [EGS03]
- Comments
- Trusted data collector

(B) c....

Interval Privacy [AS00]

[Agrawal and Srikant; SIGMOD 2000]

Idea: Clients share randomized version of their data.

Intuition: Randomized response.

Randomization:

• For a numerical attribute value x, share value z=x+y, where y is drawn from some known distribution

SIGKDD 2006 Tutorial, August 2006

(B) come





Interval Privacy: Example

- Add a random value between -30 and +30 to age.
- If randomized value is 60
 - We know with 90% confidence that age is between 33 and 87.
- Interval width is the amount of privacy.
 - Example:
 - Interval width 54 with 90% confidence
 - Interval width 60 with 100% confidence

SIGKDD 2006 Tutorial, August 2006

(B) Canad

15













Reconstruction: Iterative Algorithm
f _x ⁰ := Uniform density j := 0 // Iteration number repeat
$f_{X}^{j+1}(a) := \frac{1}{n} \sum_{i=1}^{n} \frac{f_{Y}((x_{i}+y_{i})-a)f_{X}^{j}(a)}{\int_{-}^{n} f_{Y}((x_{i}+y_{i})-a)f_{X}^{j}(a)}$
until (stopping criterion met)
 Other approach: Assume parametric distribution
 Perform MLE of distribution parameters through the EM Algorithm [AA01]
SIGKDD 2006 Tutorial, August 2006











- A random variable distributed uniformly between [0,1] has half as much privacy as if it were distributed in [0,2]
- In general: If $f_B(x)=2f_A(2x)$ then B offers half as much privacy as A
 - Think of A as B stretched out at twice the length

SIGKDD 2006 Tutorial, August 2006

(B) com

(B) ---

(B) ---

• Need a privacy measure that captures this intuition



- Differential entropy h(X): $h(X) = -\int_{\Omega X} f_X(x) \log f_X(x) dx$
- Examples:
 - X is uniformly distributed between 0 and 1: h(X)=0. • X is uniformly distributed between 0 and a: $h(X) = \log_2(a)$.
- Random variables with less uncertainty than U[0,1] have negative differential entropy
 Random variables with more uncertainty than U[0,1] have positive differential entropy

SIGKDD 2006 Tutorial, August 2006

Proposed Measure

• Propose $\Pi(X)=2^{h(X)}$ as measure of privacy for attribute X

- Examples:
- Uniform U between 0 and 1: $\Pi(U)=2^{\log_2(1)}=2^0=1$ Uniform U between 0 and a: $\Pi(U)=2^{\log_2(a)}=a$
- In general, II(A) denotes the length of an interval over which a uniformly distributed random variable has as much uncertainty as A.
- Example:
 - $\Pi(X)=2$: X has as much privacy as a random variable distributed uniformly in an interval of length 2

Conditional Privacy

• Conditional privacy takes the additional information in perturbed values into account:

 $h(X \mid Z) = -\int_{\Omega X, Z} f_{X, Z}(x, z) \log f_{X \mid Z = z}(x) dx dz$

• Average conditional privacy of X given Z: $\Pi(X|Z)=2^{h(X|Z)}$

SIGKDD 2006 Tutorial, August 2006

(B) com







SIGKDD 2006 Tutorial, August 2006

Caveat: Privacy Preserved Only On Average Example:

- Example: $f_x(x) = 0.5, 0 \le x \le 1$ $f_x(x) = 0.5, 4 \le x \le 5$ $f_y(x) = 0, \text{ otherwise}$ Uniform noise Y in [0,1] Assume sensitive property: "X<= 0.01." (prior probability: 0.5%) If Z in [-1, -0.99], the posterior probability P[X <= 0.01 | Z = Z] = 1. However, Z in [-1, -0.99] is unlikely (only one in 100,000 records) \rightarrow not much privacy loss
- Caveat:
 - Every time this occurs the property "X <= 0.01" is fully disclosed.
 The mutual information, being an average measure, does not notice this rare disclosure.

SIGKDD 2006 Tutorial, August 2006



- Untrusted data collector
 - Randomized response [W65]
 - The search for a good privacy definition

SIGKDD 2006 Tutorial, August 2006

- Interval privacy [AS00]
- Mutual information [AA01]
- (α,β) privacy breach [EGS03]
- Comments
- Trusted data collector

(B) ----

(B) c...



































































Theorem:

• If randomization operator **R** is at most γamplifying, and if: B 1 a γ

$$\alpha < \frac{\beta}{\alpha} \cdot \frac{1-\alpha}{1-\beta}$$

• Then, revealing R(X) to the server will never cause an α -to- β privacy breach.

SIGKDD 2006 Tutorial, August 2006

(B) -----





Constit



































The Unbiased Estimators

• Given randomized partial supports, we can estimate original partial supports:

$$\vec{s}_{est} = Q \cdot \vec{s}'$$
, where $Q = P^{-1}$

• Covariance matrix for this estimator:

$$\operatorname{Cov} \vec{s}_{\mathsf{est}} = \frac{1}{|T|} \sum_{l=0}^{k} s_l \cdot Q D[l] Q^T,$$

where
$$D[I]_{i,j} = P_{i,l} \cdot \delta_{i=j} - P_{i,l} \cdot P_{j,l}$$

ite it, substitute s_l with $(s_{est})_l$.

(B) ---

To estimate it, substitute s_l with (s_{est})_l.
 Special case: estimators for support and its variance

32







- Untrusted data collector
 - Randomized response [W65]
 - The search for a good privacy definition

SIGKDD 2006 Tutorial, August 2006

- Interval privacy [AS00]
- Mutual information [AA01]
- (α , β) privacy breach [EGS03]
- Comments
- Trusted data collector

(B) card



Extensions: (s,α,β) Privacy Breach [AST05]

[Agrawal, Srikant, Thomas; SIGMOD 2005]

- Consider the following class of randomization operators:
 - Each attribute value is retained with probability p and replaced with probability (1-p) with a value selected from a replacing distribution

Example: Uniform perturbation

• Replacing distribution is the uniform distribution on the domain

SIGKDD 2006 Tutorial, August 2006

🛞 c...

(B) ---

(s,α,β) Privacy Breach (Contd.)

- Consider the following probabilities:
 - $P_f[X \text{ in } S] = p_S$, where P_f is the a priori distribution • $P_g[Y \text{ in } S] = m_S$, where P_g is the replacing distribution.
- Define the *relative a priori probability* of event S as p_s/m_s.
- Intuition: How frequent is S in its a priori distribution compared to the replacing distribution?





(\mathbf{s},α,β) Privacy Breach (Contd.)

• Theorem [AST05]: Uniform perturbation applied to a single column is secure against a (s, α, β) privacy breach if

$$s < \frac{(\beta - \alpha)(1 - p)}{(1 - \beta)p}$$

(Recall: p is probability not to pick from randomizing distribution)

SIGKDD 2006 Tutorial, August 2006

🛞 c...







• Observation:

- Original data could be correlated.
- Noise is not correlated.
- Similar observation by Kargupta and Datta [ICDM 2003]

SIGKDD 2006 Tutorial, August 2006

Const









What Happened?

Original data:

- Correlated.
- If we remove half the attributes, the actual information loss might be much smaller
- Noise:
- Uncorrelated
- Variance evenly distributed across attributes
- If we remove half the attributes, the actual loss in noise should be 50%

SIGKDD 2006 Tutorial, August 2006

(B) ---



SIGKDD 2006 Tutorial, August 2006

(B) com









Disclosure Limitations

- Ideally, we want a solution that discloses as much statistical information as possible while preserving privacy of the individuals who contributed data.
- How do we design algorithms that compute the "largest" set of queries that can be disclosed while preserving data privacy?

SIGKDD 2006 Tutorial, August 2006

(B) come

(B) --

(B) c...

• How do we measure privacy?

Goals

- Safe from attackers who try to learn customers' identities or sensitive information
- Useful for a wide range of statistical analyses
- Easy for users to analyze with standard statistical methods
 - Just load the published dataset into your favorite analysis tool

[Reiter, Chance 17(3), 2004]

Why is Disclosure Bad?

- Violation of laws and thus subject to legal action
- Lose the trust of the public (no future participants)
- Data of dubious quality (since participants are afraid that their privacy is threatened)

[Reiter, Chance 17(3), 2004]

				1
SSN	Zip	Age	Nationality	Disease
631-35-1210	13053	28	Russian	Heart
051-34-1430	13068	29	American	Heart
120-30-1243	13068	21	Japanese	Viral
070-97-2432	13053	23	American	Viral
238-50-0890	14853	50	Indian	Cancer
265-04-1275	14853	55	Russian	Heart
574-22-0242	14850	47	American	Viral
388-32-1539	14850	59	American	Viral
005-24-3424	13053	31	American	Cancer
248-223-2956	13053	37	Indian	Cancer
221-22-9713	13068	36	Japanese	Cancer
615-84-1924	13068	32	American	Cancer



		•		
Zip	Age	Nationality	Disease	Medical Pecords of a
13053	28	Russian	Heart	hospital near Ithaca
13068	29	American	Heart	serving patients from
13068	21	Japanese	Viral	 Freeville (13068)
13053	23	American	Viral	 Dryden (13053)
14853	50	Indian	Cancer	 Ithaca (14850, 1485
14853	55	Russian	Heart	• Iulaca (17050, 1705
14850	47	American	Viral	
14850	59	American	Viral	
13053	31	American	Cancer	
13053	37	Indian	Cancer	
13068	36	Japanese	Cancer	
13068	32	American	Cancer	











Zip	Age	Nationality	Disease	Base Table:
13053	28	Russian	Heart	Medical Records of a
13068	29	American	Heart	hospital near Ithaca
13068	21	Japanese	Viral	Serving patients from
13053	23	American	Viral	Dryden (13053), and
14853	50	Indian	Cancer	Ithaca (14850, 14853
14853	55	Russian	Heart	
14850	47	American	Viral	 The combination
14850	59	American	Viral	{Zip, Age, Nationality
13053	31	American	Cancer	identifier
13053	37	Indian	Cancer	Disease is the
13068	36	Japanese	Cancer	sensitive attribute
13068	32	American	Cancer	



K-Anonymity [Sweeney02]

- Generalize, modify, or distort quasi-identifier values so that no individual is uniquely identifiable from a group of *k*
- In SQL, table T is k-anonymous if each SELECT COUNT(*) FROM T GROUP BY Quasi-Identifier is ≥ k
- Parameter k indicates the "degree" of anonymity

```
SIGKDD 2006 Tutorial, August 2006
```

(B) ---

K-Anonymity

• There are at least k tuples sharing the same values for each combination of the quasi-identifiers.

SIGKDD 2006 Tutorial, August 2006

(B) come

- Techniques
 - Generalizing non-sensitive attributes
 - Tuple Suppression
 - Data Swapping
 - Randomization













• Generalization Property: If T is k-anonymous with respect to a set of attributes, then it is kanonymous with respect to any generalization of these attributes.

Hospital Patients

1/21/76 Male 5371* Flu 1/21/76 Male 5370* Broken Arm 2/28/76 Male 5370* Bronchitis 4/13/86 Female 5371* Hepatitis 4/13/86 Female 5370* Sprained Anklee 2/28/76 Female 5370* Hang Nail	DOB	Sex	Zipcode	Disease
1/21/76 Male 5370* Broken Arm 2/28/76 Male 5370* Bronchitis 4/13/86 Female 5371* Hepatitis 4/13/86 Female 5370* Sprained Anklee 2/28/86 Female 5370* Hang Nail	1/21/76	Male	5371*	Flu
2/28/76 Male 5370* Bronchitis 4/13/86 Female 5371* Hepatitis 4/13/86 Female 5370* Sprained Ankle 2/28/86 Female 5370* Hang Nail	1/21/76	Male	5370*	Broken Arm
4/13/86 Female 5371* Hepatitis 4/13/86 Female 5370* Sprained Ankle 2/28/86 Female 5370* Hang Nail	2/28/76	Male	5370*	Bronchitis
4/13/86 Female 5370* Sprained Ankle 2/28/86 Female 5370* Hang Nail	4/13/86	Female	5371*	Hepatitis
2/28/86 Female 5370* Hang Nail	4/13/86	Female	5370*	Sprained Ankle
	2/28/86	Female	5370*	Hang Nail

Some Simple Observations

- Generalization Property
- Rollup Property: If attribute set P is a generalization of Q, counts grouped by P can be computed directly from the counts grouped by Q.

DOB	Sex	Zipcode	Disease
1/21/76	Male	5371*	Flu
1/21/76	Male	5370*	Broken Arm
2/28/76	Male	5370*	Bronchitis
4/13/86	Female	5371*	Hepatitis
4/13/86	Female	5370*	Sprained Ankle
2/28/86	Female	5370*	Hang Nail

SIGKDD 2006 Tutorial, August 2006

Some Simple Observations

- Generalization Property
- Rollup Property
- Subset Property: If T is k-anonymous with respect to attribute set Q, then T is k-anonymous with respect to $P \subseteq Q$.

Hospital Patients

DOB	Sex	Zipcode	Disease
1/21/76	Male	537**	Flu
1/21/76	Male	537**	Broken Arm
2/28/76	Male	537**	Bronchitis
4/13/86	Female	537**	Hepatitis
4/13/86	Female	537**	Sprained Ankle
2/28/86	Female	537**	Hang Nail
S	SIGKDD 200	6 Tutorial, A	ugust 2006

Some Simple Observations	
 Generalization Property Rollup Property Subset Property → Frequent Itemsets 	9
SIGVED 2006 Tuberol August 2006	Consel University





Basic Incognito Algorithm

- Finds all k-anonymous full-domain generalizations
- Begins by checking k-anonymity with respect to single-attribute subsets of quasi-identifier. Then iteratively checks larger subsets. (*Subset Property*)
- Each iteration has two phases:
 - Breadth-first search (*Rollup Property*)
 - Candidate graph construction

SIGKDD 2006 Tutorial, August 2006

(B) --





SIGKDD 2006 Tutorial, August 2006

Constitu



Example	Micro	odata	a		
	Zip	Age	Nationality	Disease	
	13053	28	Russian	Heart	
	13068	29	American	Heart	
	13068	21	Japanese	Viral	
	13053	23	American	Viral	
	14853	50	Indian	Cancer	
	14853	55	Russian	Heart	
	14850	47	American	Viral	
	14850	59	American	Viral	
	13053	31	American	Cancer	
	13053	37	Indian	Cancer	
	13068	36	Japanese	Cancer	
	13068	32	American	Cancer	
	SI	GKDD 2006	6 Tutorial, August 2006		Constit University



Zip	Age	Nationality	Diseas
130**	<30	*	Heart
130**	<30	*	Heart
130**	<30	*	Viral
130**	<30	*	Viral
1485*	>40	*	Cancer
1485*	>40	*	Heart
1485*	>40	*	Viral
1485*	>40	*	Viral
130**	30-40	*	Cancer
130**	30-40	*	Cancer
130**	30-40	*	Cancer
130**	30-40	*	Cancer



Attacks on K-Anonymity [Ohrn, Ohno-Machado; Artif Intell Med. 15(3), 1999] [Machanavajjhala, Gehrke, Kifer, Venkitasubramaniam; ICDE 2006] • K-Anonymity does not protect against some simple attacks

Hon	noge	eneity Atta	ack	
Zip	Age	Nationality	Disease	 Alice's neighbor Bob is
130**	<30	*	Heart	in the hospital.
130**	<30	*	Heart	 Alice knows Bob is 35
130**	<30	*	Viral	vears old and is from
130**	<30	*	Viral	Dryden (13053)
1485*	>40	*	Cancer	Diyden (15055).
1485*	>40	*	Heart	
1485*	>40	*	Viral	 Alice learns that Bob
1485*	>40	*	Viral	has cancer.
130**	30-40	*	Cancer	
130**	30-40	*	Cancer	
130**	30-40	*	Cancer	
130**	30-40	*	Cancer	Alice
		SIGKDD 2	2006 Tutorial, /	August 2006

SIGKDD 2006 Tutorial, August 2006

🛞 Constit



kgro	und Knov	wledge	Attack
Age	Occupation	Disease	
<30	*	Heart	
<30	*	Heart	
<30	*	Viral	
<30	*	Viral	Allee
>40	*	Cancer	 Alice's friend Umeko is in the table
>40	*	Heart	 Alice knows Umeko is 24, a
>40	*	Viral	Japanese, living in Freeville
>40	*	Viral	(13068)
30-40	*	Cancer	Japanese have extremely low
30-40	*	Cancer	incidence of heart disease
30-40	*	Cancer	Alice learns Umeko
30-40	*	Cancer	has a viral infection
	Age <30	Age Occupation <30	Age Occupation Disease <30



Data Publishing Desiderata Need to defend against attacks based on background knowledge Need to permit efficient sanitization algorithms Guarantee understood by a lay person



SIGKDD 2006 Tutorial, August 2006

(B) ---



Privacy Definition (1)

- Positive Disclosure: Posterior Belief > $1-\delta$
- Negative Disclosure: Posterior Belief < δ

BUT:

- Not all positive disclosures are bad
 OK to disclose Bob is healthy
- Not all negative disclosures are bad
 OK to disclose Bob does not have Ebola

Privacy Definition (2)

• Bayes-optimal privacy: After publishing we have Posterior belief ~ prior belief

SIGKDD 2006 Tutorial, August 2006

🛞 canal

(B) ---

- Example instantiation: α-to-β privacy breach definition
 Prior Belief < α and posterior Belief > β OR
 Prior Belief >1- α and posterior Belief <1-β
- Automatically eliminates homogeneity attack
 Homogeneity → Posterior belief = 1

SIGKDD 2006 Tutorial, August 2006

Bayes-Optimal Privacy– Drawbacks Insufficient knowledge Nobody knows the complete joint distribution Adversary's knowledge unknown Data publisher does not know how much the adversary knows Omputational intractability Checking for every (q,s) pair ...













3-D)iver:	se Microo	data	
Zip	Age	Nationality	Disease	 Bob is 35 years old
1306*	<=40	*	Heart	and is from Dryden
1306*	<=40	*	Viral	(13053)
1306*	<=40	*	Cancer	(19099).
1306*	<=40	*	Cancer	
1485*	>40	*	Cancer	• Umeko is 24, a
1485*	>40	*	Heart	Japanese from
1485*	>40	*	Viral	Freeville (13068)
1485*	>40	*	Viral	 Japanese have
1305*	<=40	*	Heart	<i>extremely low</i> <i>incidence</i> of heart
1305*	<=40	*	Viral	
1305*	<=40	*	Cancer	disease
1305*	<=40	*	Cancer	
		SIGKDD	2006 Tutorial,	Condi Detendry August 2006



L-Diversity Revisited

- L -Diversity: Every group has at least L *well represented* groups
- <u>Note</u>: L-diversity does not protect against adversaries having arbitrary background knowledge.

Q	S
q *	S
q *	<i>.</i>
q *	×Z
q *	85
q *	84

(B) ---

(B) ---

• <u>But</u>: L-diversity increases the bar.

L-Diversity: Summary

- Defends against background knowledge attacks and homogeneity attacks
 - L-Diversity ensures diversity
 - Gives guarantees against "unknown" background knowledge

SIGKDD 2006 Tutorial, August 2006

- Can model don't care values ("person is healthy")
- Guarantee understood by a lay person
- "At least L different values"
- Permits efficient sanitization algorithms
 - Bayes-optimal definition is not monotone
- L-Diversity and (c,k)-recursive L-Diversity are monotone
 Experiments show that little utility is lost compared to k-
- Experiments show that little utility is lost compared to kanonymity



[Wong, Li, Fu, and Wang; KDD 2006] Defends against homogeneity attacks

- Dataset is α-deassociative for a value s: Relative frequency of s within its group is <= α.
- (α ,k)-anonymity: Dataset is k-anonymous and α -deassociative for all values in the domain of a sensitive attribute

SIGKDD 2006 Tutorial, August 2006

🛞 canal

(B) c...

(B) ---

What About Other Knowledge?

• If Carol and David are both sick and if Carol has the flu, then David also has the flu: $t_{Carol}[Disease] = Influenza \rightarrow t_{David}[Disease] = Influenza$

- Other types of knowledge?
- Language for background knowledge?
- Complexity, guarding against worst-case disclosure?

SIGKDD 2006 Tutorial, August 2006

The Curse of Dimensionality [A05]

[Aggarwal; VLDB 2005]

- Curse of dimensionality
- Formal analysis that shows with increasing dimensionality all information in the data is lost in order to achieve k-anonymity



- M(S): Maximum Euclidean distance between any pair of points in S
- M(D): Maximum Euclidean distance between any pair of points in whole database S
- Relative condensation loss L(S) through k-anonymization L(S) = M(S)/M(D)
- <u>Theorem [A05]</u>: For any set S of points to be kanonymous, the relative condensation loss goes to 1 with increasing dimensionality:

 $\lim_{d\to\infty} E[M(S)/M(D)] = 1$ SIGKDD 2006 Tutorial, August 2006

Condit

(B) com

(B) com

Protection Against An Adversary

[Aggarwal, Pei, and Zhang; KDD 2006]

• Problem: Any attribute might be sensitive; need to defend against inference attacks based on rules learned from the data

• Example: [Type = Manager and DEP = Toy] → Salary > 100k; Confidence of rule: 100% Simple suppression of private values insufficient.

• Approach: Make strong rules weaker

SIGKDD 2006 Tutorial, August 2006

Open Problems

- Tradeoff of utility versus privacy
 See Kifer et al, SIGMOD 2006, Levefre et al, KDD 2006, Xu et al., KDD 2006
- Re-publication
- Theory of learning from summaries
- Multi-round protocols
- Formalization of classes of background knowledge
- Location privacy



Thanks

Rakesh Agrawal, Chris Clifton, Wenliang Du, Cynthia Dwork, Alexandre Evfimievski, Ashwin Machanavajjhala, Daniel Kifer, Lucja Kot, Kristen Lefevre, David Martin, Kobbi Nissim, Muthuramakrishnan Venkitasubramaniam, Ramakrishnan Srikant, Walker White

For an annotated list of references for all the topics see (soon :-) <u>http://www.cs.cornell.edu/database/privacy</u>

SIGKDD 2006 Tutorial, August 2006

(B) c.....

Questions? johannes@cs.cornell.edu.