KDE Resampling Algorithm

Numerical Results

Conclusions

Privacy-Preserving Sharing of Horizontally-Distributed Private Data for Constructing Accurate Classifiers

Vincent Y. F. Tan¹ See-Kiong Ng²

¹Department of Electrical Engineering and Computer Science Massachusetts Institute of Technology

²Institute for Infocomm Research (I²R), Singapore

PinKDD 2007: 1st ACM SIGKDD Int. Workshop on Privacy, Security and Trust in KDD

KDE Resampling Algorithm

Numerical Results

・ コ ト ・ 雪 ト ・ 目 ト ・

э

Conclusions

Outline

Introduction

- Motivation
- Related Work
- Problem Statement
- 2 KDE Resampling Algorithm
 - Resampling from Reconstructed PDF
 - Properties of Representative Samples
- 3 Numerical Results
 - Performance Metrics, Datasets, Classification Techniques
 - Results of Distributed Experiments
- 4 Conclusions
 - Summary
 - Further Work and Acknowledgments

KDE Resampling Algorithm

Numerical Results

・ロット (雪) ・ (日) ・ (日)

э

Conclusions

Outline

Introduction

- Motivation
- Related Work
- Problem Statement
- 2 KDE Resampling Algorithm
 - Resampling from Reconstructed PDF
 - Properties of Representative Samples
 - 3 Numerical Results
 - Performance Metrics, Datasets, Classification Techniques
 - Results of Distributed Experiments
- 4 Conclusions
 - Summary
 - Further Work and Acknowledgments

KDE Resampling Algorithm

Numerical Results

・ロット (雪) ・ (日) ・ (日)

э

Conclusions

Outline



- Motivation
- Related Work
- Problem Statement
- 2 KDE Resampling Algorithm
 - Resampling from Reconstructed PDF
 - Properties of Representative Samples
- 3 Numerical Results
 - Performance Metrics, Datasets, Classification Techniques
 - Results of Distributed Experiments
- 4) Conclusion
 - Summary
 - Further Work and Acknowledgments

KDE Resampling Algorithm

Numerical Results

・ ロ ト ・ 雪 ト ・ 雪 ト ・ 日 ト

э

Conclusions

Outline



- Motivation
- Related Work
- Problem Statement
- 2 KDE Resampling Algorithm
 - Resampling from Reconstructed PDF
 - Properties of Representative Samples
- 3 Numerical Results
 - Performance Metrics, Datasets, Classification Techniques
 - Results of Distributed Experiments
- 4 Conclusions
 - Summary
 - Further Work and Acknowledgments

 KDE Resampling Algorithm

Numerical Results

Conclusions

Motivation

Motivation

• Example: Two hospitals seek to construct a global classifier based on existing private patient data.



- Pooling data results in a more accurate global classifier.
- Instead of building local classifiers on limited data at each hospital site.

 KDE Resampling Algorithm

Numerical Results

Conclusions

Motivation

Motivation

• Example: Two hospitals seek to construct a global classifier based on existing private patient data.



- Pooling data results in a more accurate global classifier.
- Instead of building local classifiers on limited data at each hospital site.

KDE Resampling Algorithm

Numerical Results

・ロト ・ 同ト ・ ヨト ・ ヨト

э

Conclusions

Motivation

Motivation

• We choose to randomize the data before transmitting the data to a centralized server.



• Can privacy + accuracy co-exist in a distributed scenario?

KDE Resampling Algorithm

Numerical Results

・ロト ・ 同ト ・ ヨト ・ ヨト

э

Conclusions

Motivation

Motivation

• We choose to randomize the data before transmitting the data to a centralized server.



• Can privacy + accuracy co-exist in a distributed scenario?

KDE Resampling Algorithm

Numerical Results

・ ロ マ ・ 雪 マ ・ 雪 マ ・ 日 マ

3

Conclusions

Motivation

Randomization in PPDM

Randomization

- Randomize locally to give sanitized data for sharing.
- Similar to horizontally-partitioned scenario (Du et al. 2004).
- Pool data to form a larger training set to construct a global classifier.
- Challenges: Two conflicting concerns.
 - Confidentiality of the private information
 - Otility of the aggregate statistics.
- New randomization algorithm Kernel Density Estimate (KDE) Resampling.

KDE Resampling Algorithm

Numerical Results

Conclusions

Motivation

Randomization in PPDM

Randomization

- Randomize locally to give sanitized data for sharing.
- 2 Similar to horizontally-partitioned scenario (Du et al. 2004).
- Pool data to form a larger training set to construct a global classifier.
- Challenges: Two conflicting concerns.
 - Confidentiality of the private information
 - **2** Utility of the aggregate statistics.
- New randomization algorithm Kernel Density Estimate (KDE) Resampling.

KDE Resampling Algorithm

Numerical Results

Conclusions

Motivation

Randomization in PPDM

Randomization

- Randomize locally to give sanitized data for sharing.
- 2 Similar to horizontally-partitioned scenario (Du et al. 2004).
- Pool data to form a larger training set to construct a global classifier.
- Challenges: Two conflicting concerns.
 - Confidentiality of the private information
 - **2** Utility of the aggregate statistics.
- New randomization algorithm Kernel Density Estimate (KDE) Resampling.

KDE Resampling Algorithm

Numerical Results

・ロ ・ ・ 一 ・ ・ 日 ・ ・ 日 ・

э

Conclusions

Related Work

Randomization Techniques

- Atallah et al. (1999) first considered the data sanitization problem.
- Agrawal and Srikant (2000) suggested adding IID noise.
- Noise addition has since been shown to be insecure (Kargupta et al. 2003).
- Distance-Preserving approaches have been suggested (Liu et al. 2006, Chen et al. 2005, Olivera et al. 2007).
 - But we can derive bounds on the data.
- We suggest a Non-Distance-Preserving randomization approach, which also has classification accuracy.

KDE Resampling Algorithm

Numerical Results

Conclusions

Related Work

Randomization Techniques

- Atallah et al. (1999) first considered the data sanitization problem.
- Agrawal and Srikant (2000) suggested adding IID noise.
- Noise addition has since been shown to be insecure (Kargupta et al. 2003).
- Distance-Preserving approaches have been suggested (Liu et al. 2006, Chen et al. 2005, Olivera et al. 2007).
 - But we can derive bounds on the data.
- We suggest a Non-Distance-Preserving randomization approach, which also has classification accuracy.

KDE Resampling Algorithm

Numerical Results

Conclusions

Related Work

Randomization Techniques

- Atallah et al. (1999) first considered the data sanitization problem.
- Agrawal and Srikant (2000) suggested adding IID noise.
- Noise addition has since been shown to be insecure (Kargupta et al. 2003).
- Distance-Preserving approaches have been suggested (Liu et al. 2006, Chen et al. 2005, Olivera et al. 2007).
 - But we can derive bounds on the data.
- We suggest a Non-Distance-Preserving randomization approach, which also has classification accuracy.

KDE Resampling Algorithm

Numerical Results

Conclusions

Related Work

Non-randomization Techniques

- k-anonymization used to generalize databases for preserving privacy (Sweeney, 1998).
- Secure Multi-Party Computation (SMC) techniques are more accurate but higher communication overhead.
- Distributed Clustering via optimization of information theoretic quantities (Merugu and Ghosh, 2005).

KDE Resampling Algorithm

Numerical Results

イロト イポト イヨト イヨト 三日

Conclusions

Problem Statement

Problem Definition and Notation

- Private data is stored in *d*-dim row vectors $\mathbf{x}_1, \ldots, \mathbf{x}_N$.
- Associated with *N* targets (class labels) $t_1, \ldots t_N$.
- *L* distributed data sites (private) and 1 centralized (untrusted) server.
- $\mathbf{x}_{(I)} \in \mathbb{R}^{N_I \times d}$ contains the N_I data vectors at site *I*.
- Assume: Vectors in $\mathbf{x}_{(l)}$ are drawn from IID RVs with PDF $f_{\mathbf{x}_{(l)}}(\mathbf{x}_{(l)})$.

KDE Resampling Algorithm

Numerical Results

Conclusions

◆□▶ ◆□▶ ◆∃▶ ◆∃▶ = のへの

Problem Statement

Problem Definition and Notation

- Private data is stored in *d*-dim row vectors $\mathbf{x}_1, \ldots, \mathbf{x}_N$.
- Associated with *N* targets (class labels) $t_1, \ldots t_N$.
- L distributed data sites (private) and 1 centralized (untrusted) server.
- $\mathbf{x}_{(I)} \in \mathbb{R}^{N_I \times d}$ contains the N_I data vectors at site *I*.
- Assume: Vectors in $\mathbf{x}_{(l)}$ are drawn from IID RVs with PDF $f_{\mathbf{x}_{(l)}}(\mathbf{x}_{(l)})$.

KDE Resampling Algorithm

Numerical Results

Conclusions

Problem Statement

Problem Definition and Notation

- Private data is stored in *d*-dim row vectors $\mathbf{x}_1, \ldots, \mathbf{x}_N$.
- Associated with *N* targets (class labels) t_1, \ldots, t_N .
- L distributed data sites (private) and 1 centralized (untrusted) server.
- $\mathbf{x}_{(I)} \in \mathbb{R}^{N_I \times d}$ contains the N_I data vectors at site *I*.
- Assume: Vectors in $\mathbf{x}_{(l)}$ are drawn from IID RVs with PDF $f_{\mathbf{x}_{(l)}}(\mathbf{x}_{(l)})$.

KDE Resampling Algorithm

Numerical Results

・ロット (雪) (日) (日) (日)

Conclusions

Problem Statement

Problem Definition and Notation

Task

Find a randomization scheme for site I, $R_{I}(\cdot)$ s.t.

$$\mathbf{y}_{(l)} = R_l(\mathbf{x}_{(l)}), \qquad 1 \le l \le L.$$

- y = {y_(l)}^L_{l=1} is then sent to the centralized server for global classification.
- We will demonstrate empirically that for our choice of $R_l(\cdot)$,

 $P_{rand}(err) \approx P_{ori}(err).$

Why non-distance-preserving?

KDE Resampling Algorithm

Numerical Results

Conclusions

Problem Statement

Problem Definition and Notation

Task

Find a randomization scheme for site I, $R_{I}(\cdot)$ s.t.

$$\mathbf{y}_{(I)} = R_I(\mathbf{x}_{(I)}), \qquad 1 \le I \le L.$$

- y = {y_(l)}^L_{l=1} is then sent to the centralized server for global classification.
- We will demonstrate empirically that for our choice of $R_{l}(\cdot)$,

 $P_{rand}(err) \approx P_{ori}(err).$

Why non-distance-preserving?

KDE Resampling Algorithm

Numerical Results

Conclusions

Problem Statement

Problem Definition and Notation

Task

Find a randomization scheme for site I, $R_{I}(\cdot)$ s.t.

$$\mathbf{y}_{(I)} = R_I(\mathbf{x}_{(I)}), \qquad 1 \le I \le L.$$

- y = {y_(l)}^L_{l=1} is then sent to the centralized server for global classification.
- We will demonstrate empirically that for our choice of $R_{l}(\cdot)$,

 $P_{rand}(err) \approx P_{ori}(err).$

Why non-distance-preserving?

Numerical Results

Conclusions

Problem Statement

Why Non-Distance-Preserving Randomization?

If Distance-Preserving Randomization (Random Projection-Based Multiplication) was used:

- The randomized data can be vulnerable to disclosure (Caetano 2004)
 - True iff (d + 1) points do not lie in a (d 1)-dimensional vector subspace.
- Furthermore in a L = 2 site scenario (details in our paper):
 - One can lower bound the norm of all the columns of the data matrix x₍₁₎ given an estimate of the norm of one of the columns of x₍₂₎.

Numerical Results

Conclusions

Problem Statement

Why Non-Distance-Preserving Randomization?

If Distance-Preserving Randomization (Random Projection-Based Multiplication) was used:

- The randomized data can be vulnerable to disclosure (Caetano 2004)
 - True iff (*d* + 1) points do not lie in a (*d* 1)-dimensional vector subspace.
- Furthermore in a L = 2 site scenario (details in our paper):
 - One can lower bound the norm of all the columns of the data matrix x₍₁₎ given an estimate of the norm of one of the columns of x₍₂₎.

Numerical Results

Conclusions

Problem Statement

Why Non-Distance-Preserving Randomization?

If Distance-Preserving Randomization (Random Projection-Based Multiplication) was used:

- The randomized data can be vulnerable to disclosure (Caetano 2004)
 - True iff (d + 1) points do not lie in a (d 1)-dimensional vector subspace.
- Furthermore in a L = 2 site scenario (details in our paper):
 - One can lower bound the norm of all the columns of the data matrix x₍₁₎ given an estimate of the norm of one of the columns of x₍₂₎.

KDE Resampling Algorithm

Numerical Results

Conclusions

Problem Statement

Why Non-Distance-Preserving Randomization?

Hence, if Distance-Preserving Randomization was used:

- Potential security breach.
- Intuition: Along with the preservation of distances, the order of the samples is also preserved.
- Hence, we propose a non-distance-preserving approach to circumvent the problem.



KDE Resampling Algorithm

Numerical Results

Conclusions

Problem Statement

Why Non-Distance-Preserving Randomization?

Hence, if Distance-Preserving Randomization was used:

- Potential security breach.
- Intuition: Along with the preservation of distances, the order of the samples is also preserved.
- Hence, we propose a non-distance-preserving approach to circumvent the problem.

KDE Resampling Algorithm

Numerical Results

Conclusions

Problem Statement

Why Non-Distance-Preserving Randomization?

Hence, if Distance-Preserving Randomization was used:

- Potential security breach.
- Intuition: Along with the preservation of distances, the order of the samples is also preserved.
- Hence, we propose a non-distance-preserving approach to circumvent the problem.

KDE Resampling Algorithm

Numerical Results

Conclusions

Resampling from Reconstructed PDF

Kernel Density Estimation

- Main Idea: Sample from the reconstructed Kernel Density Estimate (Parzen, 1962) (Devroye, 1985).
- The KDE is a non-parameteric estimate of the PDF.
- Nice properties of the new, representative samples, including asymptotic independence and consistency.

KDE Resampling Algorithm

Numerical Results

Conclusions

Resampling from Reconstructed PDF

Kernel Density Estimation

- Main Idea: Sample from the reconstructed Kernel Density Estimate (Parzen, 1962) (Devroye, 1985).
- The KDE is a non-parameteric estimate of the PDF.
- Nice properties of the new, representative samples, including asymptotic independence and consistency.

KDE Resampling Algorithm

Numerical Results

Conclusions

Resampling from Reconstructed PDF

Kernel Density Estimation

For each data site *I*, we will construct and sample from an estimate of the PDF using x(*I*).

$$\hat{f}_{\mathbf{X}_{(l)}}\left(\mathbf{x}_{(l)};\mathbf{x}_{(l,1)},\ldots,\mathbf{x}_{(l,N_l)}\right) = \frac{1}{N_l}\sum_{j=1}^{N_l} K\left(\mathbf{x}_{(l)} - \mathbf{x}_{(l,j)};\mathbf{h}_l\right),$$



KDE Resampling Algorithm

Numerical Results

Conclusions

Resampling from Reconstructed PDF

Resampling

The random vector

$$\mathbf{X}_{(l)} = \frac{1}{N_l} \sum_{j=1}^{N_l} \mathbf{X}_{(l,j)}$$

is a mixture density.

 KDE does not have to be explicitly constructed for sampling.

KDE Resampling Algorithm

Numerical Results

Conclusions

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Resampling from Reconstructed PDF

Resampling

The random vector

$$\mathbf{X}_{(l)} = \frac{1}{N_l} \sum_{j=1}^{N_l} \mathbf{X}_{(l,j)}$$

is a mixture density.

 KDE does not have to be explicitly constructed for sampling.

KDE Resampling Algorithm

Numerical Results

Conclusions

Resampling from Reconstructed PDF

KDE Resampling Algorithm





KDE Resampling Algorithm

Numerical Results

・ ロ ト ・ 雪 ト ・ 雪 ト ・ 日 ト

э

Conclusions

Properties of Representative Samples

Asymptotic Independence

- Randomized samples are independent of the Original samples as the number of samples $N_I \rightarrow \infty$.
- Probabilistic inference cannot be performed based on the randomized samples y_(l) if N_l is sufficiently large.

• Dependent on how we select **h**_l (Scott's rule).

$$h_{l,i} = \left(\frac{4}{d+2}\right)^{1/(d+4)} N_l^{-1/(d+4)} \hat{\sigma}_{l,i}.$$

KDE Resampling Algorithm

Numerical Results

イロト 不良 とくほ とくほう 二日

Conclusions

Properties of Representative Samples

Asymptotic Independence

- Randomized samples are independent of the Original samples as the number of samples N_I → ∞.
- Probabilistic inference cannot be performed based on the randomized samples y_(l) if N_l is sufficiently large.

• Dependent on how we select **h**_l (Scott's rule).

$$h_{l,i} = \left(\frac{4}{d+2}\right)^{1/(d+4)} N_l^{-1/(d+4)} \,\hat{\sigma}_{l,i}.$$

KDE Resampling Algorithm

Numerical Results

Conclusions

Properties of Representative Samples

Asymptotic Independence

- Randomized samples are independent of the Original samples as the number of samples N_I → ∞.
- Probabilistic inference cannot be performed based on the randomized samples y_(l) if N_l is sufficiently large.
- Dependent on how we select h_l (Scott's rule).

$$h_{l,i} = \left(\frac{4}{d+2}\right)^{1/(d+4)} N_l^{-1/(d+4)} \hat{\sigma}_{l,i}.$$



KDE Resampling Algorithm

Numerical Results

A D > A D > A D > A D >

ж

Conclusions

Properties of Representative Samples

Consistency and Tractability

$$\lim_{N_{l}\to\infty}\mathbb{E}\left[\int\left|\hat{f}_{l}-f_{l}\right|\right]=0,\quad 1\leq l\leq L.$$

- As N_l becomes large, the KDE $\hat{f}_l(\cdot)$ becomes increasingly accurate.
- Treat $\{\mathbf{y}_{(l)}\}_{l=1}^{L}$ as the training data without compromising accuracy.
- KDE approximation algorithm is tractable.

KDE Resampling Algorithm

Numerical Results

Conclusions

Properties of Representative Samples

Consistency and Tractability

$$\lim_{N_{l}\to\infty}\mathbb{E}\left[\int\left|\hat{f}_{l}-f_{l}\right|\right]=0,\quad 1\leq l\leq L.$$

- As N_l becomes large, the KDE f̂_l(·) becomes increasingly accurate.
- Treat $\{\mathbf{y}_{(l)}\}_{l=1}^{L}$ as the training data without compromising accuracy.
- KDE approximation algorithm is tractable.

KDE Resampling Algorithm

Numerical Results

Conclusions

Properties of Representative Samples

Consistency and Tractability

$$\lim_{N_{l}\to\infty}\mathbb{E}\left[\int\left|\hat{f}_{l}-f_{l}\right|\right]=0,\quad 1\leq l\leq L.$$

- As N_l becomes large, the KDE f̂_l(·) becomes increasingly accurate.
- Treat $\{\mathbf{y}_{(l)}\}_{l=1}^{L}$ as the training data without compromising accuracy.
- KDE approximation algorithm is tractable.

KDE Resampling Algorithm

Numerical Results

Conclusions

Properties of Representative Samples

Consistency and Tractability

$$\lim_{N_{l}\to\infty}\mathbb{E}\left[\int \left|\hat{f}_{l}-f_{l}\right|\right]=0,\quad 1\leq l\leq L.$$

- As N_l becomes large, the KDE $\hat{f}_l(\cdot)$ becomes increasingly accurate.
- Treat $\{\mathbf{y}_{(l)}\}_{l=1}^{L}$ as the training data without compromising accuracy.
- KDE approximation algorithm is tractable.

KDE Resampling Algorithm

Numerical Results

・ロット (雪) (日) (日)

Conclusions

Performance Metrics, Datasets, Classification Techniques

Distributed Aggregate Privacy Loss \mathcal{DAPL}

- We define privacy loss is a function of the degree of independence between new, randomized samples and original samples.
- A sufficient condition for asymptotic independence is:
- The expected l_1 distance between \hat{f}_l and f_l tends to zero with $N_l \rightarrow \infty$ (Devroye, 1985).

Definition

The Distributed Aggregate Privacy Loss \mathcal{DAPL} is defined as:

$$\mathcal{DAPL} \stackrel{\triangle}{=} \frac{1}{2} \left(\sum_{l=1}^{L} c_l \mathbb{E} \left[\int \left| \hat{f}_l - f_l \right| \right] \right), \quad c_l \stackrel{\triangle}{=} N_l / N.$$

KDE Resampling Algorithm

Numerical Results

・ コ ト ・ 雪 ト ・ 目 ト ・

Conclusions

Performance Metrics, Datasets, Classification Techniques

Distributed Aggregate Privacy Loss \mathcal{DAPL}

- We define privacy loss is a function of the degree of independence between new, randomized samples and original samples.
- A sufficient condition for asymptotic independence is:
- The expected l_1 distance between \hat{f}_l and f_l tends to zero with $N_l \rightarrow \infty$ (Devroye, 1985).

Definition

The Distributed Aggregate Privacy Loss \mathcal{DAPL} is defined as:

$$\mathcal{DAPL} \stackrel{\triangle}{=} \frac{1}{2} \left(\sum_{l=1}^{L} c_l \mathbb{E} \left[\int \left| \hat{f}_l - f_l \right| \right] \right), \quad c_l \stackrel{\triangle}{=} N_l / N.$$

KDE Resampling Algorithm

Numerical Results

(日)

Conclusions

Performance Metrics, Datasets, Classification Techniques

Distributed Aggregate Privacy Loss \mathcal{DAPL}

- We define privacy loss is a function of the degree of independence between new, randomized samples and original samples.
- A sufficient condition for asymptotic independence is:
- The expected l_1 distance between \hat{f}_l and f_l tends to zero with $N_l \rightarrow \infty$ (Devroye, 1985).

Definition

The Distributed Aggregate Privacy Loss \mathcal{DAPL} is defined as:

$$\mathcal{DAPL} \stackrel{\triangle}{=} \frac{1}{2} \left(\sum_{l=1}^{L} c_{l} \mathbb{E} \left[\int \left| \hat{f}_{l} - f_{l} \right| \right] \right), \quad c_{l} \stackrel{\triangle}{=} N_{l} / N.$$

KDE Resampling Algorithm

Numerical Results

・ロット (雪) (日) (日)

Conclusions

Performance Metrics, Datasets, Classification Techniques

Deterioration of Classification ϕ

We define the classification error as:

$$P(err) \stackrel{ riangle}{=} 1 - \sum_{i=1}^{|\mathcal{C}|} \int_{\Omega_i} p(\xi|\mathcal{C}_i) P(\mathcal{C}_i) d\xi.$$

Definition

The *Deterioration of Classification* ϕ is defined as:

$$\phi \stackrel{\triangle}{=} P_{rand}(err) - P_{ori}(err),$$

• *P*_{ori}(err): Error with original samples for training.

• *P_{rand}(err)*: Error with randomized samples for training.

KDE Resampling Algorithm

Numerical Results

・ロット (雪) (日) (日)

Conclusions

Performance Metrics, Datasets, Classification Techniques

Deterioration of Classification ϕ

• We define the classification error as:

$$P(err) \stackrel{ riangle}{=} 1 - \sum_{i=1}^{|\mathcal{C}|} \int_{\Omega_i} p(\xi|\mathcal{C}_i) P(\mathcal{C}_i) d\xi.$$

Definition

The Deterioration of Classification ϕ is defined as:

$$\phi \stackrel{\triangle}{=} P_{rand}(err) - P_{ori}(err),$$

- *P*_{ori}(err): Error with original samples for training.
- *P_{rand}(err)*: Error with randomized samples for training.

KDE Resampling Algorithm

Numerical Results

Conclusions

Лij

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ - 三 - のへぐ

Performance Metrics, Datasets, Classification Techniques

Datasets

Dataset	#Class	#Dim(<i>d</i>)	#Trg(N)	#Test
Iris	2	4	120	30
SVMGuide1	2	4	3089	4000
Diabetes	2	8	576	192
Breast-Cancer	2	10	512	171
Ionosphere	2	34	263	88

Table: Our five datasets from LIBSVM and the UCI ML Repository.

KDE Resampling Algorithm

Numerical Results

Conclusions

Performance Metrics, Datasets, Classification Techniques

Classification Techniques

- Artificial Neural Networks (ANN) by trained by error backpropagation.
- k-Nearest Neighbors classifier (kNN).
- Naïve Bayes classifier (NB) with each attribute assumed to be Gaussian.



KDE Resampling Algorithm

Numerical Results

イロト イ理ト イヨト イヨト

Conclusions

Results of Distributed Experiments

SVMGuide1 Dataset



(a) (Key: x - ANN, o - kNN, + - NB); (b) DAPL (Key: x - C1, o - C2)

KDE Resampling Algorithm

Numerical Results

イロト イ理ト イヨト イヨト

Conclusions

Results of Distributed Experiments

Breast-Cancer Dataset



(a) (Key: x - ANN, o - kNN, + - NB); (b) DAPL (Key: x - C1, o - C2)

KDE Resampling Algorithm

Numerical Results

Conclusions

Results of Distributed Experiments

Diabetes Dataset



(a) (Key: x - ANN, o - kNN, + - NB);

(b) DAPL (Key: x - C1, o - C2)

イロト イ理ト イヨト イヨト

KDE Resampling Algorithm

Numerical Results

・ ロ ト ・ 雪 ト ・ 雪 ト ・ 日 ト

э

Conclusions

Results of Distributed Experiments

Observations

• Accuracy:

- No correlation between *L* and the Deterioration of Classification ϕ .
- Randomized data is still amenable to data mining tasks regardless of *L*.

• Privacy:

- Increasing trend for the \mathcal{DAPL} because $N_I \downarrow$ with $L \uparrow$.
- A compromise between L and \mathcal{DAPL} .
- Improving the sampling algorithm by selecting optimal h₁.

KDE Resampling Algorithm

Numerical Results

・ロ ・ ・ 一 ・ ・ 日 ・ ・ 日 ・

э

Conclusions

Results of Distributed Experiments

Observations

• Accuracy:

- No correlation between *L* and the Deterioration of Classification ϕ .
- Randomized data is still amenable to data mining tasks regardless of *L*.
- Privacy:
 - Increasing trend for the \mathcal{DAPL} because $N_l \downarrow$ with $L \uparrow$.
 - A compromise between L and \mathcal{DAPL} .
- Improving the sampling algorithm by selecting optimal h_l.

KDE Resampling Algorithm

Numerical Results

・ロ ・ ・ 一 ・ ・ 日 ・ ・ 日 ・

3

Conclusions

Results of Distributed Experiments

Observations

• Accuracy:

- No correlation between *L* and the Deterioration of Classification ϕ .
- Randomized data is still amenable to data mining tasks regardless of *L*.
- Privacy:
 - Increasing trend for the \mathcal{DAPL} because $N_l \downarrow$ with $L \uparrow$.
 - A compromise between L and \mathcal{DAPL} .
- Improving the sampling algorithm by selecting optimal h₁.

KDE Resampling Algorithm

Numerical Results

Conclusions

Summary

- Proposed an algorithm (KDE Resampling) for data sanitization to share private data for distributed classification.
 - Asymptotically independent (Privacy).
 - Consistent (Accuracy).
- max $|P_{rand}(err) P_{ori}(err)| < 3\%$ for all the datasets.
- A malicious intruder cannot establish bounds on the private data using KDE Resampling.
- Only one-way communication tractable.

KDE Resampling Algorithm

Numerical Results

Conclusions

Summary

- Proposed an algorithm (KDE Resampling) for data sanitization to share private data for distributed classification.
 - Asymptotically independent (Privacy).
 - Consistent (Accuracy).
- max $|P_{rand}(err) P_{ori}(err)| < 3\%$ for all the datasets.
- A malicious intruder cannot establish bounds on the private data using KDE Resampling.
- Only one-way communication tractable.

KDE Resampling Algorithm

Numerical Results

Conclusions ●○○

Summary

- Proposed an algorithm (KDE Resampling) for data sanitization to share private data for distributed classification.
 - Asymptotically independent (Privacy).
 - Consistent (Accuracy).
- max $|P_{rand}(err) P_{ori}(err)| < 3\%$ for all the datasets.
- A malicious intruder cannot establish bounds on the private data using KDE Resampling.
- Only one-way communication tractable.

KDE Resampling Algorithm

Numerical Results

Conclusions

Summary

- Proposed an algorithm (KDE Resampling) for data sanitization to share private data for distributed classification.
 - Asymptotically independent (Privacy).
 - Consistent (Accuracy).
- max $|P_{rand}(err) P_{ori}(err)| < 3\%$ for all the datasets.
- A malicious intruder cannot establish bounds on the private data using KDE Resampling.
- Only one-way communication tractable.

KDE Resampling Algorithm

Numerical Results

・ロ ・ ・ 一 ・ ・ 日 ・ ・ 日 ・

э

Conclusions

Further Work and Acknowledgments

Further Work

- Re-examine the issue of the privacy metric.
- \mathcal{DAPL} only quantifies the l_1 distances between the 2 distributions.
- Assumption that the data records is IID may not be realistic.
 - But our algorithm requires this assumption.
 - This shortcoming will be addressed in our future work.

KDE Resampling Algorithm

Numerical Results

Conclusions

Further Work and Acknowledgments

Further Work

- Re-examine the issue of the privacy metric.
- DAPL only quantifies the l₁ distances between the 2 distributions.
- Assumption that the data records is IID may not be realistic.
 - But our algorithm requires this assumption.
 - This shortcoming will be addressed in our future work.

KDE Resampling Algorithm

Numerical Results

Conclusions

Further Work and Acknowledgments

Acknowledgments

• I would like to thank the support of the Agency for Science, Technology and Research (A*STAR), Singapore.



• Homepage: http://web.mit.edu/vtan/www.

