# COMP 551: Applied Machine Learning

Project 4 Presentation :

Methods for Responsible Machine Learning
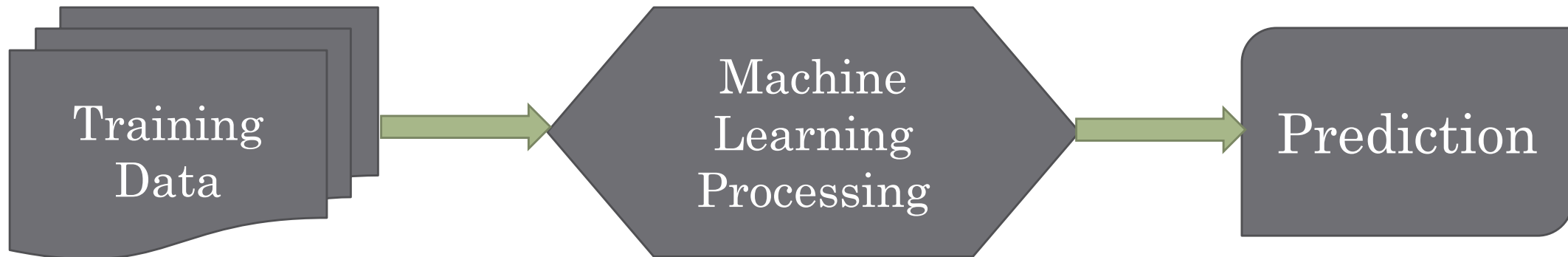
Topic: Privacy Preserving for Machine Learning in Banking
Applications

Hui Lee Ooi
Raylene MacDonald
Amèn Memmi

# Introduction: Motivation

- Banking industry, process huge amount of data, tackle problems at large scale.

- Automated processing and decision making.

- Security

Training Data → Machine Learning Processing → Prediction

# Introduction: Motivation

- Too little information: insufficient knowledge for processing the data.

- Too much information : privacy infringement

- Tradeoff?

- Maximize the amount of useful information while preventing sensitive information from being inferred.

|  | Gender | Age | Profession | Annual Earning |
|---|---|---|---|---|
| Ali | Male | 34 | Doctor | … |
| Bob | Male | 55 | Engineer | … |
| Cecilia | Female | 42 | Teacher | … |
| Dan | Male | 23 | Programmer | … |

# Introduction: Dataset

- Dataset on default of credit card in Taiwan.

- Dataset containing 30000 clients with 23 features from UCI Machine Learning Repository.

- Attribute information:
  - $X_1$: amount of given credit
  - $X_2$: gender
  - $X_3$: education
  - $X_4$: marital status
  - $X_5$: age
  - $X_6$-$X_{11}$: history
  - $X_{12}$-$X_{17}$: amount of bill statement
  - $X_{18}$-$X_{23}$: amount of previous payment

  Binary output variable: default on next bill payment (1 or 0)

# SVD (Singular Value Decomposition)

- SVD is a matrix factorization method generally used for dimensionality reduction; however, it can also be used for data distortion.

The singular value decomposition of the matrix $A$ is [14]

$$A = U \Sigma V^T,$$

where $U$ is an $n \times n$ orthonormal matrix, $\Sigma = \text{diag}[\sigma_1, \sigma_2, \ldots, \sigma_s]$ $(s = \min\{m, n\})$ is an $n \times m$ diagonal matrix whose nonnegative diagonal entries are in a descending order, and $V^T$ is an $m \times m$ orthonormal matrix. The number of nonzero diagonals of $\Sigma$ is equal to the rank of the matrix $A$.

[1]

- We choose the k largest diagonals of $\Sigma$, representing the highest variance among the attributes, to obtain $A_k = U_k \Sigma_k V_k$, an n x k matrix approximating the original n x m data matrix A.

# Impact of SVD on Classification Accuracy

- Using scikit-learn's Random Forest classifier on the original data, we obtained, through 5-fold cross-validation, an accuracy of 0.81. We then experimented with different values of k for SVD.

| Data distortion | Accuracy |
|---|---|
| Untouched data | 0.81 |
| SVD with k = 15 | 0.79 |
| SVD with k = 10 | 0.76 |
| SVD with k = 5 | 0.73 |

- Even for small values of k, the accuracy remains reasonable. This indicates the output variable is closely correlated to a small subset of the original 23 features.

# Data Perturbation with Random Noise

- We normalize the values of our data matrix A to the [0, 1] range

- We compute $A_u = N_u + A$, where the values of $N_u$ are uniformly distributed in the [0, 1] range

- We compute $A_n = N_n + A$, where the values of $N_n$ are normally distributed with parameters $\mu$ and $\sigma$

| Data distortion | Accuracy of Random Forest Classifier |
|---|---|
| Original, normalized data | 0.81 |
| With uniform noise from [0, 1] | 0.78 |
| With Gaussian noise, $\mu = 0$ and $\sigma = 0.1$ | 0.76 |
| With Gaussian noise, $\mu = 0$ and $\sigma = 0.5$ | 0.75 |
| With Gaussian noise, $\mu = 0.5$ and $\sigma = 0.5$ | 0.74 |

# How secure are these methods?

We use the Value Difference metric to measure how much the original data, A, has been altered to yield the distorted data matrix A' :

$$VD = \| A - A' \| \, / \, \| A \|$$

| Data Distortion | Value Difference |
| --- | --- |
| SVD with k = 15 | 0.001 |
| SVD with k = 10 | 0.049 |
| SVD with k = 5 | 0.170 |
| With uniform noise from [0, 1] | 2.203 |
| With Gaussian noise, μ = 0 and σ = 0.1 | 0.393 |
| With Gaussian noise, μ = 0.5 and σ = 0.5 | 2.625 |

A suggested [1] Value Difference of at least 0.15 is deemed sufficient for privacy-preserving methods; however, this depends on the domain.

# K-Anonymity

- Attributes are altered until each row is identical with at least k-1 other rows:
  - *Suppression* – can replace individual attributes with a *.
  - *Generalization* – replace individual attributes with a broader category according to a certain <u>hierarchy</u>. Example: (Age: 26 => Age: [20-30]; Occupation: Police officer => Governmental )

- K-Anonymity thus prevents definite database linkages. At worst, the data released narrows down an individual entry to a group of k individuals

- Unlike Output Perturbation models, K-Anonymity guarantees that the data released is accurate

# Implementation

- ARX anonymization tool: open source software used to anonymize data

- We defined hierarchy for each attribute

- ARX scans solution space and chooses optimal Suppression – Generalization combo to minimize risk.

# Input vs Output Data

# Anonymity gain – Identification risk reduction

**Before 3- anonymity:**



**After 3- anonymity:**

# Effect on ACC and AUC

| K-anonymity order | Original data | K=2 | K=5 | K=10 | k=20 | K=50 | K=100 | K=500 | K=1000 |
|---|---|---|---|---|---|---|---|---|---|
| Logistic Regression | 81.06 | 80.16 | 80.77 | 79.74 | 78.12 | 78.44 | 79.12 | 78.45 | 78.44 |
| Naïve Bayes | 80.93 | 80.43 | 81.2 | 79.56 | 79.56 | 79.76 | 79.5 | 79.26 | 78.81 |
| Random Forest | 82.16 | 81.13 | 80.63 | 80.63 | 80.53 | 79.46 | 80.2 | 78.26 | 78.33 |

Accuracy(ACC) vs K-anonymity order K for 3 classifiers

| K-anonymity order | Original data | K=2 | K=5 | K=10 | k=20 | K=50 | K=100 | K=500 | K=1000 |
|---|---|---|---|---|---|---|---|---|---|
| Logistic Regression | 61.14 | 60.55 | 60.75 | 60.12 | 59.88 | 59.74 | 59.48 | 60.15 | 60.01 |
| Naïve Bayes | 72.90 | 71.92 | 73.75 | 70.05 | 70.56 | 70.76 | 70.65 | 70.94 | 70.66 |
| Random Forest | 77.01 | 75.01 | 72.70 | 72.35 | 72.12 | 70.40 | 72.92 | 66.79 | 66.5 |

Area Under Curve (AUC) vs K-anonymity order K for 3 classifiers

# References

1. S. Xu, J. Zhang, D.Han, J.Wang, Data Distortion for Privacy Protection in a Terrorist Analysis System. University of Kentucky, 2005.

2. S. Han, W. Keong Ng, P. S. Yu, Privacy-Preserving Singular Value Decomposition. IEEE International Conference on Data Engineering, 2009

3. M. Naga Lakshmi1 & K Sandhya Rani, SVD based Data Transformation Methods for Privacy Preserving Clustering. International Journal of Computer Applications Volume 78 – No.3, September 2013.