

小口研究室 研究紹介 (2016年度)

(お茶の水女子大学理学部情報科学科)

完全準同型暗号を用いたAprioriアルゴリズムの並列分散計算による高速化手法の検討 (研究担当: 宇佐美 文梨)

Introduction | セキュアなビッグデータからの知識獲得

ビッグデータ活用の推進により、企業や病院などで蓄積されているビッグデータについても注目が集まっている。大量のデータを解析し新たな知識を得ることで、今までにない手法や着眼点を得るためには、年々膨れ上がるデータを捌き、要求により複雑化していく計算量の大きなアルゴリズムを扱えるマシンパワーが必要となる。

データマイニングを行う際、クラウドサービス等を提供する第三者の計算資源が一般的に用いられるが、高い機密性を有するデータを外部のストレージに預けることで、委託先に勤務する従業員等による流出等、新たなセキュリティ上の問題が生じる。

頻出パターンマイニングの例 | Aprioriアルゴリズム

POSレジなどに蓄積された購買情報から、お客さんがよく買う商品ボタンを知りたい!

アイテム (商品)	D	B	I	A
T1	○	○	○	○
T2	○	○	○	×
T3	○	×	×	×
T4	○	○	×	○
T5	×	○	○	×
T6	×	○	○	○
T7	○	○	×	×

①探索対象となる商品の集合(ボタン)を決める
初期値は ボタン長=1

②その商品ボタンの購入者数(サポート)を数える

③閾値と比べ、頻出ボタンを決定する

ボタン長	1	2	3
①	{D},{B},{I},{A}	{D,B},{B,I},{I,D}	{D,B,I}
②	5,6,4,3	4,4,2	2
③	{D},{B},{I}	{D,B},{B,I}	∅

最低サポート数(閾値)=4

該当なしなら終了

プライバシー保護データマイニング

高い機密性を有するデータから、適切なセキュリティ管理のもと知識を獲得するための技術の総称。ここでは①匿名化、②秘匿計算を紹介する。

①匿名化
真の値にノイズを混ぜ、統計的な手法で規則性を知る

②秘匿計算(準同型暗号方式)
暗号化したままAND・OR演算ができる完全準同型暗号によって暗号したデータを使用し委託計算

統計的な手法を用いにくいセキュアなデータからマイニング結果を得るには、準同型暗号を用いた秘匿計算が有効だが、一方で計算量が膨大になってしまう。

Method | 完全準同型暗号を用いた委託データマイニング

◆ Liuら(2015)^[1]の手法

◆ 本提案手法

目的: Aprioriの完全準同型暗号を用いた委託計算の高速化
方法: 広域分散環境を構築 & Boost.MPIを用いた分散並列計算を実現

大規模データ分散プラットフォームApache Sparkを用いた分散並列機械学習 (研究担当: 加藤 香澄)

研究背景
ライフログの取得・蓄積と利用の拡大
カメラやセンサーの発達
クラウドコンピューティングの普及

研究概要
各家庭から収集した動画データデータをSparkのストリーミング機能を用いて複数のワーカーに分散させ、動画データの機械学習処理を並列化

実験構成
MNISTをRDDに変換
Spark
Worker
Chainer
データ(MNIST)
マスターでPythonプログラムを実行

実験結果
Sparkによる並列機械学習処理部分の実験を実施
変化するパラメータ
読み込ませるデータのpartition数
ワーカーのノード数

Fairness Index
最も割り当てタスク数

まとめと今後の課題
今後の課題
分散環境のさらなる調査
MNIST以外のデータを用いて実験

場所と時間の制約条件に対応するSNSデータを用いたタイムリーなイベント提示システム (研究担当: 工藤 瑠璃子)

研究背景
2020年の東京オリンピック開催を受けインバウンドが増加しているが、観光情報の発信は不足
恒常的に人気のあるスポットはガイドブックなどから取得可能だが、今まさに開催されているイベントは情報を取得しにくい
地理的・時間的な制約がある旅行者などが必要とする「その時」「その場」で役立つ情報配信は少ない

提案システム
1) ツイートの抽出
2) ツイートのイベント分類
3) 配信準備

ツイートの分類評価
SVM (support vector machine)
ランダムフォレスト

配信準備
インバウンド対応のため多言語翻訳
情報の優先順位付け

今後の課題
ツイートの分類精度、テキスト情報整理の精度の向上
地名・日付・時間を含むツイートのみが対象になっているため解析対象のツイートの種類が少ない
この3点の情報が欠けている場合に補う手法を提案し、解析するツイートの種類を増やす