

東海大学大学院2011年度 修士論文

摂動化によってプライバシーを保護した
情報推薦方式

A Privacy-Preserving Recommendation Method
using Perturbation

指導教員 菊池 浩明 教授

東海大学大学院 工学研究科 情報理工学専攻

0BDRM039 望月 安菜

目次

第 1 章	序論	2
1.1	背景	2
1.2	目的	4
1.3	論文構成	5
第 2 章	情報推薦	6
2.1	情報推薦とは	6
2.2	情報推薦に利用するデータの種類の種類	7
2.2.1	内容ベースフィルタリング (Content based filtering)	7
2.2.2	協調フィルタリング (Collaborative Filtering)	7
2.2.3	Slope One	9
2.3	類似度の求め方	10
2.3.1	アイテム間 (Item-to-Item)	10
2.3.2	ユーザ間 (User-to-User)	10
2.4	考えられる特徴	11
2.4.1	ユーザ	11
2.4.2	アイテム	12
2.5	一般的なデータセットの特徴	14
2.6	従来研究	15
2.6.1	暗号化 (Encryption) によるアプローチ	15
2.6.2	摂動化 (Perturbation) によるアプローチ	15
第 3 章	協調フィルタリングを用いた情報推薦	17
3.1	摂動化協調フィルタリング	17
3.2	準備	18
3.2.1	関連研究	18
3.2.2	摂動化と再構築法	19
3.2.3	アイテムベース協調フィルタリング	21
3.3	提案方式	22
3.3.1	アイデア	22

3.3.2	提案方式 – 期待値による協調フィルタリング	22
3.3.3	数値例	24
3.3.4	考察	26
第 4 章	Slope One を用いた情報推薦	27
4.1	摂動化 Slope One	27
4.2	Slope One	28
4.2.1	概要	28
4.3	準備	30
4.3.1	関連研究	30
4.3.2	摂動化と再構築	30
4.4	提案方式	31
4.4.1	アイデア	31
4.4.2	提案方式 - 摂動化 Slope One	31
4.4.3	数値例	34
4.4.4	結果	35
4.4.5	考察	36
第 5 章	アイテム依存を考慮した摂動化	39
5.1	アイテム依存を考慮した摂動化	39
5.2	維持確率	40
5.2.1	数値例	41
5.3	提案方式	42
5.3.1	アイデア	42
5.3.2	提案方式 - アイテム依存維持確率	43
5.3.3	実験	45
5.3.4	実験結果	46
5.3.5	考察	46
第 6 章	評価実験	47
6.1	実験データ	47
6.2	維持確率	48
6.3	精度評価	53
6.3.1	平均絶対誤差 (Mean Absolute Error)	53
6.3.2	Slope One	53

第 7 章 結論と今後の課題	54
7.1 結論	54
7.2 課題	54
7.2.1 アイテム依存維持確率	54
参考文献	55
謝辞	60
付 録 A ベイズ推定による摂動化アルゴリズム	61
A.1 ベイズ推定による摂動化アルゴリズムの実装	61
A.1.1 摂動化の概要	61
A.1.2 再構築アルゴリズムの原理	62
A.1.3 実験結果	63
A.1.4 結論	66
付録	61

第1章 序論

1.1 背景

現在，我々は家から出ることなく様々な買い物を楽しむことができる．Amazon¹⁾に始まり，楽天²⁾など多くのインターネットショッピングサイトが我々の生活に定着している．いつでも，家にいながら簡単にボタンをクリックするだけで欲しい商品を注文することができ，次の日には手元に届く．そこで目にするのが，「おすすめ商品」や「セール情報」などの広告である．これらの情報は，ショッピングサイトだけにかかわらず，多くのウェブサイトで見かけることができる．いわゆる，情報推薦システムである．

この情報推薦は，カスタマーのユーザ登録時に入力した情報や，商品の閲覧履歴，検索ワードなどの「好み」を基にデータマイニングをすることによって実現されている．データマイニングとは，商品の閲覧履歴や購入履歴などの蓄積されるデータを解析し，データを見ただけでは分からないような「モノ」のパターンや相関関係を探し出す技術である．その代表的な手法の一つに協調フィルタリング (Collaborative Filtering) がある．協調フィルタリングは複数のユーザによって複数のアイテムが評価付けされているデータベースにおいて，他のユーザの値をもとに，評価されていないアイテムの評価値を予測する技術である [15, 23, 24, 25]．この協調フィルタリングより容易に計算が行える手法として *Slope One* がある．*Slope One* は，協調フィルタリングと同様に評価されていないアイテムの評価値を予測する方法であるが，計算が容易に行え，コストも小さく抑えることができる [7, 8]．本論文は協調フィルタリングと *Slope One* の双方を扱う．

一方で，クラウドコンピューティング技術が発達し Google³⁾を始めとする様々な企業が，新しいサービスを世の中に発信している．その結果，多くのユーザが様々なデータをインターネット上に保存している．クラウドコンピューティングとは，インターネットを経由して様々なソフトウェアやハードウェアなどのコンピュータ資源を利用することができるサービスである．

例えば，手元のコンピュータに文書作成ソフトがインストールされていなくても，イン

¹⁾アメリカの大手インターネットショッピングサイト．本から家電，衣料品まで様々な標品を取り扱っている．アメリカや日本だけでなく現在7カ国で利用されている．(2012年3月27日)

²⁾日本の大手インターネットショッピングサイト．幅広いジャンルをカバーしており，個人で出店することも可能．(2012年3月27日)

³⁾アメリカの大手検索サイト．検索エンジンだけにかかわらず，メールや，カレンダー，OSなど様々なサービスを基本的に無料で提供している．(2012年3月27日)

ターネットに接続すれば文書作成ソフトが使用することが可能になる。他にも、手元のコンピュータには少しの保存容量がない場合でも、オンラインストレージと呼ばれるクラウドサービスを利用することによって、インターネット上にある保存領域を使用することができる。有名なサービスとして電子ノートブックサービスを提供している Evernote⁴⁾や、オンラインストレージサービスを提供している Dropbox⁵⁾などが存在する。

⁴⁾電子ノートブックのクラウドサービス。ちょっとしたメモから、ウェブサイトのクリップ、PDF、画像など様々な情報を保存することができ、後で参照や検索が可能である。OCR 機能も備わっており、画像の中に書いてある文字も検索することができる。

⁵⁾オンラインストレージサービス。無料で最大 4GB の保存領域を利用することができる。

1.2 目的

現在の情報推薦は、自組織が管理する顧客のデータのみを活用して行われている。もし、単一の組織が管理する顧客の秘密情報だけではなく、クラウドコンピューティング技術によってインターネット上に保存されている情報や、他の組織の管理するユーザの情報を学習することができれば、よりの確に商品の推薦やサービスの提供などを行うことができるのではないかと考えられる。

しかし、これらの情報は組織にとって、外部に知られたくはない重要な秘密の情報であり、プライバシーを保護しなければならない対象でもある。インターネット上に保管されている個人情報、内容が外部に知られないように、暗号化などが行われ秘匿されている。このように秘匿する場合には、暗号が使われていることが多いが、暗号化は計算コストが多くかかる。そのため本研究では、摂動化方式を提案する。評価値に、乱数を加える方法や、維持確率を定め攪乱する方法である。

そこで、本論文では摂動化によってプライバシーを守り、精度の高い情報推薦を試みる..

1. 摂動化協調フィルタリング

プライバシーを保持したまま、摂動化したデータから再構築を行う新しい情報推薦方式である。摂動化により一定確率で評価値を `Randmize` することでプライバシーを保護する。また、ベイズ推定による再構築によって、アイテム間類似度をオリジナルデータへ近似させることができ、協調フィルタリングを使用した情報推薦の精度を向上させることができる方式である。

2. 摂動化 *Slope One*

情報推薦の主流の協調フィルタリングではなく、アイテム間類似度を評価値の差分で定めるアイテムベース推薦方式の *Slope One* を用いることで `Randmize` で生じる誤差を小さくすることができる方式である。

また、維持確率を変化させることで、プライバシー保護と推薦精度の向上を試みる。それぞれについて、推薦制度を評価する。

1.3 論文構成

本論文の構成は次の通りである。

第2章で情報推薦に関する基礎や、方式の種類、情報推薦を行うにあたって考慮すべき点などについて説明し、第3章で、摂動化によってプライバシーを保護し、協調フィルタリングによって情報推薦を行う提案方式について説明する。第4章では、第3章の協調フィルタリングに変え、Slope One を使用した情報推薦を提案する。そして、第5章で、第4章の精度を向上させるための維持確率を提示する。第3章、第4章、第5章の実験と結果を第7章で述べる。最後に第7章で本論文の結論を述べ、今後の課題を示す。

第2章

情報推薦

2.1 情報推薦とは

情報推薦とは、まだ自分が所持していない、または知らないであろう商品や情報を予測し、推薦することである。

インターネットを利用すれば分かるように、膨大な情報や商品が溢れている。これらの増え続けていく膨大な情報を隅から隅まで見ることは不可能であり、その中から自分にとって有益なものを見つけ出すことは非常に困難だといえる。そこで、利用されているシステムが推薦システム (Recommender System) である。

最も身近な例は、インターネットショッピングサイトの Amazon でよくみられる、おすすめ商品である。Amazon のページにアクセスをすると、“閲覧履歴からお勧め”(図 2.1) や“これにも注目”(図 2.2) などの項目に商品が表示される。これらは、閲覧履歴や以前購入した商品から推測された情報である。



図 2.1: Amazon による“閲覧履歴からお勧め”された商品

この第2章では、この情報推薦(システム)にはどのような種類や、方法があるのかを説明していく。



図 2.2: Amazon の推薦によって“これにも注目”に表示された商品

2.2 情報推薦に利用するデータの種類の種類

推薦システムで利用する方法は、大きく内容ベースフィルタリング (Content based Filtering) と、協調フィルタリング (Collaborative Filtering), *Slope One* の三種類に分けられる。この節では、この三つの手法について簡単に説明する。

2.2.1 内容ベースフィルタリング (Content based filtering)

内容ベースフィルタリングは、はじめにユーザがどのような情報を欲しがっているかを入力してもらい、その内容にあったフィルタリング¹⁾結果を提示する方式である。例えば、あるユーザが「1990年代」、「アクション映画」という情報を提示した場合、それにマッチした「Toy Story(1995)」や「Jumanji (1995)」、「GoldenEye (1995)」をフィルタリング結果として提示する。

この内容ベースフィルタリングの利点は、ユーザの要望を基に情報をフィルタリングするため、求めている情報を提示し易いといえる。逆に欠点は、ユーザがわざわざ情報を提示しないとフィルタリングが出来ないという点や、検索される範囲が限定されるため、別ジャンルの意外な推薦が出来ないという点がある。

2.2.2 協調フィルタリング (Collaborative Filtering)

協調フィルタリングとは内容ベースフィルタリングとは異なり、ユーザ間、またはアイテム²⁾間の類似度³⁾を利用してフィルタリングを行う方式である。この方法は、自分と嗜好が

¹⁾膨大な情報に対して、不要な情報を取り除く為、本論文ではこの操作を「フィルタリング」と呼ぶこととする。

²⁾ここで、アイテムとは団体が取り扱っている商品や情報などの事を指す。

³⁾類似度とは対象 A と対象 B がどのくらい似ているか (類似しているか) を示す指標である。一般的に数値が高ければ高いほど似ているとみなされる。

似ているユーザを類以度によって探し出し、その人がお勧めするアイテムをフィルタリング結果として返す事によって実現され、最も有名な手法として広く利用されている。つまり、自分と同じ嗜好を持った人が居たとして、その人が好きなアイテムは「おそらく」自分も好きであろうと予測するのである。

メモリベース法 (Memory-Based Methods)

メモリベース法は、ユーザが評価した値をそのまま利用して予測する方法である。

ユーザが評価した値の例を表 2.1 に示す。 u_j は j 番目のユーザを示し、 i_ℓ は ℓ 番目のアイテムを示す。例えば、ユーザ u_1 がアイテム i_6 に与えた評価は 5 となる。

この方法には、アイテム間の類以度に注目したアイテムタイプと、ユーザ間の類以度に注目したユーザタイプの二種類がある。表 2.1 を映画に対する評価だとする。 u_2 に新しい映画を推薦しようとした時に、まず u_2 と類似しているユーザを探す。最も、類似しているユーザが u_4 だったとすると、 u_2 がまだ観たことがなく、なおかつ u_4 が高い評価をつけている映画は i_7 になる。最終的に推薦システムは u_2 に i_7 を推薦することが出来る。

これらの詳細については、2.3 節で具体的な例を用いて紹介する。

表 2.1: メモリベース法で利用されるデータ例

	i_1	i_2	i_3	i_4	i_5	i_6	i_7
u_1	1	1	3	4	2	5	
u_2	3	2		5	5	1	
u_3		2	3	2			1
u_4	4			4		2	5

モデルベース法 (Model-Based Methods)

モデルベース法は、クラスタリングなどを行い、先にモデルを構築しておき、推薦対象者をそのクラスタに当てはめ、そのクラスタ内で推薦する方式である。

クラスタリングとは、ある集合 U を重複部分が無いように分類することで、分類されたそれぞれのグループをクラスタと呼ぶ。クラスタリングの手法に関しては最短距離法や最長距離法による階層型クラスタリングや、 k -means や SOM などの非階層型クラスタリングがある。

例えば、ある集合 U を 3 つのクラスタ C_1, C_2, C_3 に分類すると、 $U = C_1 \cup C_2 \cup C_3$ となり、 $\phi = C_1 \cap C_2 \cap C_3$ となる。ユーザ u_1 と各クラスタとの距離 (類以度) を求め、最も近いクラスタ無いで評価が高いアイテムを推薦する方法である。

2.2.3 Slope One

Slope One は、D. Leniel ら [7] によって提案された方式である。シンプルなアルゴリズムであるが、高い性能で商用にも採用されている。Slope One とは、アイテム間の相関に傾き 1 の一次式、 $f(x) = x + b$ を用いているところからその名が付いている。特異値分解などの既存の推薦方式と比較してアイテム間平均差分に基づいて推薦を行うので実装も容易で処理性能も高い。また、協調フィルタリングと比較しても、計算量の多い類似度計算を行わないため、コストが少ない。詳細は、第 4 章に示す。

2.3 類以度の求め方

協調フィルタリングのメモリベース法において、類以度を計算する対象として、ユーザ間とアイテム間の二つがある。

2.3.1 アイテム間 (Item-to-Item)

アイテム間類以度を利用するタイプは、アイテムに与えられた評価値をアイテムのベクトルとみなし、二つのアイテムベクトル間の類以度を計算する方法である。

これは、好きなアイテムと似ているアイテムであれば、そちらも好きであろうというというアイデアである。Amazon でよくみられる“これにも注目”などの、既に関連したアイテムと似たような（類以度の高い）アイテムが表示されるのは、この方式ではないかと考える。

ここで $s_{j,\ell}$ はアイテム i_j とアイテム i_ℓ の類以度を表す。つまり、 i_1 と i_2 の類以度 $s_{1,2}$ は 0.5 となる。

表 2.2: アイテム間の類以度

	i_1	i_2	i_3
u_1	1	1	5
u_2	1	1	4
u_3	2	2	2
u_4	4	4	1
u_5	5	4	1
s_{u_1, u_k}	-	0.5	0.02

2.3.2 ユーザ間 (User-to-User)

ユーザ間類以度を利用するタイプは、ユーザに与えられた評価値をユーザのベクトルとみなし、二つのユーザベクトル間の類以度を計算する方法である。

これは、似たような嗜好を持ったユーザが好きなモノであれば、そのユーザも好きであろうというというアイデアである。Amazon で言うと、“この商品を買った人はこんな商品も買っています”に相当するのではないかと考える。

アイテム間と同様に、変形ユークリッド距離を利用して、 u_1 とその他のユーザの類以度を計算した例を表 2.3 に示す。

ここで $s_{j,\ell}$ はユーザ u_j とユーザ u_ℓ の類以度を表す。つまり、 u_1 と u_2 の類以度 $s_{1,2}$ は 0.5 となる。

表 2.3: ユーザ間の類似度

	i_1	i_2	i_3	$s_{1,l}$
u_1	1	1	5	-
u_2	1	1	4	0.5
u_3	2	2	2	0.08
u_4	4	4	1	0.03
u_5	5	4	1	0.02

2.4 考えられる特徴

2.4.1 ユーザ

情報推薦をするにあたって考慮しなくてはならないポイントの一つに、ユーザ特徴がある。ここで言う、ユーザの特徴とは評価値の付け方や購買率を言う。以下に代表的なユーザの特徴例を示す。

1. 辛口ユーザ

どんなアイテムも辛口に評価するため、評価値の平均が低い。

2. 甘口ユーザ

どんなアイテムも甘口に評価するため、評価値の平均が高い。

3. 評価が極端

お気に入りのアイテムは極端に高く、好きではないアイテムは極端に低いため、評価値の分散が大きい。

4. 購買率が高い

多くのアイテムを購入しているユーザで、そのユーザに関するスパース率が低い。

スパース率とは、全アイテム中いくつ評価されているかという、データベースの密度を表す。データベースが疎である则この値は高く、密であると低い値を示す。

5. 購買率が低い

アイテムの購入数が少ないため、そのユーザに関するスパース率が高い。

表 2.4 に評価値にみられるユーザの特徴例を示す。ここで、 σ はそのユーザが評価値の標準偏差を表す。

表 2.4: 評価値にみられるユーザーの特徴例

	i_1	i_2	i_3	i_4	i_5	i_6	i_7	mean	σ	sparsity
辛口ユーザー	1	1	2	3	2	3		2	0.81	85%
甘口ユーザー	4	3		5	5	3		4	0.89	71%
評価が極端		1	5	1	4	5	1	2.83	1.86	85%
購買率が高い	4	2	2	4	3	2	5	3.14	1.12	100%
購買率が低い	4				3			3.5	0.5	28%

2.4.2 アイテム

ユーザーの特徴と同時にアイテムの特徴もある。アイテムの特徴は、ユーザーの特徴と同様に評価値の付け方、購買率、そしてユーザーからの人気を言う。以下に代表的なアイテムの特徴例を示す。

1. 人気アイテム

どのようなユーザーも高い評価を付けるため、評価値の平均が高い。

2. 不人気ユーザー

どのようなユーザーも低い評価を付けるため、評価値の平均が低い。

3. 評価が極端

ユーザーによって好き嫌いが分かれるアイテムのため、評価値の分散が大きい。

4. 購買率が高い

多くのユーザーが購入しているので、スパース率が低い。

スパース率とは、全アイテム中いくつ評価されているかという、データベースの密度を表す。データベースが疎であるとこの値は高く、密であると低い値を示す。

5. 購買率が低い

購入しているユーザーが少なく、スパース率が高い。

表 2.4 に評価値にみられるユーザーの特徴例を示す。ここで、 σ はそのユーザーが評価値の標準偏差を表す。

表 2.5: 評価値にみられるユーザの特徴例

	人気アイテム	不人気アイテム	評価が極端	購買率が高い	購買率が低い
i_1	5	1	1	3	
i_2	5	3		5	5
i_3	5	1	5	1	
i_4	4	2	2	4	
i_5	5		5	1	3

2.5 一般的なデータセットの特徴

Movie Lens Dataset[19] や Netflix[20] のデータセットにみられる特徴は、非常にスパース（疎）であるということである [12]。おそらく、Amazon で扱われている商品は数えきれないほど存在する一方で、各ユーザが購入している商品はほんの 1% にも満たない。実際に、自分のケースに当てはめてみても、Amazon で販売されている商品を買った個数は 20 個にも満たない。

これらのスパースなデータベースに対応するため、Breese らによって標準投票（Default Voting）などが提案されている [14]。この手法は、デフォルト値を決定しておき、欠損値を補完する手法で、最も中立的な値である投票値の平均が利用される。

また、高島らによって Prediction Voting[26] が提案されている。Default Voting による欠損値の補完を、平均値ではなく相関係数法による予測結果にするというアイデアである。

しかし、評価方法のばらつきやノイズが存在する為、そのまま適応すると予測性能が低下してしまうという問題がある。そこで、平山らは Prediction Voting において、いかに良いフィードバックを行うかを提案し、その精度向上を図っている [15]。

2.6 従来研究

多くの情報推薦に関する研究がなされてきた．ここでは，プライバシーを保護した従来研究を紹介する．

プライバシー保護には，大きく分けて暗号による保護と摂動化による保護の二つのアプローチが研究されている．

暗号による保護は，準同型暗号を使い，値を暗号化したまま計算できるため精度に優れているが，暗号文生成や暗号計算になど膨大な計算コストがかかる．また，暗号文のサイズが大きいため通信コストも大きくなってしまう．

摂動化による保護は，値に乱数などを加え計算をするため，暗号を使ったアプローチに比べ，正確な計算をすることができないが，計算コストが高い暗号計算をする必要がなく，コストが安いメリットを持っている．

本研究では，摂動化によるアプローチを行っている．

2.6.1 暗号化 (Encryption) によるアプローチ

Canny は特異値分解 (Singular Value Decomposition) を準同型暗号を用いて，情報を秘匿したまま協調フィルタリングを行う手法を提案している．木澤らは，コサイン類似度と準同型暗号を用いた秘匿協調フィルタリング [16] や，秘匿性集合プロトコルを利用してプライバシーを保護した協調フィルタリングを提案している [17]．これらは水平分割された P2P モデルを想定している手法である．

また，多田らは秘匿関数計算を利用し水平分割，アイテムタイプの秘匿協調フィルタリングを提案している [18]．

佐久間らは P2P 環境で非同期に値を秘匿したまま平均を計算プロトコル [55] や，値を秘匿したまま k -means クラスタリングを行う手法を提案している [56]．

Vaidya らは，水平分割で Naïve Bayes を利用したクラスタリングを秘匿したまま行う手法 [28] や，秘匿アソシエーションルールマイニングを提案している [29]．

2.6.2 摂動化 (Perturbation) によるアプローチ

Agrawal らは，Randomized Response を利用して決定木学習を行う方式を提案している [1]．この手法は，分割されたデータに対して独立に摂動化によるランダム化を適用した後，これらをマージしてアルゴリズムを適用するため，任意の分割モデルに対して利用可能である．彼らが用いた決定木では，数値属性を扱う．葉ノード以外のノードはある属性 A と閾値 θ および二本の枝を持つ．各枝は，属性 A の属性値 α について，それぞれ条件式

$a \leq \theta, a > \theta$ を保持する．葉ノードは分類を決定するクラス属性を持つ．決定木は再構築アルゴリズムによって構築される．

また，Polat らによって加法摂動化による協調フィルタリング方式 [3] が提案され，Huang らがそのアタック方式と，改良方式を提案している [?] ．

第3章

協調フィルタリングを用いた情報推薦

3.1 摂動化協調フィルタリング

情報推薦の主流は、複数のユーザによって複数のアイテムが評価付けされているデータベースにおいて、他のユーザの値を基に評価されていないアイテムの評価値を予測する協調フィルタリング (Collaborative Filtering) である。

しかし、これらの情報推薦には、不正なサービス事業者によるプライバシー漏洩の課題がある。そこで、プライバシーを守るために、準同型性を満たした公開鍵暗号を使った個人情報を秘匿する研究がある [4]。しかし、暗号は、プライバシー保護は出来るが、大きな計算コストがかかる。そこで、本研究では、個人のプライバシーを保護しながら、暗号化をせずにユーザに応じた情報推薦を行うことを目的とする。

提案方式は、摂動化と呼ばれるランダムイズアルゴリズムと協調フィルタリングからなる。既存研究として、Agrawal ら [1] の研究が挙げられる。彼らは、データマイニング時に人工的に加えたノイズの影響をベイズ推定によって取り除き、決定木学習を実現している。本研究では、彼ら同様、摂動化したデータを用いて、ベイズ推定によってノイズを除去するが、決定木学習の代わりに協調フィルタリングを適用する点に特徴がある。情報提供者の持つ情報にノイズを加える。このデータの解析を行った後に、ノイズ除去の処理を施し、解析結果を同定する。同定する過程を再構築という。

摂動化を使用したこれら手法は、暗号化と比較し計算コストが小さい。また暗号化に比べ実装が容易である。しかし、暗号化がほぼ厳密に正しい結果を得ることに対して、摂動化の学習結果は近似解である。安全性の保証観点からは、暗号化の方が厳密である。

本研究では、摂動化を実行するプログラムを実装した。主な成果としては、以下の通りである。(1) 既存研究によって提案されていたいくつかの有用な摂動化の性質を明らかにした。(2) 新たな方式として摂動化協調フィルタリング法を提案した。(3) 摂動化協調フィルタリング法を使用し、公開データベースを使用し実験を行い、精度や性能を評価した。

3.2 準備

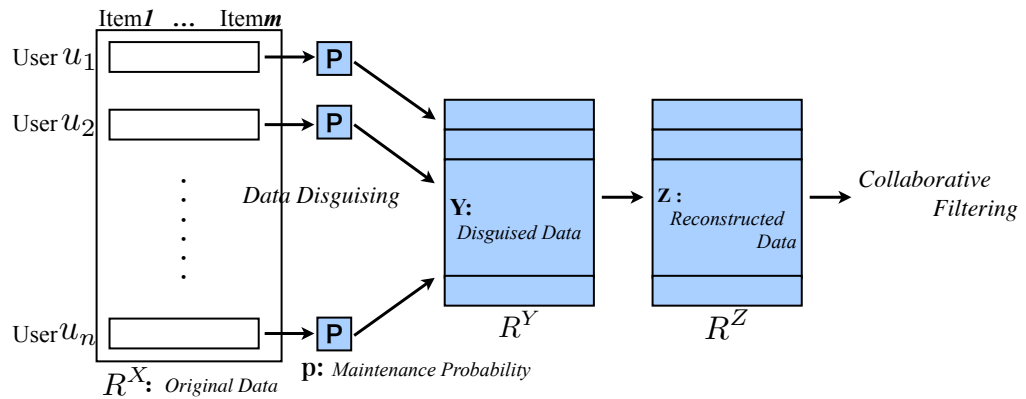


図 3.1: プライバシー保護した協調フィルタリングの全体構造

3.2.1 関連研究

H. Polat ら [3] は、加法摂動化による協調フィルタリング方式を提案している。彼らの研究では、オリジナルデータ X に一様分布の乱数 R を加えた $Y = X + R$ について、平均値 $\sum_i Y_i = \sum_i X_i + \sum_i R_i \approx \sum_i X_i$ であることを仮定したナイーブな推薦方式である。従って、 Y を、主成分分析 (PCA) することで加えた乱数ノイズを取り除くことが出来ることが指摘されており [2]、その安全性は低い。

そこで、本研究では、単純な PCA による解析が困難なランダムイズドレスポンス方式を用いて、摂動化を行う。安全性は向上するが、[?] の様な単純な協調フィルタリングでは精度が期待できない。以上の関係を表 3.1 に整理する。

表 3.1: プライバシー保護協調フィルタリング方式の比較

	H.Polat and W.Du[?]	提案方式
ランダムイズ 安全性 協調フィルタリング	加法摂動化 × (PCA) 容易	Randomized Response (目標)

3.2.2 摂動化と再構築法

再構築問題 (Reconstruction Problem) とは、摂動化された $Y_1 = X_1 + R_1$ 、確率変数 Y から、真の値 X の確率分布を見積もる問題である。R. Agrawal and R. Srikant[1] によって最初に発表された摂動化アルゴリズムである。秘匿したい情報に意図的にランダムノイズを乗せて、格納されたデータのプライバシーを保護する。

例えば、年齢が $x = 20$ 代であるという個人情報をそのまま渡す代わりに、一様分布（またはガウス分布）の乱数 r を加え、 $y = x + r = 30$ の様に歪んだ値 y を登録する。30 という属性値を持った顧客がいても、本当に 30 代なのか乱数で 40 代から歪まされたのか、第三者には区別がつかない。

暗号化による方法と異なり、時間のかかる暗号化はなく、計算も各パーティで独立に計算できる。通信効率も計算効率も高い。大規模なデータベースにおいても適用可能である。

摂動化

簡単な数値例を用いて再構築アルゴリズムの原理を示す。確率変数 A が表 3.2 の分布に従って与えられているとする。

表 3.2: 真の確率分布 $P(A)$

a	0	1	2	3
$P(A = a)$	0.1	0.3	0.1	0.5

ここで、 A の分布を秘匿する為に、表 A.2 の条件付確率 $P(B|A)$ に従って、 A の値を変化（摂動化）させた結果を B とおく。

ここで、維持確率 $p = 0.4$ は、 A を変化させない確率の大きさであり、変化させるときは一律な確率で分布させることにする。こうして摂動化した結果を表 A.3 で示す。オリジナルの分布では $A = 3$ が最頻度で生じていたのに対して、値の差が小さくなりどの値も同じくらい確からしい。

表 3.3: 条件付確率 $P(B|A)$, 維持確率 $p = 0.4$

$B \setminus A$	0	1	2	3
0	0.4	0.2	0.2	0.2
1	0.2	0.4	0.2	0.2
2	0.2	0.2	0.4	0.2
3	0.2	0.2	0.2	0.4

表 3.4: 摂動化した確率分布 $P(B)$

b	0	1	2	3
$P(B = b)$	0.22	0.26	0.22	0.3

再構築アルゴリズム

再構築アルゴリズムは, この $P(B|A)$ と摂動化後の確率分布 $P(B)$ だけを与えて, オリジナルの分布 $P(A)$ を近似することを目的とする. 初期値を $P^0(A) = P(B)$ で与える. 事後確率の i 番目の近似値は,

$$\begin{aligned} P^i(A|B) &= \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P^{i-1}(A)}{\sum_{a \in A} P(B|A=a)P^{i-1}(A=a)} \\ &= \frac{P(B|A)P(A)}{P(B|A=0)P^{i-1}(A=0) + \dots + P(B|A=3)P^{i-1}(A=3)} \end{aligned}$$

と近似され, この値を用いて A の事後確率の第一近似値は

$$P^1(A) = \sum_{b \in B} P^0(A|B=b)P(B=b)$$

で与えられる. こうして, 逐次的に近似を繰り返し, $P^{i+1}(A) = P^i(A)$ と収束した分布を再構築された P^* とする. この数値例の場合の第二近似値までの結果を表 A.4 で示す. 徐々に真の分布へ近づいていることが分かる.

表 3.5: 再構築された確率分布の第一近似 $P^1(A)$ と第二近似 $P^2(A)$

a	0	1	2	3
$P^1(A = a)$	0.22	0.26	0.22	0.31
$P^2(A = a)$	0.21	0.26	0.21	0.33

3.2.3 アイテムベース協調フィルタリング

協調フィルタリングは、ユーザ間あるいはアイテム間の類似度に基づき、未知のアイテムに対する評価値を予測するアルゴリズムである。ユーザ間の類似度から評価値を予測するユーザベース方式と、アイテム間の類似度を計算し、評価値を予測するアイテムベース方式がある。ここでは、アイテム間類似度を使用する。

ユーザ u のアイテム i についての評価値を $r_{u,i} \in V$ と表す。ユーザ数 n 、アイテム数 m の評価値行列を R とする。例えば、表 4.2 は、 $n = 4$ 、 $m = 5$ の例である。ここで、* の評価値を

$$r_{u,i} = \frac{\sum_j^m s_{j,i} r_{u,j}}{\sum_j |s_{j,i}|} \quad (3.1)$$

で予測する (アイテム間の平均に差がないことを仮定して、正規化を省略している)。ここでは、 $s_{i,j}$ は、アイテム i と j 間の類似度であり、本稿ではコサイン尺度により

$$s_{j,i} = \frac{\sum_{k=1}^n r_{k,i} r_{k,j}}{\sqrt{r_{1,i}^2 + \dots + r_{n,i}^2} \sqrt{r_{1,j}^2 + \dots + r_{n,j}^2}} \quad (3.2)$$

と定める。

表 3.6: Original Data (R^X)

	i_1	i_2	i_3	i_4	i_5
u_1	2	2	3	1	
u_2	1	3	2		3
u_3	2		3	3	2
u_4	3	2	*	2	2

表 3.7: Disguised Data (R^Y)

	i_1	i_2	i_3	i_4	i_5
u_1	2	3	1	1	
u_2	1	1	2		1
u_3	1		3	3	
u_4	3	2	*	2	3

3.3 提案方式

3.3.1 アイデア

評価値行列 R^X を真のデータ X とみなし、維持確率 p についてランダムイズドレスポンスによって摂動化を行い、偽データ行列 R^Y を作成する。この偽データ R^Y と p についてベイズ推定を行い、真のデータ X の再構築を行う。分布 Y から真の分布 X は予測できても、摂動化された行列 R^Y から、真の行列 R^X を求めることはできない。そこで、再構築の過程で求められた条件付き確率 $P(X|Y)$ を用いて、 R^Y からアイテム間類似度の期待値を求めて、推薦アルゴリズムに適用する。

3.3.2 提案方式 – 期待値による協調フィルタリング

提案の全体構造を図 3.1 に示す。各ユーザは、自分の評価値を決められた確率で摂動化(ランダムイズ)して集計サーバに送り、集めた R^Y を公開する。

真の評価値 $r_{u,i}$ と維持確率 p について、摂動化 y を

$$y = P(y) = \begin{cases} r_{u,i} & w/p = p \\ v \in V - \{r_{u,j}\} & \text{otherwise} \end{cases} \quad (3.3)$$

と定める。例えば、表 4.2 の行列 R^X を維持確率 $p = 0.4$ で摂動化した行列を表 4.3 の R^Y とする。

次に、 R^Y と p から定まる条件付き確率 $P(Y|X)$ に真の評価値 X についての予測値 Z と条件付き確率 $P^*(X|Y)$ を求める。

表 3.8: 条件付確率 $P(Y|X)$, 維持確率 $p = 0.4$

$Y \setminus X$	0	1	2	3
0	0.37	0.18	0.23	0.22
1	0.19	0.36	0.23	0.22
2	0.18	0.17	0.44	0.21
3	0.18	0.17	0.22	0.43

最後に、次のアルゴリズム CF-E により、任意のアイテムに対する評価値を予測する。

Algorithm 期待値を用いた協調フィルタリング (CF-E)Input: 摂動化評価値行列 $R^Y, P(X|Y)$ Output: 推薦値 $r_{u,i}^E$

Step 1 異なるアイテム間の2つの摂動化評価値 Y_1, Y_2 が与えられた時, それらの積を取る確率変数を $W(= Y_1 \cdot Y_2)$ とする. W の確率分布は

$$P(W|Y_1, Y_2) = \sum_{W=\alpha\beta} P(X = \alpha|Y_1) \cdot P(X = \beta|Y_2)$$

で求められる.

Step 2 積 W の期待値を求める. すなわち

$$E[W|Y_1, Y_2] = \sum_{\gamma \in V_2} P(W = \gamma|Y_1, Y_2)$$

ここで, V_2 は2つの V の要素からなる集合とする. $V = \{1, \dots, v\}$ の時, $V_2 = \{1, \dots, v^2\}$.

Step 3 摂動化行列 R^Y が与えられた条件の下で, アイテム i と j 間の類似度 $s_{i,j}^E$ の期待値を式 (3.2) のコサイン尺度で

$$s_{i,j}^E = E[S_{i,j}|R^Y] = \frac{E[\sum_u r_{u,i}^X \cdot r_{u,j}^X | R^Y]}{E[\sqrt{\sum_u (r_{u,i}^X)^2} \sqrt{\sum_u (r_{u,j}^X)^2}]} = \frac{\sum_u E[W|Y_1 = r_{u,i}^Y, Y_2 = r_{u,j}^Y]}{\sqrt{\sum_u (r_{u,i}^Y)^2} \sqrt{\sum_u (r_{u,j}^Y)^2}}$$

により求める. ここで, 分母は R^X におけるノルムを R^Y で近似している. 分子は step 2 の期待値で与えられる.

Step 4 式 (3.1) により期待類似度 s^E によるユーザ u のアイテム i の予測値は,

$$r_{u,i}^E = \frac{\sum_j S_{i,j}^E \cdot r_{u,j}^X}{\sum_j S_{i,j}^E}$$

で与えられる.

表 3.9: 各評価値の積 W の期待値 $E[W|Y_1, Y_2]$

$Y_2 \backslash Y_1$	0	1	2	3	sum
0	1.69	1.92	2.18	2.47	8.26
1	1.92	2.19	2.49	2.81	9.41
2	2.18	2.49	2.82	3.19	10.68
3	2.47	2.81	3.19	3.16	11.63

表 3.10: 類似度の期待値 $E[S_{i,j}|R^Y]$

	i_1	i_2	i_3	i_4	i_5
i_1	—	0.60	0.46	0.66	0.54
i_2	0.60	—	0.46	0.65	0.56
i_3	0.46	0.46	—	0.50	0.44
i_4	0.66	0.65	0.50	—	0.62
i_5	0.54	0.56	0.44	0.62	—

3.3.3 数値例

本実験では、表 4.2 のオリジナルデータ R^X を評価行列に、維持確率 $p = 0.4$ で摂動化した偽データ R^Y を表 4.3 とする、ユーザ数 $n = 4$ 、アイテム数 $m = 5$ 、評価値 $V = \{1, 2, 3\}$ とする。また、未評価のものを $r_{i,u} = 0$ とする。

各評価値の分布による再構築を行った結果を図 5.7 に示す。真の R^X の分布と $p = 0.4$ の維持確率により摂動化を行った偽データの分布 Y を示す。ベイズ推定を利用し、50 回再構築を行った結果 Z を示す。50 回行った再構築によって十分な精度で近似出来ることが分かる。再構築回数を増やしていくことで、更に近似する。

アルゴリズム CF-E により推薦を行う。摂動化された R^Y の 2 つの評価値が $Y_1 = 2, Y_2 = 3$ と観測された時、その積 W の確率分布を図 3.2 に示す。 $W = 0$ (未評価) が最大、 $W = 6 = 2 \cdot 3$ がその次に高い確率を持つ。

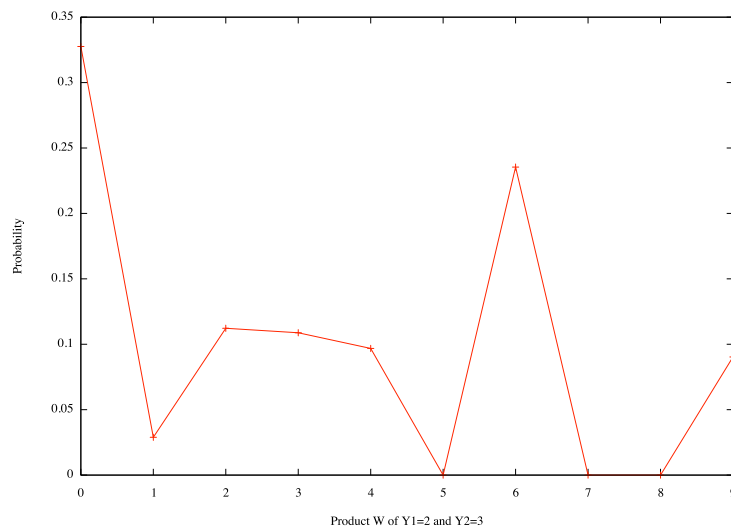


図 3.2: $Y_1 = 2, Y_2 = 3$ が与えられた時の積 $W (= Y_1 \cdot Y_2)$ の評価値の分布 $P(W|Y_1 = 2, Y_2 = 3)$

これらを平均すると、 $E[W|Y_1 = 2, Y_2 = 3] = 3.19$ であった。単純な R^Y の積 $2 \cdot 3 = 6$ より

も小さな値で見積もられる．

以上を全ての $Y_1 = 0, 1, 2, 3$, $Y_2 = 0, 1, 2, 3$ の組み合わせについて求めた結果を表 3.9 に示す．Step 3 により，求めたアイテム間の類似度の期待値 $s^E = E[S|Y_1, Y_2]$ を表 3.10 に示す．

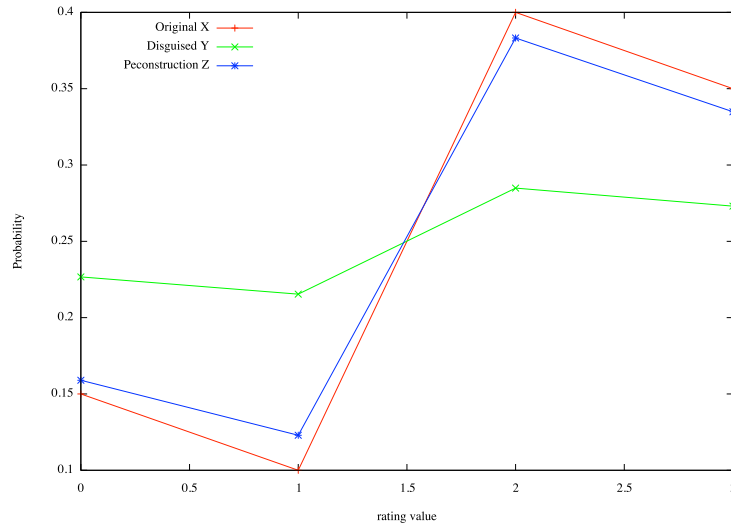


図 3.3: 再構築された評価値の分布

オリジナルデータ R^X ，偽データ R^Y ，再構築データについて，それぞれの求めた類似度の分布 $P(W|Y_1 = 2, Y_2 = 3)$ を図 3.4 に示す．オリジナルデータ X を基準としてみると，偽データ Y がオリジナルデータに近似している事を示している．

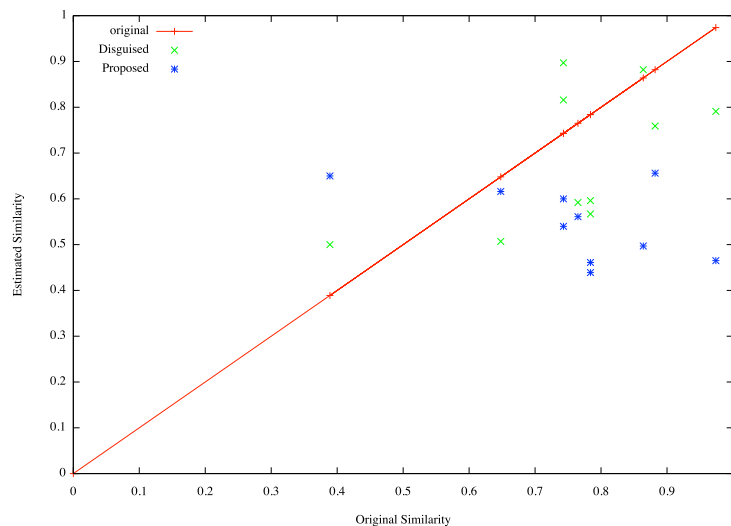


図 3.4: 予測類似度の分布 (s^X, s^Y, s^E)

オリジナルデータ X を協調フィルタリングした予測値を基準とし，摂動化偽データ Y ，再構築データ Z にそれぞれ協調フィルタリングを適用して推薦した値の分布を図 3.5 に示す．

提案手法である再構築を行ったデータがオリジナルデータの推薦結果に近似していることが分かる。

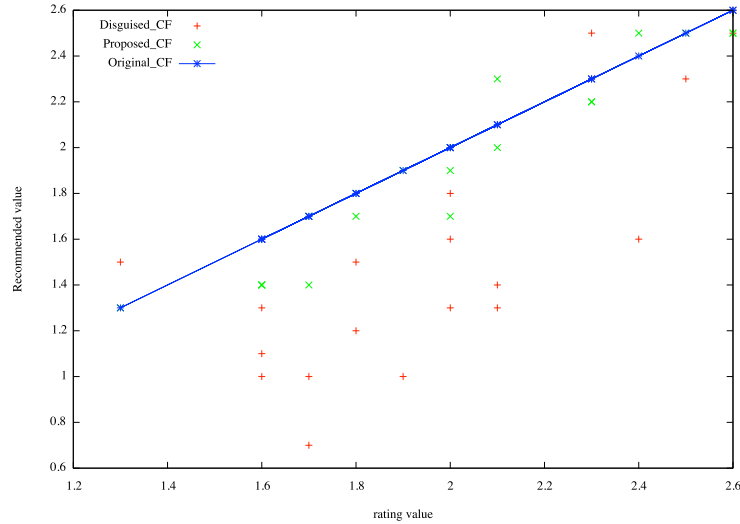


図 3.5: 協調フィルタリングによる推薦値の分布 (r^X, r^Y, r^E)

オリジナルデータ, 偽データ, 期待値を利用したデータで協調フィルタリングを行い予測した値の MAE (Mean Absolute Error) と標準偏差によって比較を行ったものを表 4.4 に示す. $MAE^E = \sum_{u,i} |r_{u,i}^X - r_{u,i}^E|$ と定める. 摂動化したデータにナイーブに CF を適用したのものは, $MAE^Y = \sum_{u,i} |r_{u,i}^X - r_{u,i}^Y|$ である. 偽データより期待値を利用した再構築データの方が誤差が少ない.

表 3.11: Mean Absolute Error

	MAE	標準偏差
Original	0.968	1.171
Disguised	1.033	1.204
Proposed	1.009	1.228

3.3.4 考察

図 3.5 より協調フィルタリングを行った際, 再構築データはオリジナルデータに近付いていることが分かる. また, 表 4.4 より, 偽データより再構築データの方が誤差が少ないが, その差は有意ではない. オリジナルの推薦値の誤差が大きいことから, 用いた評価行列が人工的で歪んでいた可能性がある. 提案方式では, 未評価値の扱いが十分ではなく, それが誤差の要因のひとつと考えられる.

第4章

Slope One を用いた情報推薦

4.1 摂動化 Slope One

情報推薦の主流は、複数のユーザによって複数のアイテムが評価付けされているデータベースにおいて、他のユーザの値を基に評価付けされていないアイテムの評価値を予測する協調フィルタリングである。しかし、協調フィルタリングは複雑な類似度の計算を行わなければならないため、評価値の差分を類似度とし、アイテムベース推薦方式の摂動化 Slope One によって推薦値を求める方式を提案する。

4.2 Slope One

D. Leniel and A. Maclachlan [7] によって提案された Slope One はアイテムベースの情報推薦アルゴリズムである。シンプルなアルゴリズムと高い性能で商用にも採用されている。Slope One とは、アイテム間の相関に傾き 1 の一次式、 $f(x) = x + b$ を用いているところからその名が付いている。特異値分解などの既存の推薦方式と比較して、アイテム間平均差分に基づいて推薦を行うので実装も容易で処理性能も高い。

4.2.1 概要

簡単な数値例を用いて Slope One アルゴリズムを解説する。ここで、評価値の定義域は 1 から 5 の離散値とし、“0” は欠損値を表す。

	I_1	I_2	I_3	I_4	I_5
U_1	-	1	4	2	1
U_2	-	-	-	-	2
U_3	-	-	1	1	-
U_4	5	1	-	-	3
U_5	2	-	2	3	5

この関係を、評価値行列

$$R_{5 \times 6} = (r_{i,j}) = \begin{pmatrix} 0 & 1 & 4 & 2 & 1 \\ 0 & 0 & 0 & 0 & 2 \\ 0 & 0 & 1 & 1 & 0 \\ 5 & 1 & 0 & 0 & 3 \\ 2 & 0 & 2 & 3 & 5 \end{pmatrix}$$

で表す。Slope One はシンプルに、差分の平均値で評価値を与える。すなわち、アイテム i_1 の評価値は、アイテム i_2 とその間の差分の平均 δ_{i_1, i_2} から、 $r_{i_1} = \delta_{i_1, i_2} + r_{i_2}$ と定義される。

$$* = \frac{\left(\frac{(4-2)+(5-2)}{2} + 1\right) + \left(\frac{(4-4)+(5-4)}{2} + 4\right)}{2} = 4.0$$
 により評価が与えられる。平均差分は、その両方のアイテムとも評価を与えているユーザについて求める。 i_2 の i_3 による類似度は i_1 によるものよりも高いので、評価値はより大きく影響を受ける。

上の例では、どちらのアイテムも同じ数のユーザによって評価されているので、単純に 2 で割って平均を取っているが、欠損値がある場合はこの限りではない。そこで、重みを考える。アイテム a と b の平均差分 $\delta_{a,b}$ を、

$$\overline{\delta_{a,b}} = \frac{\Delta_{a,b}}{\phi_{a,b}} = \frac{\sum_i \delta_{i,a,b}}{\phi_{a,b}} = \frac{\sum_i (r_{i,a} - r_{i,b})}{\phi_{a,b}} \quad (4.1)$$

で与える．平均差分行列 **average difference matrix** は，

$$\overline{\Delta}_{5 \times 5} = (\overline{\delta_{i,j}}) = \begin{pmatrix} 0 & 4 & 0 & -1 & -0.5 \\ -4 & 0 & -3 & -1 & -1 \\ 0 & 3 & 0 & 0.33 & 0 \\ 1 & 1 & -0.33 & 0 & -0.5 \\ 0.5 & 1 & 0 & 0.5 & 0 \end{pmatrix}$$

と定義する．ここで共生起数 $\phi_{a,b}$ は両方のアイテムを評価しているユーザの数である．相対共生起行列 **relative occurrence matrix**:

$$\Phi_{5 \times 5} = (\phi_{i,j}) = \begin{pmatrix} 2 & 1 & 1 & 1 & 2 \\ 1 & 2 & 1 & 1 & 2 \\ 1 & 1 & 3 & 3 & 2 \\ 1 & 1 & 3 & 3 & 2 \\ 2 & 2 & 2 & 2 & 4 \end{pmatrix}$$

で与える．

この時，ユーザ u のアイテム x に対する Slope One では，

$$\begin{aligned} r_{u,x} &= \frac{\sum_{a|a \neq x} (\overline{\delta_{x,a}} + r_{u,a}) \phi_{x,a}}{\sum_{a|a \neq x} \phi_{x,a}} \\ &= \frac{\sum_{a|a \neq x} (\Delta_{x,a} + r_{u,a} \phi_{x,a})}{\sum_{a|a \neq x} \phi_{x,a}} \end{aligned} \quad (4.2)$$

により予測値を求める．

$r_{4,2}$ について Slope One を行うと

$$r_{4,2} = \frac{\sum_{j=1,5} (\overline{\delta_{2,j}} + r_{4,j}) \phi_{2,j}}{\sum_{j=1,5} \phi_{2,j}} = 1.67$$

となり，その誤差は，平均絶対誤差 (Mean Absolute Error:MAE) を用い 0.67 となる．

4.3 準備

4.3.1 関連研究

H. Polat ら [3] は、加法摂動化による協調フィルタリング方式を提案している。彼らの研究では、オリジナルデータ X に一様分布の乱数 R を加えた $Y = X + R$ について、平均値 $\sum_i Y_i = \sum_i X_i + \sum_i R_i \approx \sum_i X_i$ であることを仮定したナイーブな推薦方式である。従って、 Y を、主成分分析 (PCA) することで加えた乱数ノイズを取り除くことが出来ることが指摘されており [?], その安全性は低い。

そこで、本研究では、単純な PCA による解析が困難なランダムイズレスポンス方式を用いて、摂動化を行う。安全性は向上するが、[3] の様な単純な協調フィルタリングでは精度が期待できない。

4.3.2 摂動化と再構築

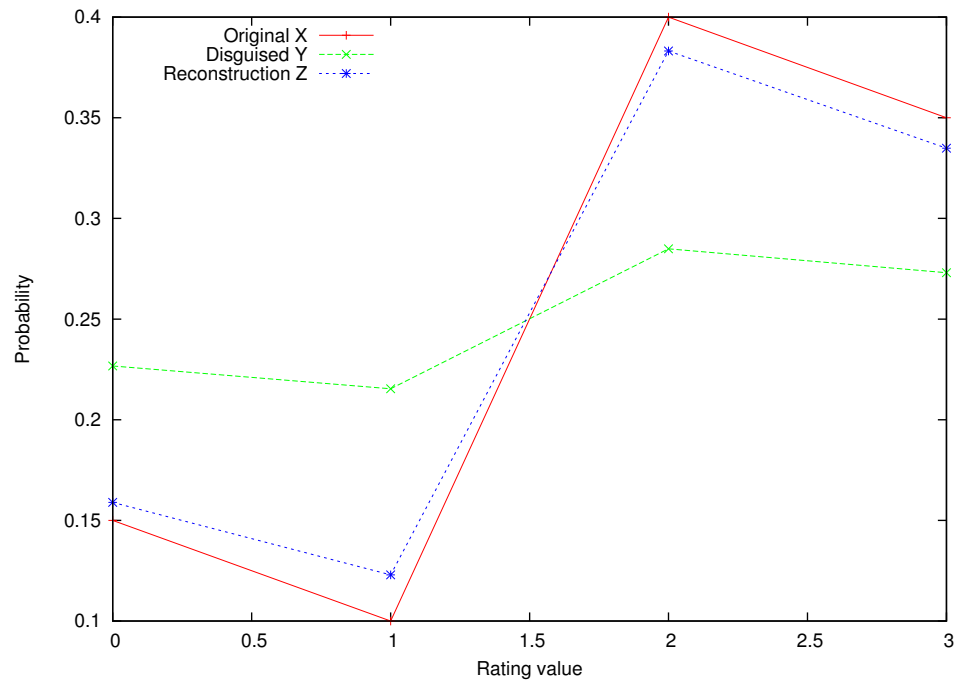
再構築問題 (Reconstruction Problem) とは、摂動化された $Y_1 = X_1 + R_1$ から、真の値 X の確率分布を見積もる問題である。R. Agrawal and R. Srikant[1] によって最初に発表された摂動化アルゴリズムである。秘匿したい情報に意図的にランダムノイズを乗せて、格納されたデータのプライバシーを保護する。

例えば、年齢が $x = 20$ 代であるという個人情報をもそのまま渡す代わりに、一様分布（またはガウス分布）の乱数 r を加え、 $y = x + r = 30$ の様に歪んだ値 y を登録する。30 という属性値を持った顧客がいても、本当に30代なのか乱数で40代から歪まされたのか、第三者には区別がつかない。

暗号化による方法と異なり、時間のかかる暗号化はなく、計算も各パーティで独立に計算できる。通信効率も計算効率も高い。大規模なデータベースにおいても適用可能である。

オリジナルデータ X をランダムイズレスポンスによって摂動化を行い、偽データ Y を作成する。評価値の集合を $V = \{1, 2, \dots, v\}$, 維持確率を p とする。評価値 x の摂動化 y は、

$$y = \begin{cases} x & \text{確率 } p \\ a \in V & \text{確率 } 1 - p \end{cases} \quad (4.3)$$

図 4.1: 再構築された評価値 Z の分布

4.4 提案方式

4.4.1 アイデア

オリジナルデータ X の評価値行列 R^X を，維持確率 p でランダムレスポンスによって摂動化を行い，偽データ Y の評価値行列 R^Y を作成する．この偽データ R^Y と p についてベイズ推定を行い，真のデータ X の再構築を行う．再構築の過程で得られた条件付き確率 $P(X|Y)$ を用いて，推薦精度を向上させることを試みる．

4.4.2 提案方式 - 摂動化 Slope One

Slope One を行う際に使用する共生起行列 $\Phi_{i,j}$ と平均差分行列 $\Delta_{i,j}$ を求める．偽データの評価値行列 R^Y のまま Slope One を行うと大きな誤差が起るため，再構築を行い R^X に近似した R^Z によって，Slope One で情報推薦を行う．

共生起行列

維持確率 p より，欠損値数の予測を行う．摂動化を行って生成した行列 R^Y ，維持確率 p よりオリジナルデータに期待欠損数を予測する． n 個の要素を持つオリジナルデータ X の欠損値数を表す確率変数を K_X ，摂動化したデータで観測した欠損値数を表す確率変数を K_Y

とする．損動化を行うと欠損値を維持する確率は p ，欠損値から評価値へと変化する確率は $1 - p$ である．同様に，評価値を維持する確率は $1 - \frac{1-p}{v}$ であり，評価値から欠損値へと変化する確率は $\frac{1-p}{v}$ である．これより，真の欠損値数 k_x の時に，損動化データに k_y 個の欠損値が生じる条件付き確率は

$$P(K_Y|k_x) = \sum_{j=0}^{K_X} \binom{K_X}{j} (1-p)^j p^{K_X-j} \binom{N-K_X}{K_Y-j} \left(1 - \frac{1-p}{v}\right)^{K_Y-j} \left(\frac{1-p}{v}\right)^{N-K_X-K_Y+j}$$

で与えられる評価値を $V = \{1, 2, 3\}$ ， $p = 0.4$ ， $n = 4$ とした時の k_y の確率分布を図 4.2 に示す． $k_x = 0$ であっても， $k_y > 0$ となる確率があることが分かる．この分布から欠損値数の期待値を $E[K_Y|K_X] = \sum_{k \in V} k P(k_Y = k|K_X)$ ，最尤値を $L[K_Y|K_X] = \arg \max_{k \in V} k P(k_Y = k|K_X)$ で求めることができる．図 4.2 から算出した期待値と最尤値を図 4.5 に示す．興味深いことに両者は必ずしも一致しない．

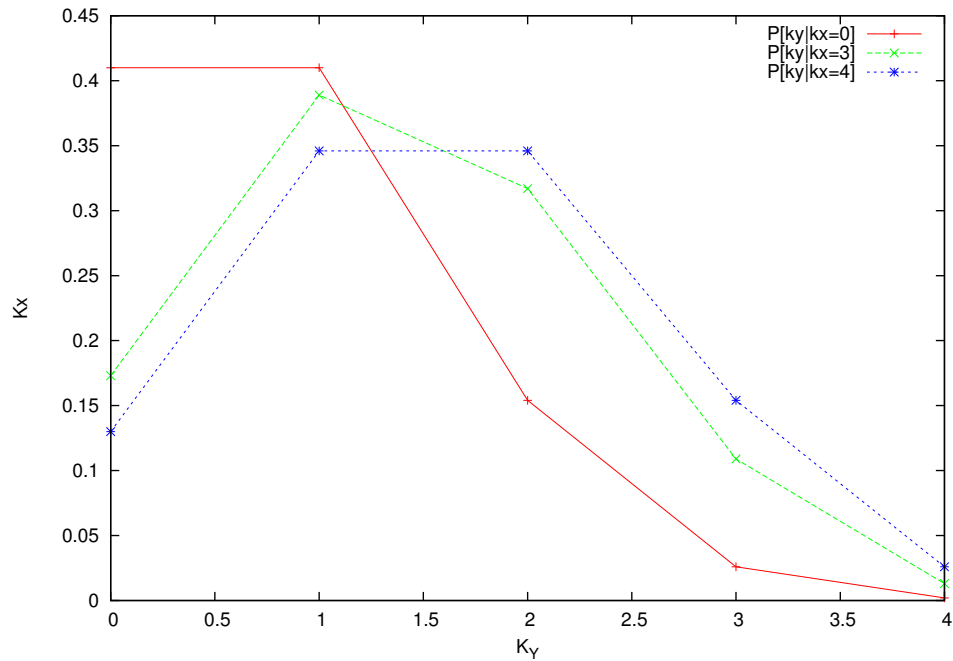


図 4.2: 損動化の欠損値数 $P(K_Y|K_X)$ の確率分布

事前確率 $P(K_Y|K_X)$ より，再構築アルゴリズムよりベイズ推定を行うことで事後確率 $P(K_X|K_Y)$ を求める．図 4.2 の分布から再構築した $P(K_X|K_Y)$ の確率分布を図 4.4 に示す．

$$P(K_X|K_Y) = \frac{P(K_Y|K_X)P(K_X)}{\sum_X P(K_Y|K_X)P(K_X)}$$

$P(K_X|K_Y)$ より，期待値 $E[K_X|K_Y]$ と最尤値 $L[K_X|K_Y]$ をそれぞれ求めたものを図 4.5 に示す．例えば， $K_Y = 1$ の時，期待値は $E[K_X] = 1.9$ であり，最尤値は $L[K_X] = 1$ である．

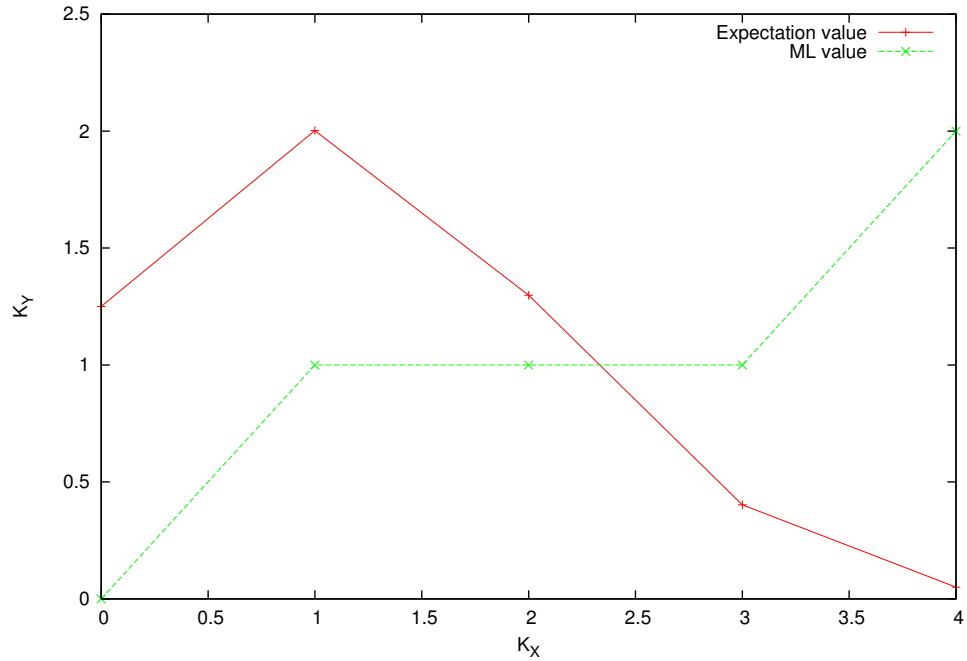


図 4.3: $P(K_Y|K_X)$ の期待値・最尤値の分布

Φ は共通に評価している欠損でないアイテム数なので、 K_Y が与えられた時の共生起数の期待値は $E[\Phi_X] = N - E[K_X|K_Y]$ で与えられる。

平均差分行列

摂動化評価値行列 R^Y と維持確率 p から定まる条件付き確率 $P(Y|X)$ より、摂動化 Slope One のための Δ_R を求める。偽データの Δ_Y より、再構築の際のベイズ推定で得られた条件付き確率 $P(X|Y)$ を用いて、ある列における差分の総和 Δ_Y から、真の差分総和 Δ_X を予測する。

$$P(\Delta_Y|\Delta_X) = \sum_{\Delta=\Delta_X} \sum_{\delta \in \Delta} P(\Delta_Y|\delta) \quad (4.4)$$

$P(\Delta_Y|\Delta_X)$ の最尤値を Δ_R と定める。

摂動化 Slope One

再構築を行った、共生起行列 Φ_Z と平均差分行列 Δ_Z が与えられた時、ユーザ u 、アイテム i の予測値は

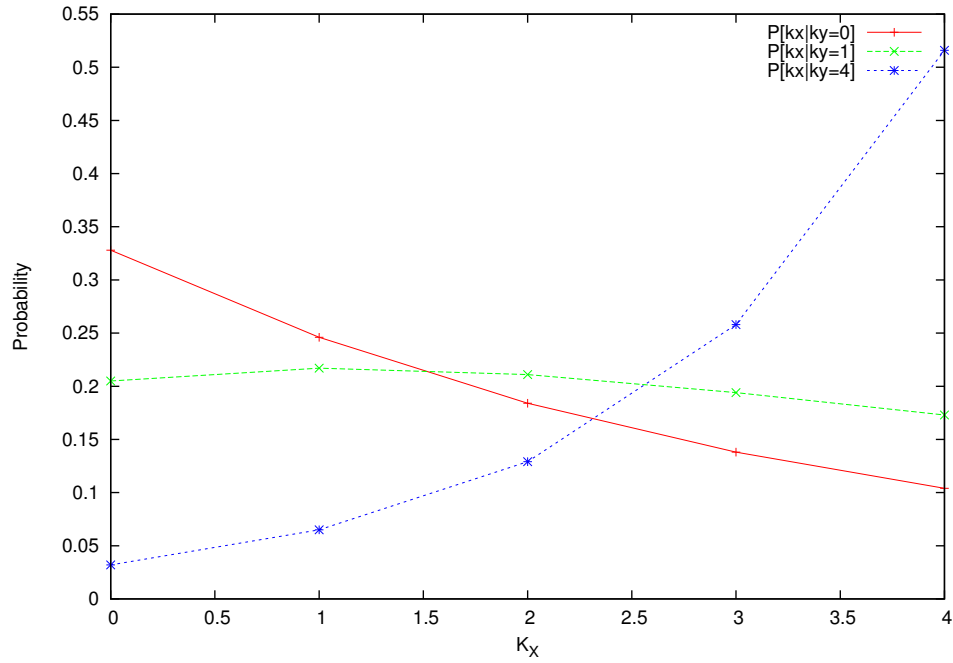


図 4.4: 真の欠損値数 $P(K_X|K_Y)$ の確率分布

表 4.1: 条件付き確率 $P(Y-X)$

$x \setminus Y$	0	1	2	3
0	0.37	0.18	0.23	0.22
1	0.19	0.36	0.23	0.22
2	0.18	0.17	0.44	0.21
3	0.18	0.17	0.22	0.43

$$\begin{aligned}
 r_{u,i}^Z &= \frac{\sum_i (\bar{\delta}_{u,i} + r_{u,i}) \phi_{u,i}}{\sum_i \phi_{u,i}} \\
 &= \frac{\sum_i (L[\Delta_i^X | \Delta_i^Y] + r_{u,i} (N - L[KX|KY]))}{\sum_i \phi_{u,i}}
 \end{aligned} \tag{4.5}$$

で与えられる .

4.4.3 数値例

本実験では、表 4.2 のオリジナルデータ R^X を評価値行列に、維持確率 $p = 0.4$ で摂動化した表 4.3 の偽データ R^Y を用いて、ユーザ数 $n = 4$ 、アイテム数 $m = 5$ 、評価値 $V \in \{1, 2, 3\}$ とする .

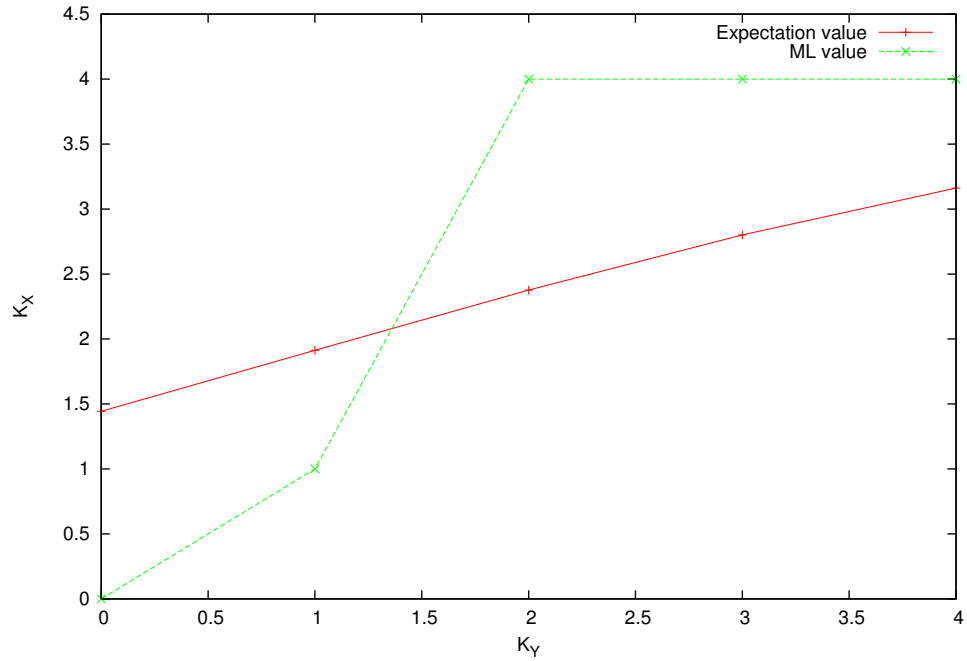


図 4.5: 欠損値数の期待値 $E[K_X|K_Y]$ ・最尤値 $L[K_X|K_Y]$ の分布

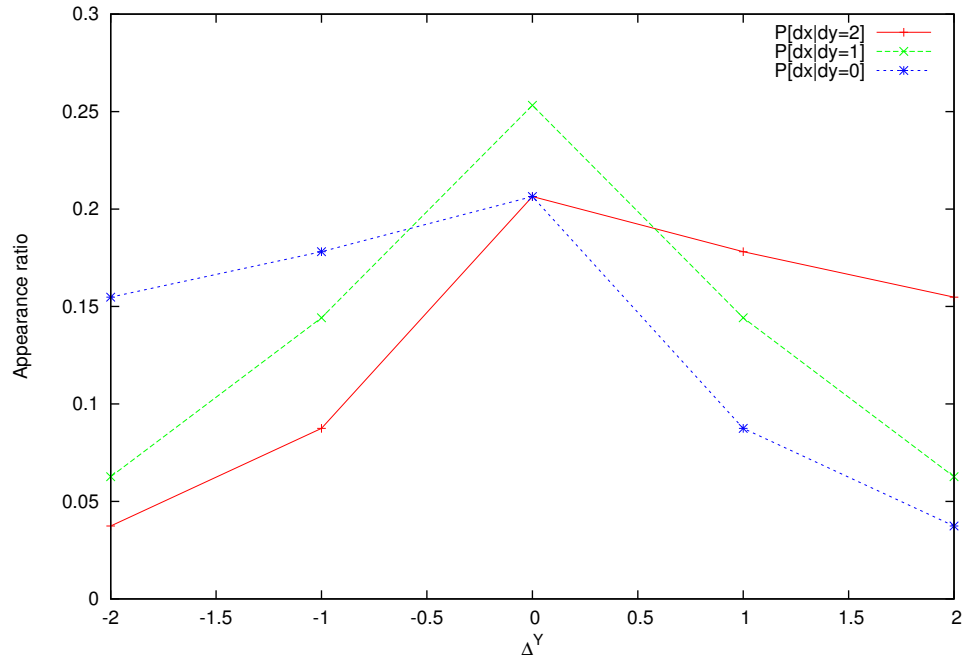
偽データの評価値行列 R^Y と $p, \text{Slope One}$ を行う際に使用する共生起行列 Φ_Y と平均差分行列 Δ_Y によって R^X を近似する R^Z を作成する。

表 4.2: Original Data (R^X)

	i_1	i_2	i_3	i_4	i_5
u_1	2	2	3	1	
u_2	1	3	2		3
u_3	2		3	3	2
u_4	3	2	3	2	2

4.4.4 結果

オリジナルデータ, 摂動化を行った偽データ, 提案手法である再構築データのそれぞれで Slope One を行い予測した値の MAE (Mean Absolute Error) と, 協調フィルタリングによって予測した値の MAE との比較を表 4.4 に示した。 $MAE^Z = \sum_{u,i} |r_{u,i}^X - r_{u,i}^Z|$ と定める。摂動化したデータにナイーブに協調フィルタリングを適用したものは, $MAE^Y = \sum_{u,i} |r_{u,i}^X - r_{u,i}^Y|$ である。偽データより提案手法である再構築データの方が誤差が 0.03 少ない。

図 4.6: $P(\Delta_X|\Delta_Y)$ の分布表 4.3: Disguised Data (R^Y)

	i_1	i_2	i_3	i_4	i_5
u_1	2	1		2	1
u_2	3	3	2		1
u_3	3	3	3	1	2
u_4	1	2	1	2	2

4.4.5 考察

図 4.8 の散布図は，Slope One の推薦値を真のデータとした時の摂動化データのみによる推薦値と提案手法による推薦値を図示している．Slope One を行った際，再構築データはオリジナルデータに近づいていることが分かる．これは，表 4.4 から分かるように，偽データより再構築データの方が誤差が 0.03 小さいことが表 4.4 から分かる．また，協調フィルタリングについても誤差を測り Slope One との比較を行った．全てのデータにおいて Slope One の方が誤差が少ない．最大で，0.24 の差がある．これら 2 点より，Slope One は，協調フィルタリングより精度が高く，摂動化との相性の良さが分かる．しかし，本結果で用いた評価行列は，オリジナルの推薦値の誤差が大きいため人工的で歪んでいた可能性がある．

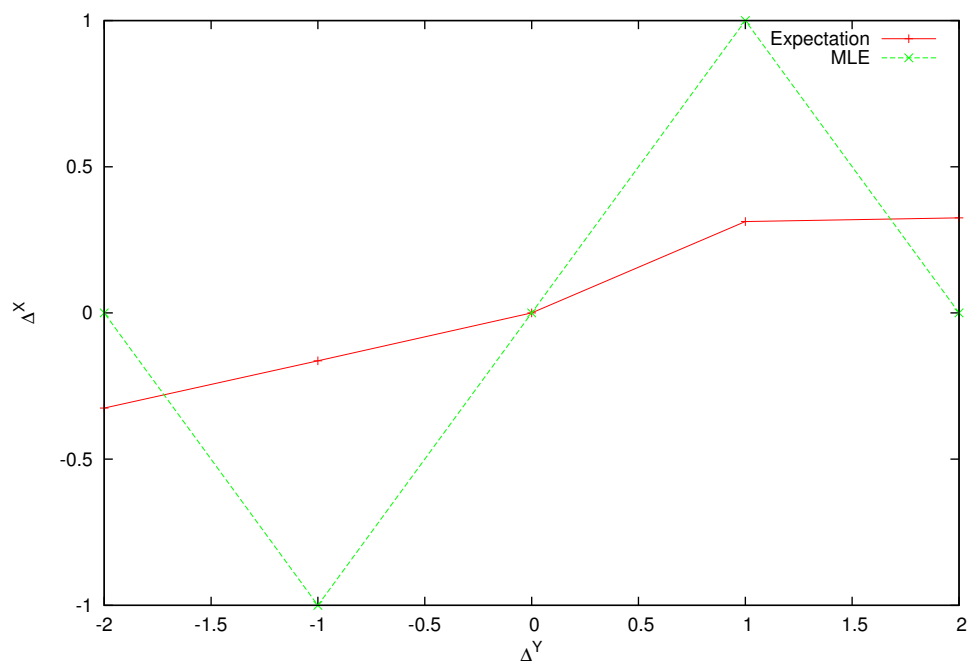


図 4.7: 真の差分 Δ_X ・ 期待値 $E[\Delta_X|\Delta_Y]$ ・ 最尤値 $L[\Delta_X|\Delta_Y]$ の分布

表 4.4: Mean Absolute Error

	CF	Slope One
Original	0.97	0.73
Disguised	1.03	0.89
Proposed	1.01	0.86

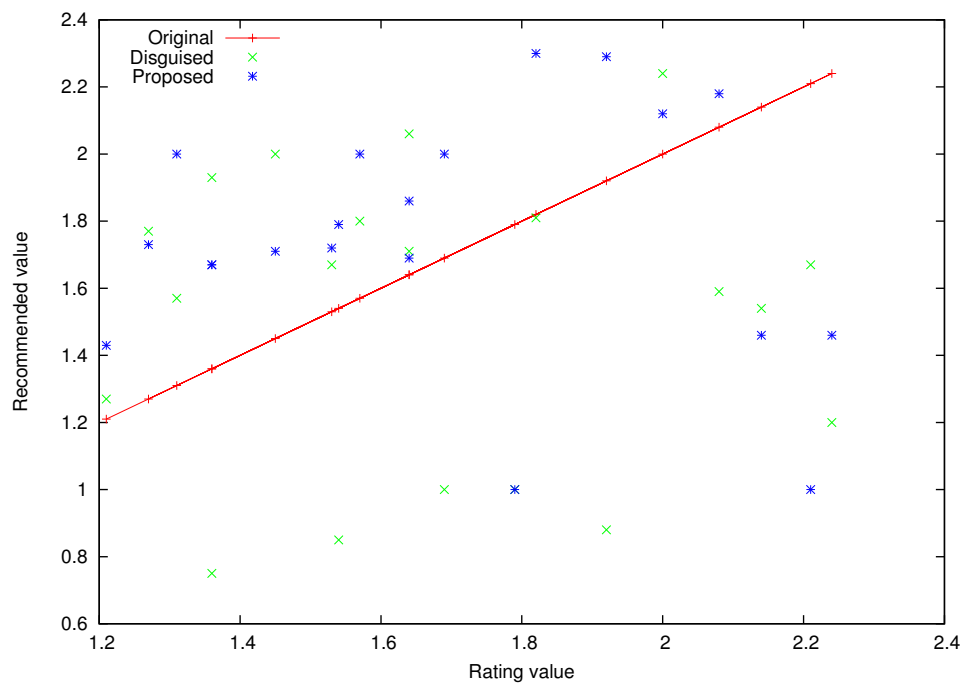


図 4.8: Slope One による推薦値の分布

第5章

アイテム依存の摂動化

5.1 アイテム依存を考慮した摂動化

H. Polat らは、加法摂動化による協調フィルタリング方式を提案している [3]。彼らの研究では、オリジナルデータ X に一様分布の乱数 R を加えた $Y = X + R$ について平均値 $\sum_i Y_i = \sum_i X_i + \sum_i R_i \approx \sum_i X_i$ であることを仮定したナイーブな推薦方式である。従って、Z. Huang らによって、 Y を主成分分析 (PCA) することで、加えた乱数ノイズを取り除くことが出来ることが指摘されており [2]、その安全性は低い。

菊池ら [10] は、分析対象のデータには、分析に重要で変化させたくない属性値（センシティブ属性）を保持することで、通常の攪乱・再構築に比べ計算量を削減する手法である。特に再構築アルゴリズムである反復ベイズ法について、センシティブ属性を含むテーブルを効率的に処理できるアルゴリズムを提案している。この手法は、センシティブ属性の値域が広い際に効果が大きい手法である。

S. Zhang ら [9] は、全てのアイテムを一様にアイテム不変に摂動化することは、精度の面でも、プライバシーの面でも効率が悪い事を指摘し、特異値分解 (SVD) を用いた摂動化方式を提案している。

本研究では、アイテムに依存した維持確率を使用し、再構築の計算コストと推薦精度を向上した情報推薦方式について検討する。

5.2 維持確率

図 5.4 のように，維持確率 p の値によって，誤差と相互情報量は変化する． p の値が小さい場合，大きく攪乱されるので，匿名化はされプライバシーは守られるが，誤差が大きくなる．一方， p の値が大きい場合は，あまり攪乱されず，誤差が少なく，個人の特定がされやすい状態となる．

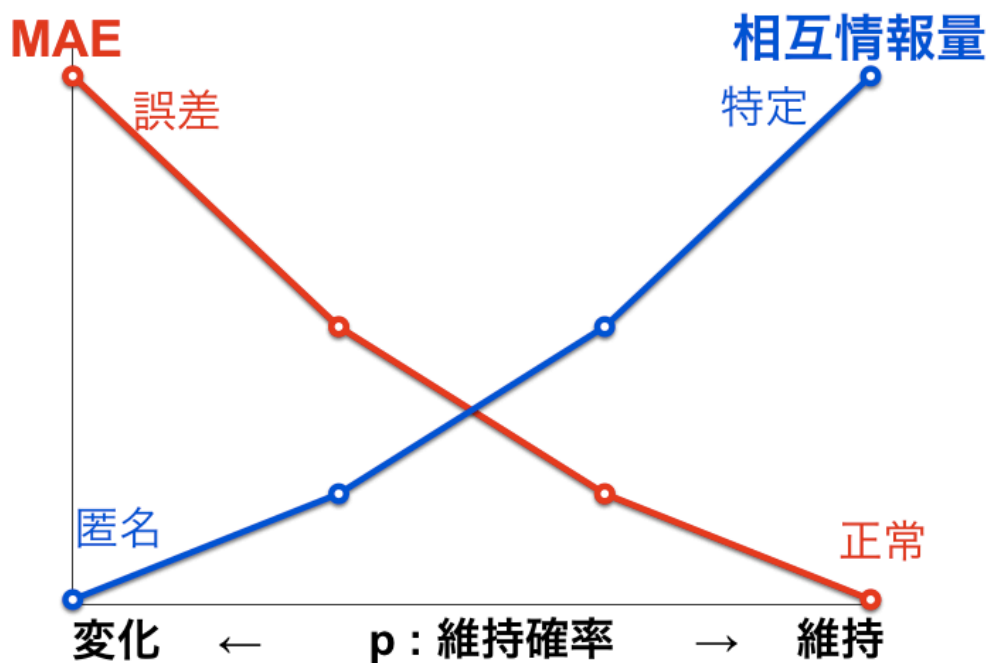


図 5.1: 2つの評価尺度

評価値データの中には，図 5.1 のように様々な評価をされているアイテムが存在する．

	欠損・小	欠損・大
一様評価	$p \rightarrow$ 大	$p \rightarrow$ 小
多様評価	$p \rightarrow$ 小	$p \rightarrow$ 中

そのため，各アイテムに適した維持確率を定めることによってプライバシーを守りながら誤差の少ない情報推薦を行うことができる．

5.2.1 数値例

図 5.2 に、数値例から MAE と、エントロピーの関係を示す。MAE は、維持確率に対して、再構築の末、最尤法で推薦値 $L(X) = 2$ とした時の分布であり、エントロピーは、 $P(X = 0) = 1/4, P(X = 1) = 0, P(X = 2) = 1/2, P(X = 3) = 1/4$ の時、維持確率 p で摂動化した時の Y の情報量の変化である。

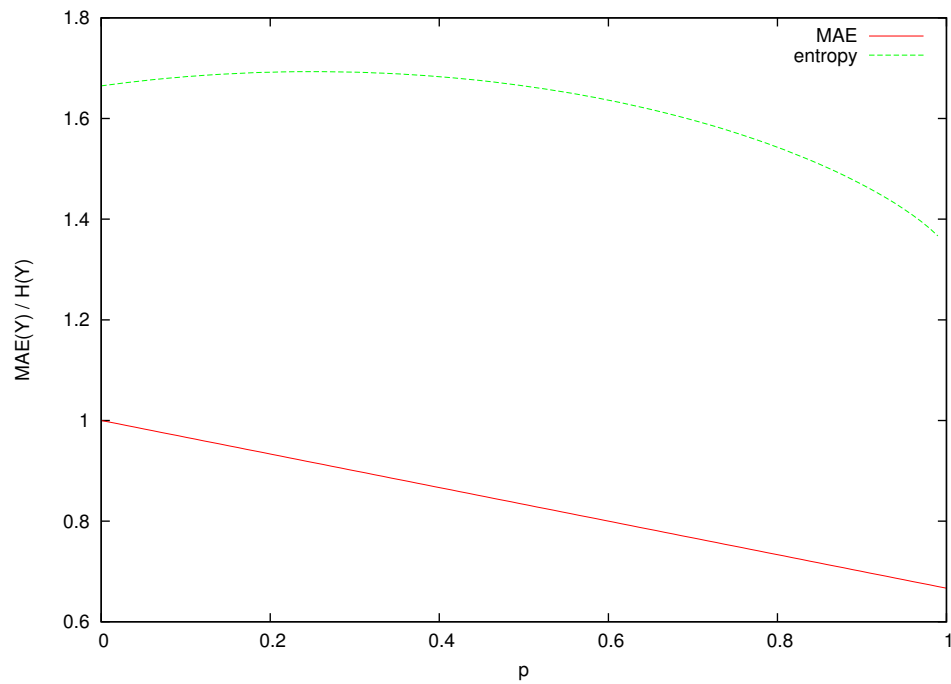


図 5.2: MAE とエントロピーの関係

5.3 提案方式

表 5.2: 2 種類のスパース率を持つ評価値データセット (Original Data)

	i_1	i_2	i_3	i_4	i_5	i_6	i_7	i_8	i_9	i_{10}
u_1	5	4	5	5		3			3	2
u_2	5	3		5	3	4			5	
u_3	5	4	5	3	1	4	4		3	
u_4	5	5	5	4		3	4			
u_5	5	3	3	5	3	4		4		2
u_6	5		5	4	3	4			5	
u_7	5	4	5	4	3	3				3
u_8	5	3	4	3	2	3	3	2	4	
u_9	4	5	5	4	3	4				4
u_{10}	4	2	4	2		4	2	5	4	
p	0.8						0.1			

5.3.1 アイデア

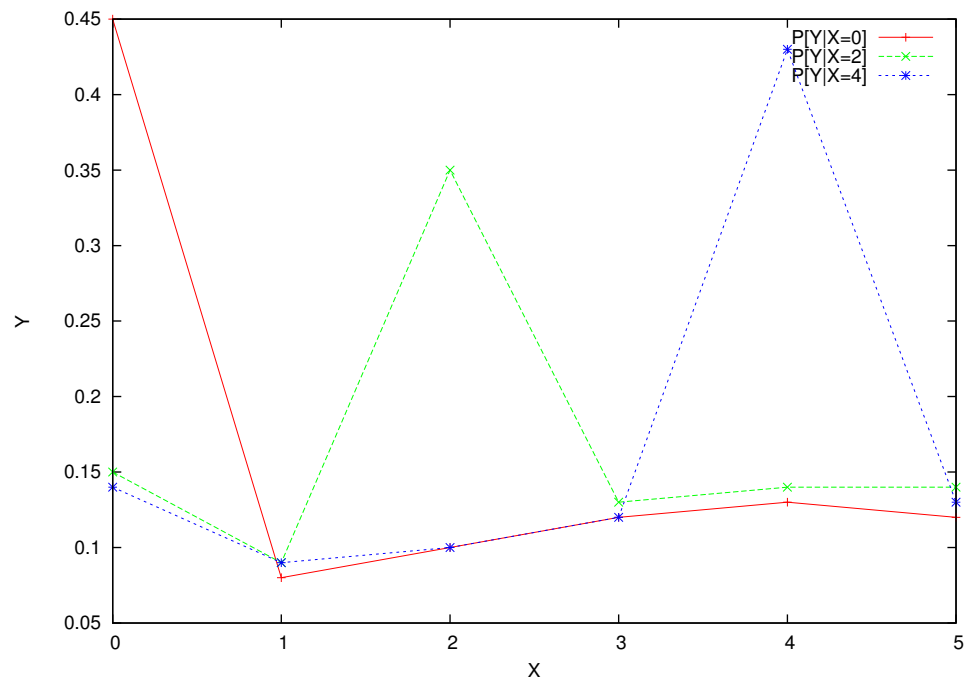
オリジナルデータ X の評価値行列 R^X について共生起行列 Φ を求め、アイテムごとの欠損値数を考慮して維持確率 p_i を決める。アイテム毎に randomized response によって摂動化を行い、摂動化データ R^Y と p についてベイズ推定を行い、真のデータ X の再構築を行う。再構築の過程で、得られた条件付き確率 $P(X|Y)$ を用いて、推薦精度を向上させることを試みる。

- アイテム不変の摂動化

全てのアイテムを同じ維持確率によって摂動化を行う。

表 5.3: アイテム不変維持確率を用いた摂動化

	i_1	i_2	i_3
u_1	4	5	
u_2	3		
u_3	5	4	
u_4			1
u_5	5	2	
p	0.5		

図 5.3: 摂動化 Y の $P(Y|X)$ の確率分布

- アイテム依存の摂動化

各アイテムのスパース率によって維持確率を変える。

表 5.4: アイテム依存維持確率を用いた摂動化

	i_1	i_2	i_3
u_1	4	5	
u_2	3		
u_3	5	4	
u_4			1
u_5	5	2	
p	0.8	0.6	0.2

5.3.2 提案方式 - アイテム依存維持確率

アイテム i のスパース率 s_i によって、アイテムごとの維持確率を定める。アイテム i のスパース率 s_i とは、アイテム i を評価済のユーザの密度で定める。データベースが疎であるとスパース率は高く、密であると低い値を示す。

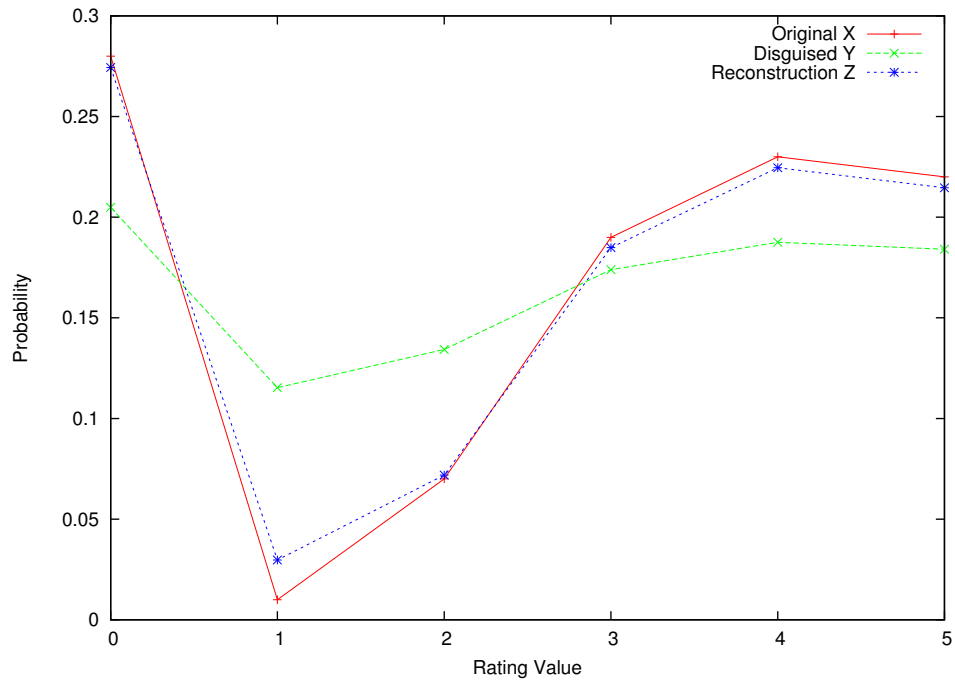


図 5.4: 摂動化と再構築による評価値分布の変化 (維持確率 $p = 0.4$)

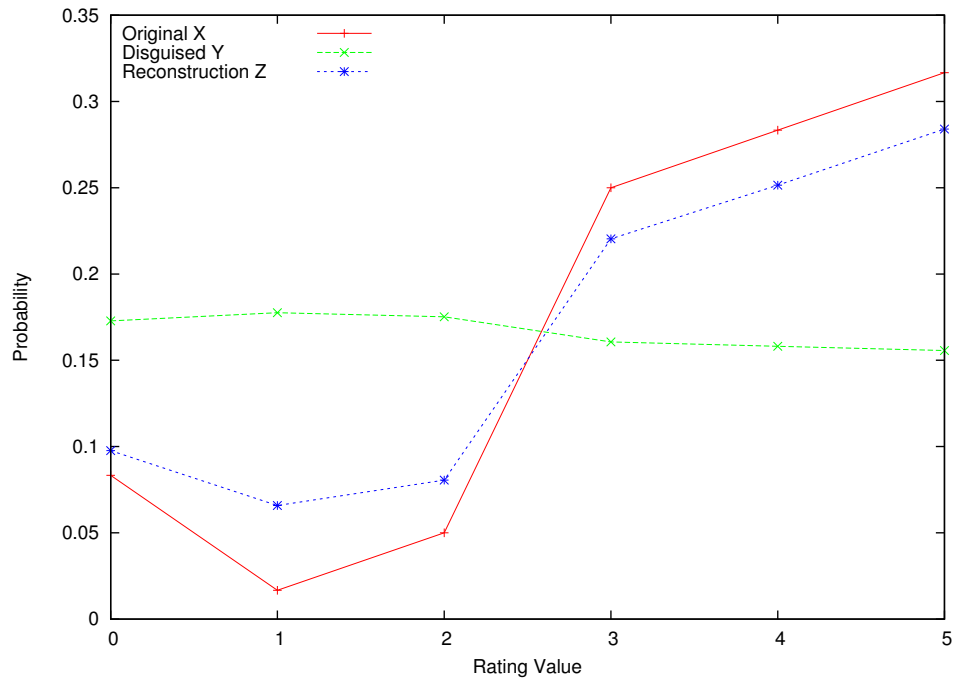


図 5.5: データセット A における評価値分布 ($p_A = 0.8$)

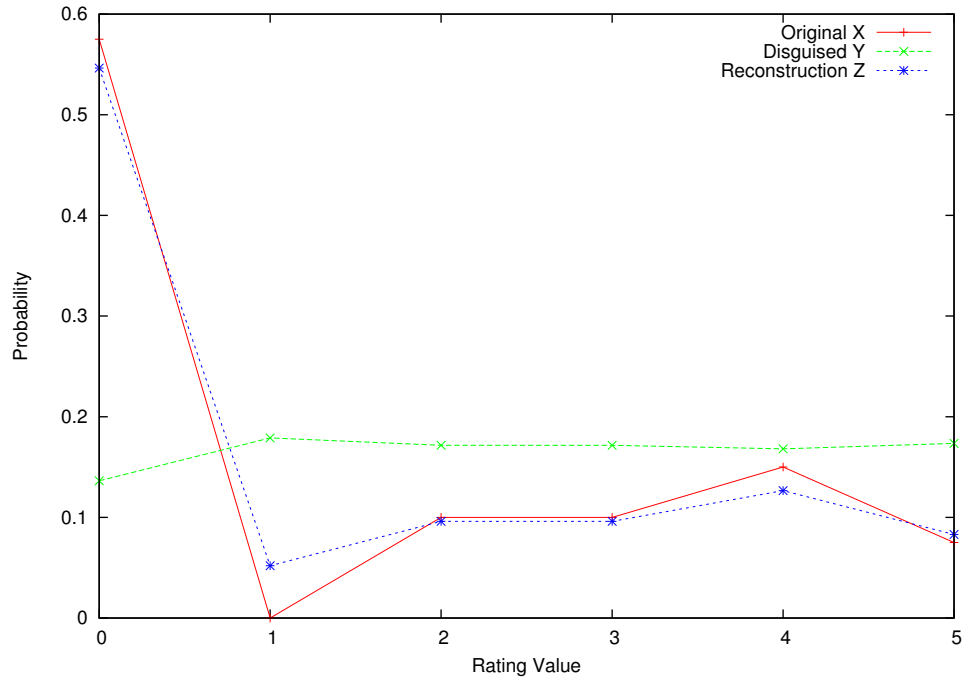


図 5.6: データセット B における評価値分布 ($p_B = 0.1$)

アイテム i のスパース率と共生起数 ϕ_i には

$$\phi_i = 1 - s_i \quad (5.1)$$

の関係がある。共生起数は、アイテム毎に安全な電子投票プロトコルを用いて求め、全員で共有して持つ。スパース率 s_i のアイテム i を摂動化する維持確率 p_i は、要求する精度と守るべきプライバシーの両方の観点から定める。全アイテムでの目標精度 MAE^* とプライバシー ϕ_Y^* を固定し、再構築したデータ Z の $MAE(Z) < MAE^*$, $\phi_Z^* \geq \phi_Y^*$ を満たす最小の p_i を決める。

5.3.3 実験

実験には、MovieLens [19] のデータセット (ユーザ数 $u = 943$, アイテム数 $i = 1682$, 評価値数 100,000 件) からランダムに抽出し利用した表 5.4 のデータを用いる。ここで、 i_1 から i_6 までを A , i_7 から i_{10} までを B とおき、各々 $p_A = 0.8$, $p_B = 0.1$ の維持確率で摂動化する。このデータをアイテム不変の維持確率 $p = 0.4$ で摂動化し、第 4 章で提案した手法で評価値の予測をした場合と、本提案のアイテム依存の維持確率で摂動化して評価した場合とを比較する。

5.3.4 実験結果

オリジナルデータ, 全てのアイテムを同じ維持確率によって摂動化を行ったアイテム不変摂動化データ (item invariant), アイテムのスパース率によって維持確率を変動させたアイテム依存摂動化データ (item dependence) のそれぞれで第4章の方式による Slope One を行い予測した値を比較する. MAE(Mean Absolute Error) による推薦値の誤差を表5.5に示す. 図5.3は, 摂動化したデータの確率分布, 図5.4は, 再構築によって補正された確率分布, 図5.5と図5.6は表5.4のAとBに対応するデータに適用した摂動化と再構築の確率分布である. 図5.7は, アイテム不変とアイテム依存の情報推薦の誤差の分布を示している.

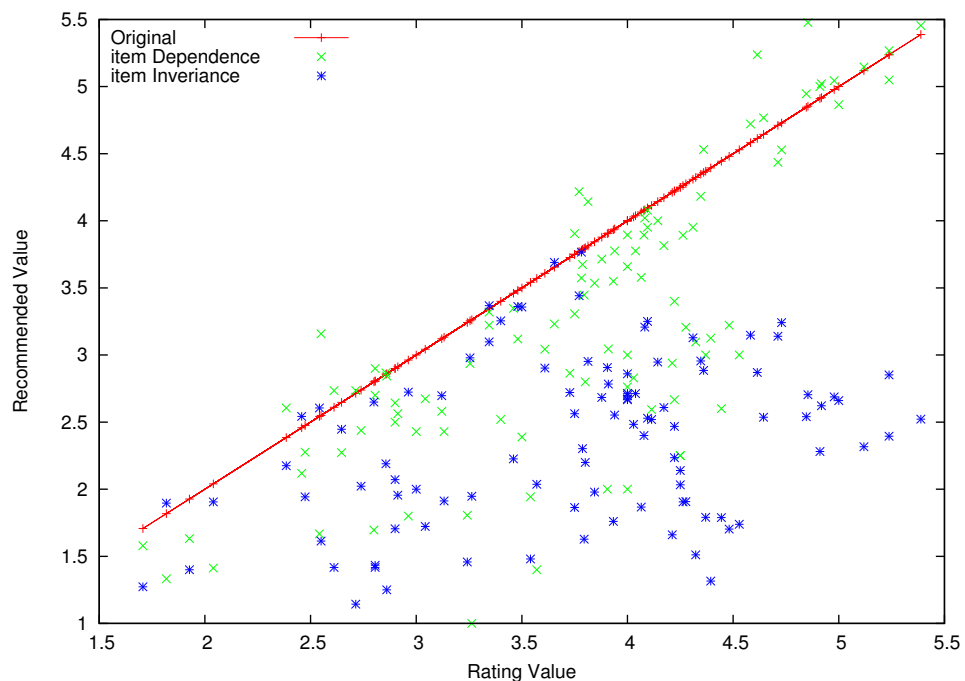
図 5.7: 再構築された評価値 Z の分布

表 5.5: Mean Absolute Error

	Original	item Invariance	item Dependence
MAE	0.57	0.81	0.61

5.3.5 考察

図5.7により, アイテム依存の維持確率を使用することで, 摂動化データはオリジナルデータに近い. 表5.5により, MAEの値もそれを裏付けている. 従って, アイテム依存の維持確率を使用した時より誤差の少ない情報推薦を行うことが示された.

第6章

評価実験

6.1 実験データ

実験に利用したデータは、MovieLens[19] のデータセット (ユーザ数 $n = 943$, アイテム数 $m = 1682$, 評価値 100,000 件) を利用した .

表 6.1: MovieLens 評価値数

評価値	評価数
欠損	1,486,126
1	6,110
2	11,370
3	27,145
4	34,174
5	21,201

6.2 維持確率

MovieLens の評価値で , 分布による摂動化・100 回の再構築を行った . 結果を図 6.1, 6.2, 6.3, 6.4, 6.5, 6.6, 6.7, 6.8, 6.9 に示す .

維持確率 p の値が小さければ , 摂動化データはオリジナルデータから離れプライバシーが保護される . 逆に , 維持確率 p の値が大きければ , あまり攪乱はされず , 再構築のための計算回数が少なくなる .

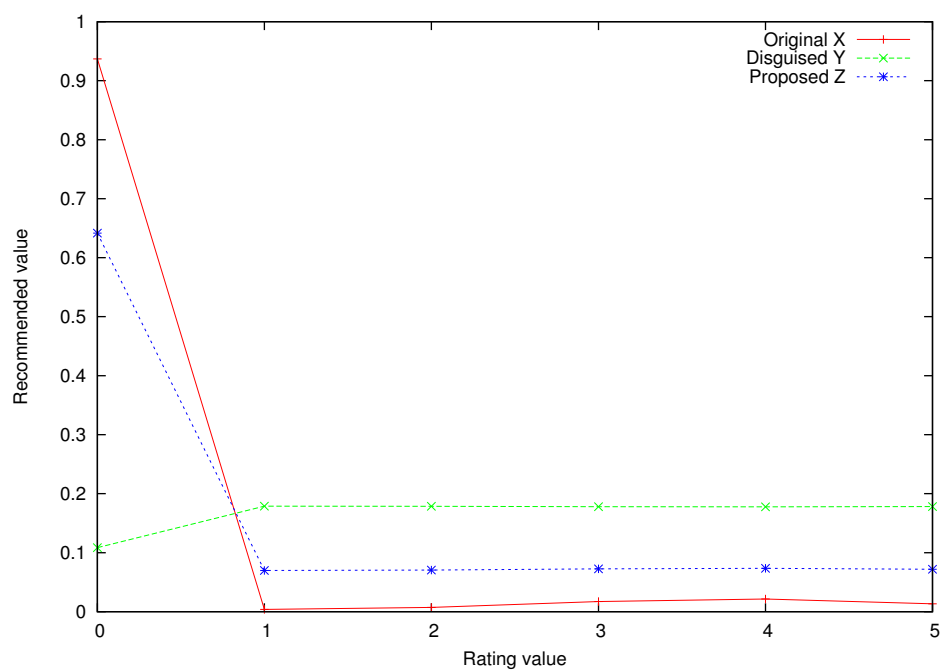
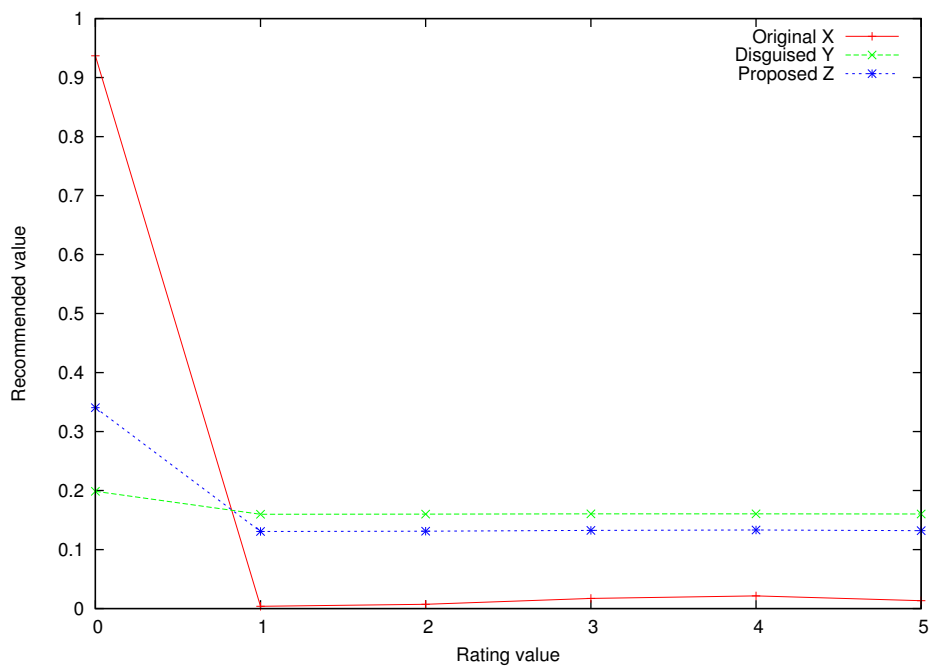
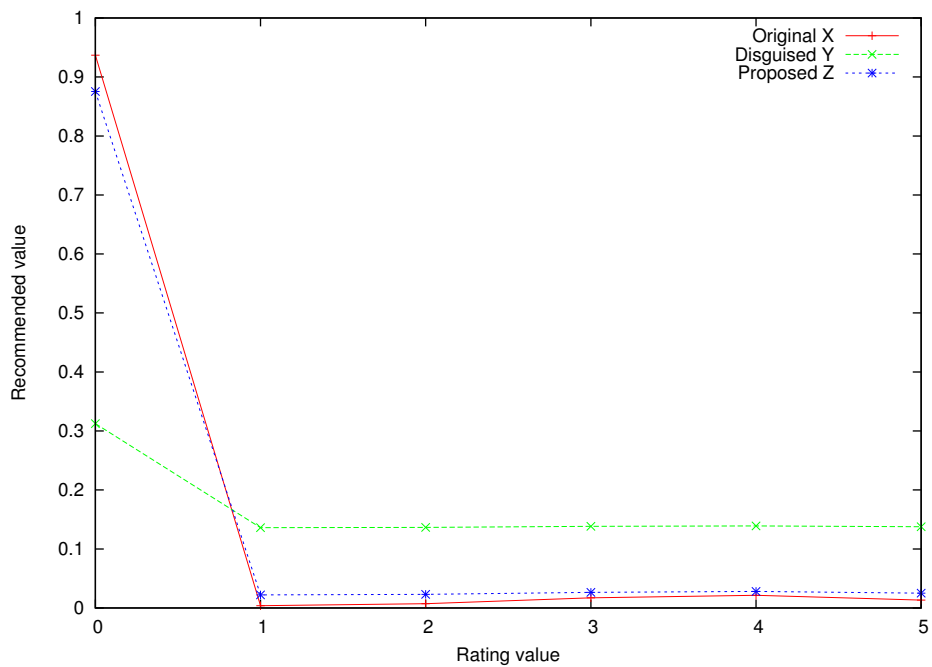
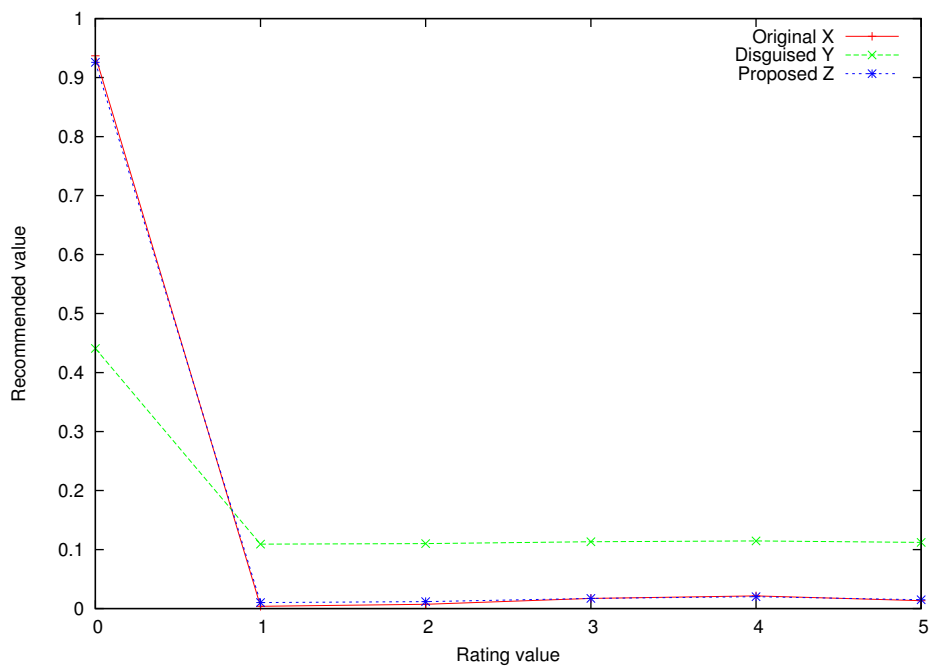
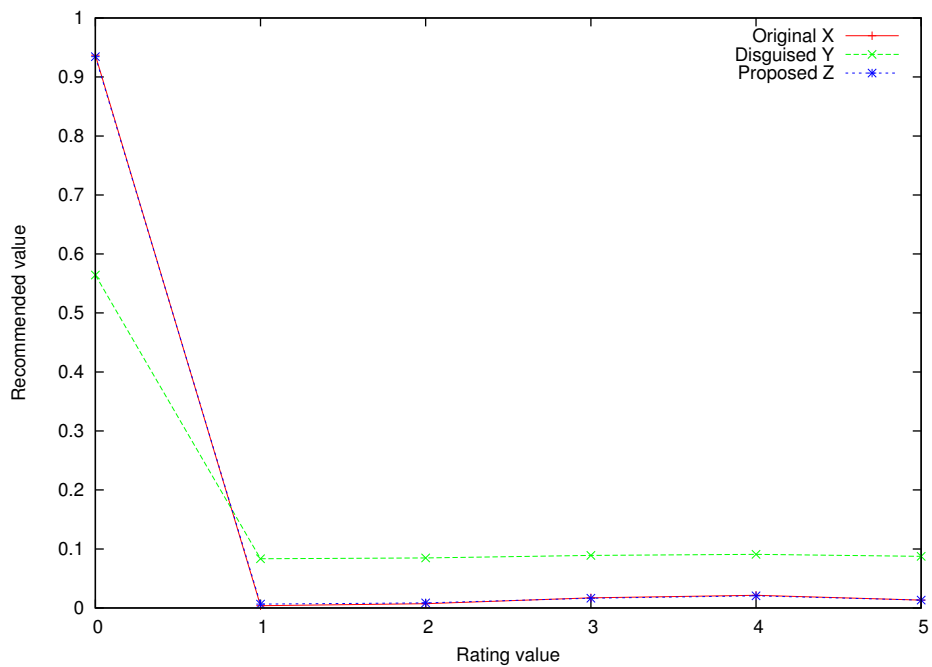
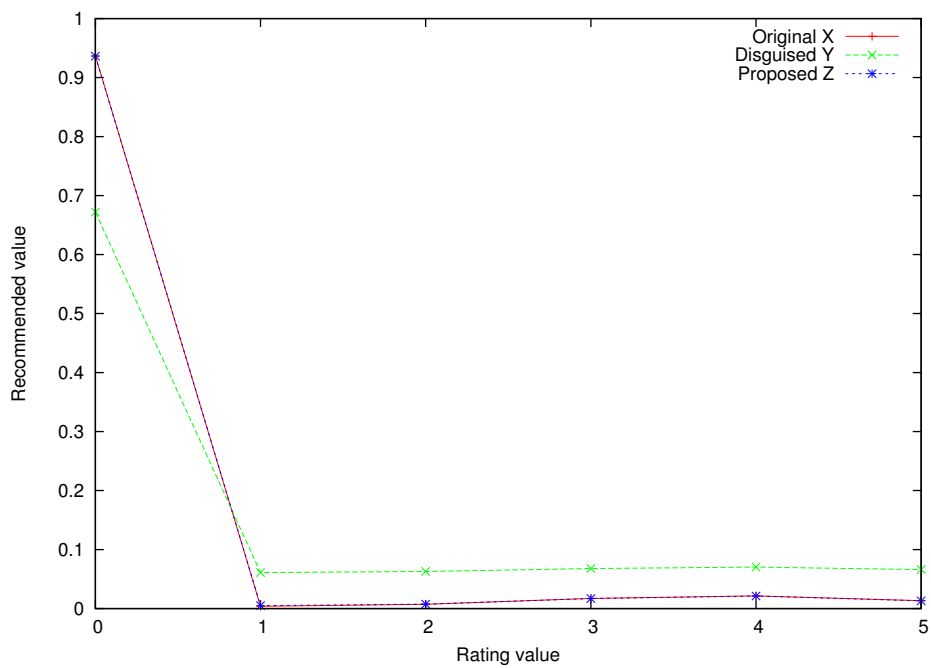
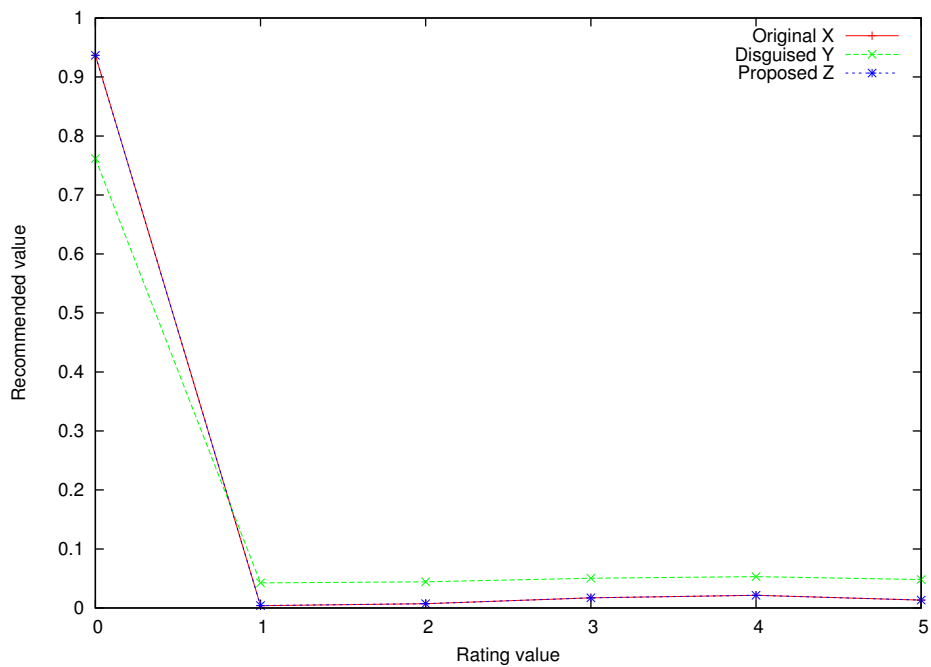
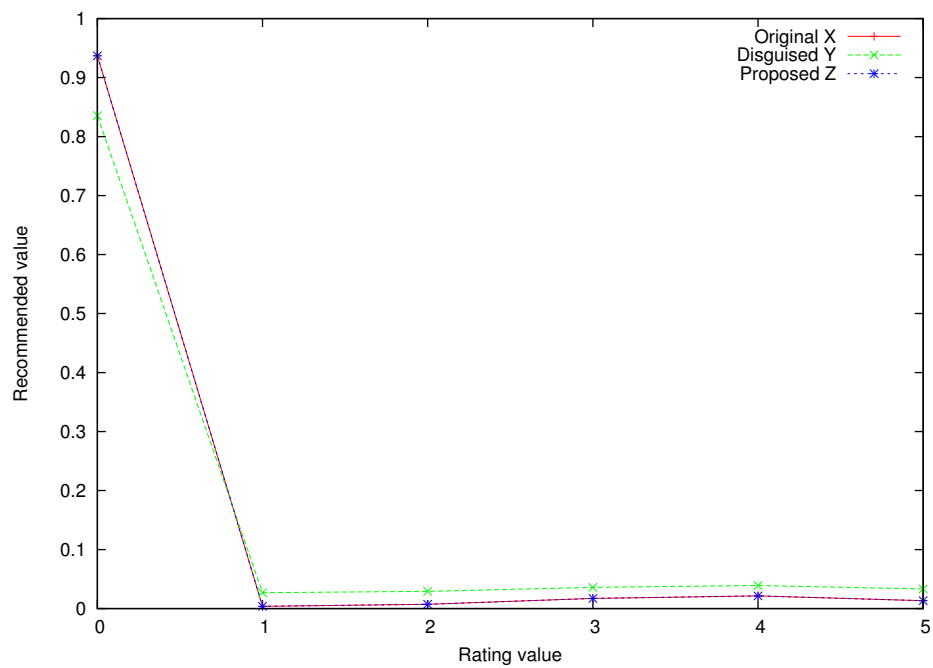
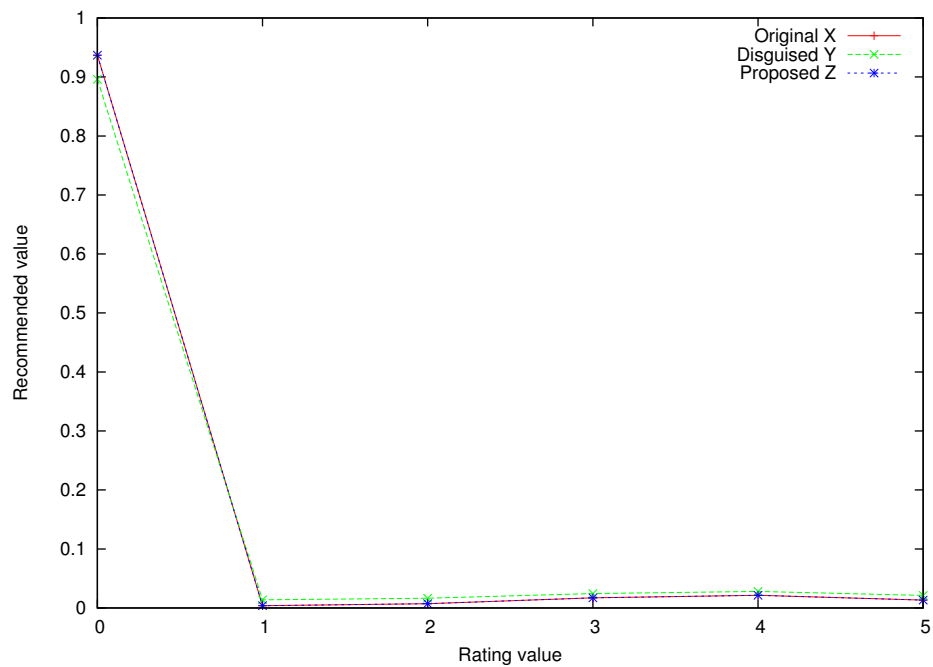


図 6.1: 維持確率 $p = 0.1$

図 6.2: 維持確率 $p = 0.2$ 図 6.3: 維持確率 $p = 0.3$

図 6.4: 維持確率 $p = 0.4$ 図 6.5: 維持確率 $p = 0.5$

図 6.6: 維持確率 $p = 0.6$ 図 6.7: 維持確率 $p = 0.7$

図 6.8: 維持確率 $p = 0.8$ 図 6.9: 維持確率 $p = 0.9$

6.3 精度評価

6.3.1 平均絶対誤差 (Mean Absolute Error)

精度評価に、平均絶対誤差 (Mean Absolute Error:MAE) を使用した。MAE は以下の式で計算される。

$$MAE = \frac{\sum_{l=1}^L |r_l - p_l|}{L} \quad (6.1)$$

ここで、 p は予測された評価値、 L はサンプル数を表す。

6.3.2 Slope One

MovieLens からランダムに抽出したデータを用いる。これを original とする。original を様々な維持確率で摂動化を行い誤差を測定する。

表 6.2: 維持確率の変化による MAE

維持確率	MAE	* との差
original(*)	0.574	—
$p = 0.1$	0.539	-0.035
$p = 0.2$	0.783	0.209
$p = 0.3$	0.856	0.282
$p = 0.4$	0.855	0.281
$p = 0.5$	0.812	0.238
$p = 0.6$	0.653	0.079
$p = 0.7$	0.644	0.070
$p = 0.8$	0.730	0.156
$p = 0.9$	0.452	-0.122

original に比べ、MAE の値が大きければ、誤差が大きいためプライバシーが守れていると言える。逆に、MAE の値が小さければ、差があまりなく、再構築が行いやすいと言える。

それぞれの MAE に対し、 p の値が小さければ MAE の値は大きく、 P の値が大きければ MAE の値は小さいという予想の基に実験を行った。しかし、実際にはバラツキがあり、より多くの評価値を扱った中で MAE を求める必要がある。

第7章

結論と今後の課題

7.1 結論

摂動化を用いてプライバシーを保護した情報推薦を協調フィルタリングと Slope One を用いて行った。暗号から摂動化に変更することによって、計算コストは格段に安くなる。また、協調フィルタリングから Slope One へと推薦方式を変えることで、誤差の少ない情報推薦を行うことが出来る。

さらに、アイテム依存の維持確率を使用することによって、プライバシーを守りながら精度の高い情報推薦を行うことができることを示した。

7.2 課題

7.2.1 アイテム依存維持確率

今回提案した維持確率はスパース率のみに着目したものであった。プライバシーをより保護するためには、評価の値も関わってくるため、より適した維持確率を求める必要がある。

参 考 文 献

- [1] R. Agrawal and R. Srikant, “Privacy-Preserving Data Mining”, ACM SIGMOD 2000, pp. 439-450, 2000.
- [2] Z. Huang, W. Du and B. Chen, “Deriving Private Information from Randomized Data”, ACM SIGMOD 2005, pp. 37-48, 2005.
- [3] H. Polat and W. Du, “Privacy-Preserving Collaborative Filtering using Randomized Perturbation Techniques”, ICDM 2003, pp. 1-15, 2003.
- [4] 青木 良樹, 菊池 浩明, “擬準同型性を満たす類似度による分散協調フィルタリングプロトコル”, 2011 年暗号と情報セキュリティシンポジウム (SCIS2011), pp. 1-6, 2011.
- [5] 麻生 英樹, 小野 智弘, 本村 陽一, 黒川 茂莉, 櫻井 彰人, “協調フィルタリングと属性ベースフィルタリングの統合について”, 信学技報 NC 2006, pp. 55-59, 2006.
- [6] 麻生 英樹, “階層ベイズによる協調フィルタリング”, 信学技報 IBISML, pp. 57-62, 2010.
- [7] D. Leniel and A. Maclachlan, “Slope One Predictors for Online Rating-Based Collaborative Filtering”, Society for Industrial Mathematics, pp. 1-5, 2005.
- [8] A. Basu, H. Kikuchi and J. Vaidya, “Privacy-preserving weighted Slope One for Item-based Collaborative Filtering”, IFIPTM 2011 Federated Workshop TP-DIS’11, pp. 1-12, 2011.
- [9] S. Zhang, J. Ford and F. Makedon, “A Privacy-preserving Collaborative Filtering Scheme with Two-way Communication”, ACM EC’06, pp. 316-323, 2006.
- [10] 菊池 亮, 五十嵐 大, 千田 浩司, 濱田 浩気, “属性値を保持する際に効果的な攪乱・再構築法”, Computer Security Symposium 2011, pp. 438-443, 2011.
- [11] 豊田, 宮川, 側高, 伊東, “匿名性グループ間の要素数の変化を比較可能な匿名化手法の実現”, Computer Security Symposium 2011, pp. 432-437, 2011.

- [12] Neal Lathia, Stephen Hailes, Licia Capra, “Private Distributed Collaborative Filtering Using Estimated Concordance Measures”, In ACM 2007 Conference on Recommender Systems (RecSys). Minneapolis, Minnesota, USA. October 19-20, 2007.
- [13] A. Agresti, “Analysis of Ordinal Categorical Data”, , 1984.
- [14] John S. Breese, David Heckerman Carl Kadie, “Empirical Analysis of Predictive Algorithms for Collaborative Filtering”, The 14th conference on Uncertainty in Artificial Intelligence, 1998.
- [15] 平山 巧馬, 小柳 滋, “協調フィルタリングにおける相関係数法の予測性能向上”, 電子情報通信学会論文誌. D, 情報・システム, Vol J90-D(2), 223-232, 2007.
- [16] 木澤, 菊池, “プライバシー協調フィルタリングにおける利用者評価行列の次元削減”, コンピュータセキュリティシンポジウム 2008(CSS2008), pp. 509-514, 2008.
- [17] 木澤, 磯崎, 菊池, “秘匿積集合プロトコルを利用したプライバシー協調フィルタリングの提案”, 2009 年暗号と情報セキュリティシンポジウム (SCIS2009), 2009.
- [18] 多田, 菊池, “秘密分散ベースの秘匿関数計算を用いたプライバシー保護情報推薦方式”, 2011 年暗号と情報セキュリティシンポジウム (SCIS2011), 2011.
- [19] Grouplens Data Sets, (<http://grouplens.org/>), 2006.
- [20] Netflix prize, (<http://www.netflixprize.com/>), 2007.
- [21] Amazon.co.jp, (<http://amazon.co.jp>)
- [22] J.Canny, “Collaborative Filtering with Privacy”, IEEE Conf. on Security and Privacy, pp. 45-47, Oakland CA, 2002.
- [23] 神鷹, “推薦システムのアルゴリズム (1)”, 人工知能学会誌, Vol. 22 No. 6, pp. 826-837, 2007.
- [24] 神鷹, “推薦システムのアルゴリズム (2)”, 人工知能学会誌, Vol. 23 No. 1, pp. 89-103, 2007.
- [25] 神鷹, “推薦システムのアルゴリズム (3)”, 人工知能学会誌, Vol. 23 No. 2, pp. 248-263, 2008.
- [26] 高島 秀佳, 山岸 英貴, 平澤 茂一, “欠損値推定による協調フィルタリング手法” 情報科学技術フォーラム一般講演論文集, Vol. 4, No. 1, pp. 15-16, 2005. 人工知能学会誌, Vol. 23 No. 2, pp. 248-263, 2008.

- [27] Jonathan L. Herlocker, Joseph A. Konstan, Al Borchers, John Riedl, “An algorithmic framework for performing collaborative filtering”, SIGIR '99 Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval, 1999.
- [28] Jaideep Vaidya and Chris Clifton, “Privacy preserving naive bayes classifier for vertically partitioned data”, In 2004 SIAM International Conference on Data mining, pp. 522-526, 2004.
- [29] Jaideep Vaidya and Chris Clifton, “Secure set intersection cardinality with application to association rule mining”, Journal of Computer Security, Vol. 13, No. 4, pp. 593-622, 2005.
- [30] R.L.Rivest,A.Shamir,and L.Adelman, “A Method for Obtaining Digital Signature and Public-key Cryptosystems”, Communications of the ACM 21,2, pp. 120-126, 1978.
- [31] T. ElGamal, “ A Public Key Cryptosystem and a Signature Scheme Based on Discrete Logarithms ”, IEEE Trans. on Information Theory, IT-31(4), pp.469-472, 1985.
- [32] P. Paillier, “Public-Key Cryptosystems Based on Composite Degree Residuosity Classes”, Proc. *EUROCRYPT'99*, LNCS 1592, pp. 223-238, 1999.
- [33] Haifeng Yu, Chenwei Shi, Kaminsky, M., Gibbons, P.B., and Feng Xiao, “DSybil: Optimal Sybil-Resistance for Recommendation Systems”, in IEEE Symp. on Security and Privacy, pp. 283-298, IEEE, 2009.
- [34] Wenliang Du and Mikhail J. Atallah, “Privacy-preserving statistical analysis”, In Proceeding of the 17th Annual Computer Security Applications Conference, pp. 10-14 2001.
- [35] Yehuda Lindell and Benny Pinkas, “Privacy preserving data mining”, Journal of Cryptology, 15(3), pp. 177-206, 2002.
- [36] Rakesh Agrawal, Alexandre Evfimievski, and Ramakrishnan Srikant, “Information shareing across private databases”, in proc. of ACM SIGMOD International Conference on Management of Data, 2003.
- [37] Michael J. Feedman, Kobbi Nissim, and Benny Pinkas, “Efficient private matching and set intersection”, in Eurocrypt 2004, IACR, 2004.

-
- [38] G. Jagannathan and R. N. Wright, “Privacy-Preserving Distributed k -Means Clustering over Arbitrarily Partitioned Data”, *ACM KDD '05*, 2005.
- [39] Andrew C. Yao, “How to generate and exchange secrets”, In Proc. of the 27th IEEE Symposium on Foundations of Computer Science, pp. 162-167, 1986.
- [40] Dahlia Malkhi, Noam Nisan, Benny Pinkas, and Yaron Sella, “Fairplay – A Secure Two-Party Computation System”, Usenix Security Symposium, 2004.
- [41] Koji Chida, Dai Ikarashi, and Katsumi Takahashi, “Tag-Based Secure Set-Intersection Protocol and Its Application”, in proc. of Computer Security Symposium (CSS 2009), IPSJ, 2009 (in Japanese).
- [42] L. F. Cranor, “I Didn’t Buy it for Myself, Privacy and E-Commerce Personalization”, WPES 2003, Washington, DC, USA, pages 111-117, 2003.
- [43] J. Canny: Collaborative Filtering with Privacy, *IEEE Conf. on Security and Privacy*, Oakland CA, May 2002.
- [44] J. S. Breese, D. Heckerman, and C. Kadie, “Empirical analysis of predictive algorithms for collaborative filtering”, In UAI, pp. 43-52, 2004.
- [45] Resnick, P., Iacovou, N., Sushak, M., Bergstrom, P., and Riedl, J. , “GroupLens: An open architecture for collaborative filtering of netnews”, Proceedings of the 1994 Computer Supported Collaborative Work Conference.
- [46] G. Morohash, et.al, “Secure Multiparty Computation for Comparator Networks”, IEICE Trans. Fundamentals, Vol. E91-A, No. 9,2008.
- [47] Katzenbeisser, S. and Petkovic, “Privacy-Preserving Recommendation Systems for Consumer Healthcare Services”, In Proceedings of the 2008 Third international Conference on Availability, Reliability and Security (ARES 2008), IEEE Computer Society, pp. 889-895, 2008.
- [48] Ahmad, W. and Khokhar, “An Architecture for Privacy Preserving Collaborative Filtering on Web Portals”, In Proceedings of the Third international Symposium on information Assurance and Security, IEEE Computer Society, pp. 273-278, 2007.
- [49] J. S. Breese, D. Heckrman, and C. Kadie, “Empirical analysis of predictive algorithms for collaborative filtering,” In UAI, pp.43-52, 2004.

-
- [50] H. Kikuchi, H. Kizawa and M. Tada, “Privacy-Preserving Collaborative Filtering Schemes”, WAIS 2009, ARES 2009 federated workshop, IEEE Press, 2009.
- [51] B. Sarwar, G. Karypis, J. Konstan, and J. Riedl “Item-Based Collaborative Filtering Recommendation Algorithms,” ACM WWW10, Hong Kong, May 2001.
- [52] 株式会社ドリーム・トレイン・インターネット, OpenBit.Net, “インターネット接続サービス利用規約”, 2011年4月1日.
- [53] G. Jagannathan, K. Pillaipakkamnatt, R. N. Wright, D. Umamo, “Communication-Efficient Privacy-Preserving Clustering”, Transaction on Data Privacy, pp. 1-25, vol. 3, 2010.
- [54] W. Kowalczyk and N. Vlassis, “Newscast EM”, Advances in Neural Information Processing Systems 17, MIT Press, 2005.
- [55] 佐久間 淳, 小林 重信, “P2P ネットワークにおけるプライバシーを保護した非同期平均計算プロトコル”, pp. 1-6, SCIS2007 3D4-1.
- [56] 佐久間 淳, 小林 重信, “P2P ネットワークにおけるプライバシーを保護した k -means クラスタリング”, pp. 1-6, SCIS2007 3D4-2.
- [57] 佐久間 淳, 荒井 ひろみ, “オンライン予測におけるプライバシー保護” 信学技報, IBISML, pp. 49-56, 2010.

謝辞

本論文を執筆するにあたり多くの方々から多大なる御指導と御援助を賜りました。

特に、研究に関わらず私を導いて下さった東海大学情報理工学部情報メディア学科菊池 浩明 教授に深く感謝を申し上げます。

また、本研究を推進するにあたって、御親切なる御教示ならびに御激励を賜りました東海大学情報理工学部情報科学科 中西 祥八郎 教授，東海大学情報理工学部情報科学科 内田 理准教授，に厚く御礼申し上げます。

2年間共に楽しみ，苦しみ，励まし合い，時には研究に対して有益な意見を与えてくれた東海大学大学院工学研究科情報理工学専攻の皆様，先生方に感謝の意を述べると共に，謝辞とさせていただきます。

最後に，家族に心より感謝致します。

付録A ベイズ推定による摂動化アルゴリズム

A.1 ベイズ推定による摂動化アルゴリズムの実装

A.1.1 摂動化の概要

Agrawal and R. Srikant によって [1] によって発表された最初の摂動化アルゴリズムである。個人情報に意図的にランダムノイズを乗せて、仮想化サーバに格納されたデータのプライバシーを保護しようとするものである。

真の値を x_1, x_2, \dots, x_n とする。ここに加える確率分布 Y の乱数を y_1, y_2, \dots, y_n とする時、再構築問題 (Reconstruction Problem) とは、 $w_1 = x_1 + y_1, w_2 = x_2 + y_2, \dots, w_n = x_n + y_n$ の値と確率変数 Y から、真の値 X の確率分布を見積もることである。例えば、年齢が $x = 20$ 代であるという個人情報をそのまま渡す代わりに、一様分布 (またはガウス分布) の乱数 r を加え、 $x + r = 30$ のように歪んだ値を登録する。30 という属性値を持った顧客がいても、本当に 30 代なのか乱数で 40 代から歪まされたのか、第三者には区別がつかない。

暗号化による方法と異なり、時間のかかる暗号化はなく、計算も各パーティで独立に計算できる。通信効率も計算効率も高い。大規模なデータベースにおいても適用可能である。

摂動化されたデータ Y と摂動化に用いた乱数の分布を表す確率密度関数 f_Y を用いると、ベイズの定理により事後密度関数を次のアルゴリズム 1 により逐次的に求めることが出来る。

Algorithm 1 再構築法 [1]

Input: 確率密度関数 f_Y

1. $f_X^0 = \text{Uniform distribution}$

2. $j = 0$

3. repeat

$$f_X^{j+1}(a) = \frac{1}{n} \sum_{i=1}^n \frac{f_Y(w_i - a) f_X^j(a)}{\int_{-\infty}^{\infty} f_Y(w_i - z) f_X^j(z) dz}$$

4. $j = j + 1$. until stopping criterion met.

A.1.2 再構築アルゴリズムの原理

簡単な数値例を用いて再構築アルゴリズムの原理を示す．確率変数 A が次の分布に従って与えられているとする．

表 A.1: 真の確率分布 $P(A)$

A	0	1	2	3
$P(A)$	0.1	0.3	0.1	0.5

ここで， A の分布を秘匿する為に，表 A.2 の条件付確率 $P(B|A)$ に従って， A の値を変化（摂動化）させた結果を B とおく．ここで，維持確率 $p = 0.4$ は， A を変化させない確率の

表 A.2: 条件付確率 $P(B|A)$, 維持確率 $p = 0.4$

$B \setminus A$	0	1	2	3
0	0.4	0.2	0.2	0.2
1	0.2	0.4	0.2	0.2
2	0.2	0.2	0.4	0.2
3	0.2	0.2	0.2	0.4

大きさであり，変化させるときは一律な確率で分布させることにする．こうして摂動化した結果が，表 A.3 で与えられたとする．オリジナルの分布では $A = 3$ が最頻度で生じていたのに対して，値の差が小さくなりどの値も同じくらい確からしい．

表 A.3: 摂動化した確率分布 $P(B)$

A	0	1	2	3
$P(A)$	0.1	0.3	0.1	0.5

再構築アルゴリズムは，この $P(B|A)$ と摂動化後の確率分布 $P(B)$ だけを与えて，オリジナルの分布 $P(A)$ を近似することを目的とする．アルゴリズム 1 に従い，初期値 $P^0(A) = P(B)$ で与える $i = 1$ において，

$$\begin{aligned} P(A|B) &= \frac{P(B|A)P(A)}{P(B)} = \frac{P(B|A)P(A)}{\sum_{a \in A} P(B|a)P(a)} \\ &= \frac{P(B|A)P(A)}{P(B|A=0)P^0(A=0) + P(B|A=1)P^0(A=1) + \dots + P(B|A=3)P^0(A=3)} \end{aligned}$$

と近似され，この値を用いて A の事後確率の第一近似値は

$$P^1(A) = \sum_{b \in B} P(A|B=b)P(B=b)$$

で与えられる．この数値例の場合の第二近似値までの結果を表 A.4 で示す．徐々に真の分布へ近づいていることが分かる．

表 A.4: 再構築された確率分布の第一近似 $P^1(A)$ と第二近似 $P^2(A)$

A	0	1	2	3
$P^1(A)$	0.22	0.26	0.22	0.31
$P^2(A)$	0.21	0.26	0.21	0.33

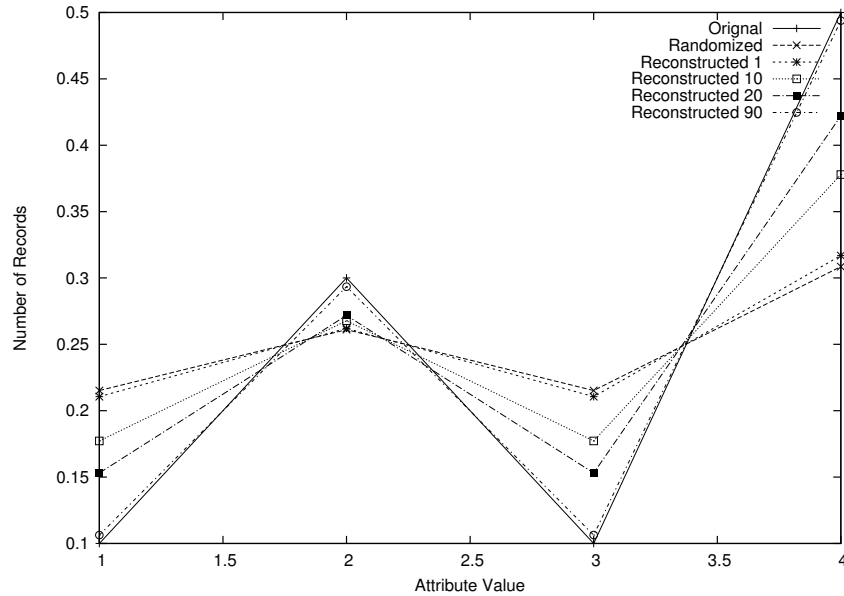


図 A.1: ノイズが除去されて再構築された確率分布 (維持確率 0.4, オリジナルの確率分布 1:3:1:5)

A.1.3 実験結果

摂動化と再構築

再構築を実行するプログラムを実装し, いくつかの数値例や公開データセットについて適用した結果を報告する.

図 A.1 に確率変数 A に維持確率 $p = 0.4$ で摂動化した結果 B と再構築の様子を示す. オリジナルの確率変数 A は, $P(A = 1) = 0.1, P(A = 2) = 0.3, P(A = 3) = 0.1, P(A = 4) = 0.5$ で分布している. 点線で示される摂動化 (randomized) の分布は, ほぼ一様であり, 十分な秘匿が行われている. 再構築による逐次処理を $j = 10, 20, 90$ 回繰返した結果を示しており, オリジナルの実線の分布に近づいている様子が示されている.

オリジナルの分布を変えてもほぼ同様に再構築が可能である. 図 A.2 は A の分布を 1:3:1:5 とした結果, 図 A.3 は A の分布を 1:2:1:3 とした結果であり, いずれも十分に近似されていく様子を示している. 複雑な分布にしても再構築は可能である.

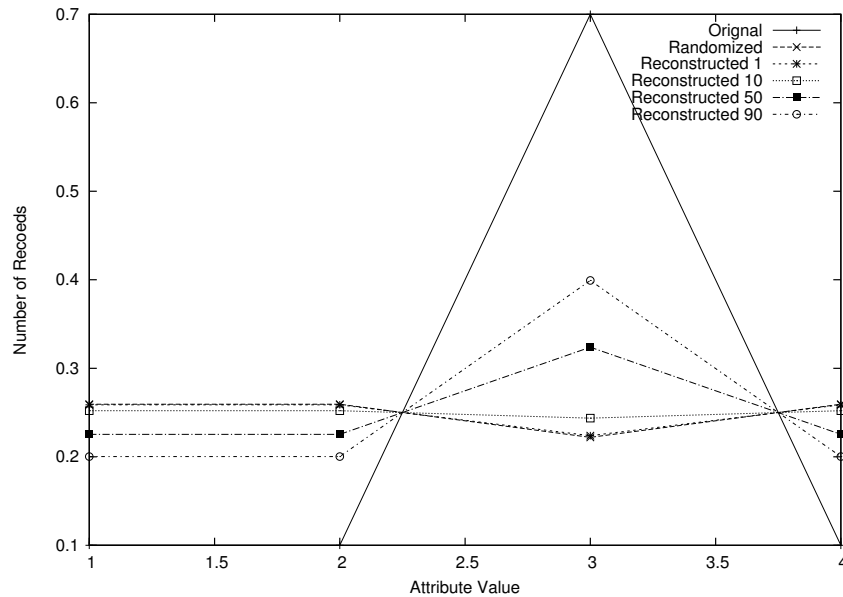


図 A.2: ノイズが除去されて再構築された確率分布 (オリジナルのデータの分布 1:3:1:5)

誤差の収束

データの定義域を2値 $\{1, 2\}$ に絞ると漸近していく様子が分かりやすい。図 A.4 は、2値に限定して、維持確率 $p = 0.13$ とした時の再構築結果を示している。4値の場合と比べて、早く収束していることが分かる。

そこで、近似の繰り返し回数 j について、誤差がどのように変化するかを調べた。図 A.5 はその結果である。オリジナルの確率分布は二値であり、 $P(A = 0) = 0.4, P(A = 1) = 0.6$ の時に、摂動化する維持確率の大きさを $p = 0.3, 0.6, 0.9$ の3通りに変化させて収束を調べている。維持確率 $p = 0.9$ の時が一番収束が早いのは当然として、 $p = 0.6$ の時のほうが $p = 0.3$ の時よりも早く収束する様子は興味深い。

収束に対する維持確率の影響を見るために、同様に4:6の比率で生起する確率変数 A に対して維持確率を変化させて正しい分布に対する絶対誤差を求めた結果を図 A.6 に示す。摂動化の誤差は、維持確率 $p = 0.3$ をピークにして p に対して単調に減少している (Randomized)。

再構築の逐次処理を繰り返すたびに、平均誤差の大きさは小さくなっていく (Reconstructed 10, ..., 50)。しかし、維持確率 $p = 0.5$ の時には誤差が全く減少していない。

維持確率 $p = 0.5$ とは、 $1/2$ の確率で値 0 と 1 とが入れ替わるので、40%の0の半分が1になり、60%の1の半分が0になるので、結局、 $P(B = 0) = P(B = 1) = 0.5$ で等確率になる。従って、ベイズの定理により、

$$\begin{aligned}
 P(A|B) &= \frac{P(B|A)P(A)}{P(B|A=0)P(A=0) + P(B|A=1)P(A=1)} \\
 &= \frac{P(B|A)/2}{P(B|A=0)/2 + P(B|A=1)/2} = \frac{P(B|A)}{p + 1 - p} = P(B|A)
 \end{aligned}$$

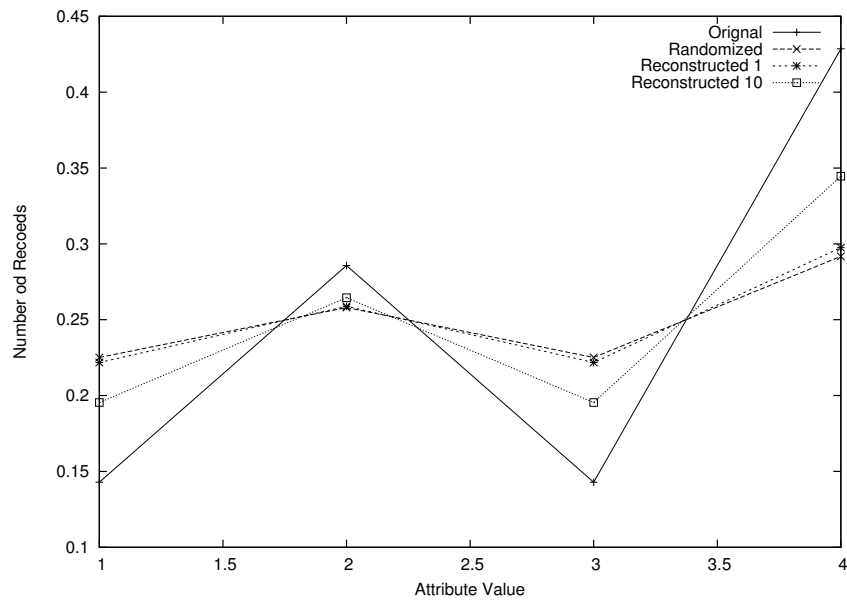


図 A.3: ノイズが除去されて再構築された確率分布 (オリジナルのデータ分布 1:2:1:3)

となり,

$$\begin{aligned}
 P^1(A) &= P(A|B=0)P(B=0) + P(A|B=1)P(B=1) \\
 &= P(B=0|A)\frac{1}{2} + P(B=1|A)\frac{1}{2} = \frac{p}{2} + \frac{1-p}{2} \\
 &= 1/2 = P^0(A)
 \end{aligned}$$

既に収束していることが示された。この様に、与えられた確率分布と維持確率によっては、真の値でないところで局所的な偽収束を引き起こしてしまうことがある。2値の場合に限らず、多値や連続値の場合でもこの可能性は否定出来ず、再構築を繰り返しても誤差を0にすることは困難な例があり得る。

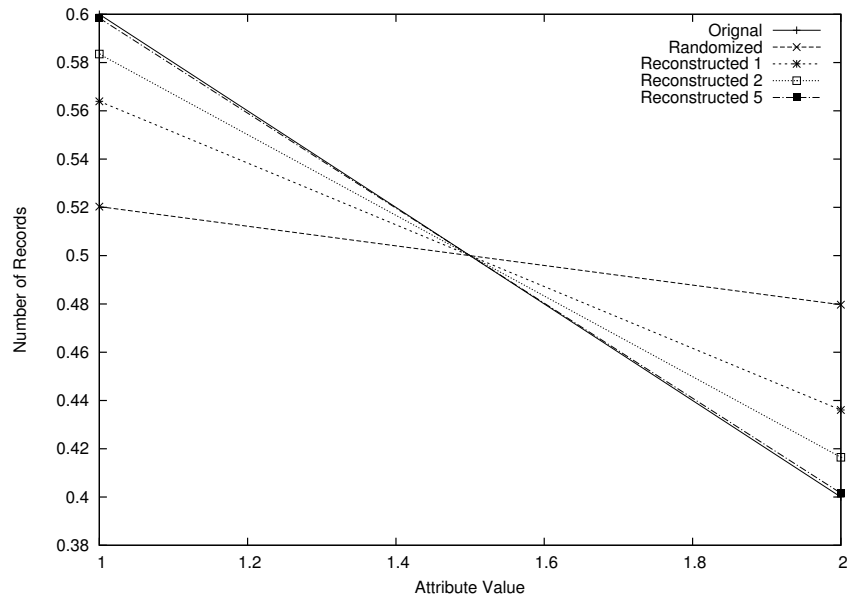


図 A.4: ノイズが除去されて再構築された確率分布 (2 値, 維持確率 0.13)

MovieLens データセットに対する摂動化実験

摂動化の有効性を検討するために, 映画に関する公開データセット MovieLense の 10,000 個の評価値に摂動化と再構築を試みる.

図 A.7 は摂動化を行った評価値の分布を表している. MovieLense データセットでは, 欠損値 0, 1 から 5 までの 5 段階の離散値を用いており, 評価値 4 が最も頻度が高い. 図より, 摂動化によって評価値の分布がより一様分布に近づいており, 300 回までの再構築を繰り返すことによりほぼオリジナルの分布を近似していることが示されている.

図 A.8 は再構築による誤差の減少を表している. 20 回までの近似回数について, 維持確率 $p = 0.4, 0.6, 0.8$ の 3 種類の平均絶対誤差 MAE の大きさを図している. 維持確率が小さければ小さいほど, 収束までに多くの繰り返しが必要なことが分かる.

A.1.4 結論

Agrawal らによって提案されている再構築アルゴリズムを実装し, 公開データセットなどのデータに適用して動作を検証した. 処理は高速に行うことが可能であり, 多くの場合はオリジナルの分布を十分に近似する確率分布を再構築できることが示された. ただし, 特殊な分布や摂動化で行うパラメータによっては収束するまでに多くの時間を要するケースや誤差が生じるケースがあることが分かってきた.

今後は, 処理時間や精度について包括的な実験を行い, 誤差を生じさせる条件などを明ら

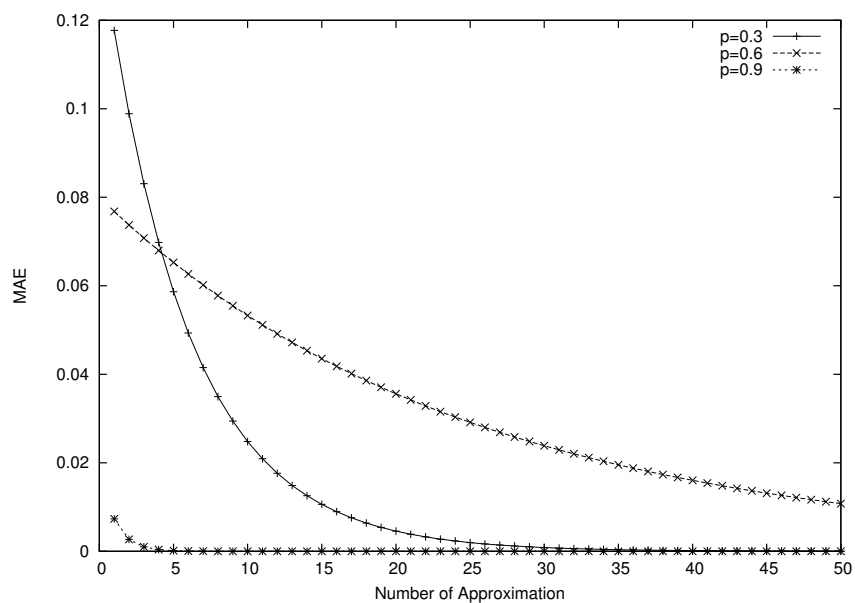


図 A.5: 近似回数に対しての誤差の変化

かにすることを試みる．実用的な摂動化の大きさやパラメータを明らかにし，大規模な仮想化環境のプライバシーを安全に保証する方式を目指す予定である．

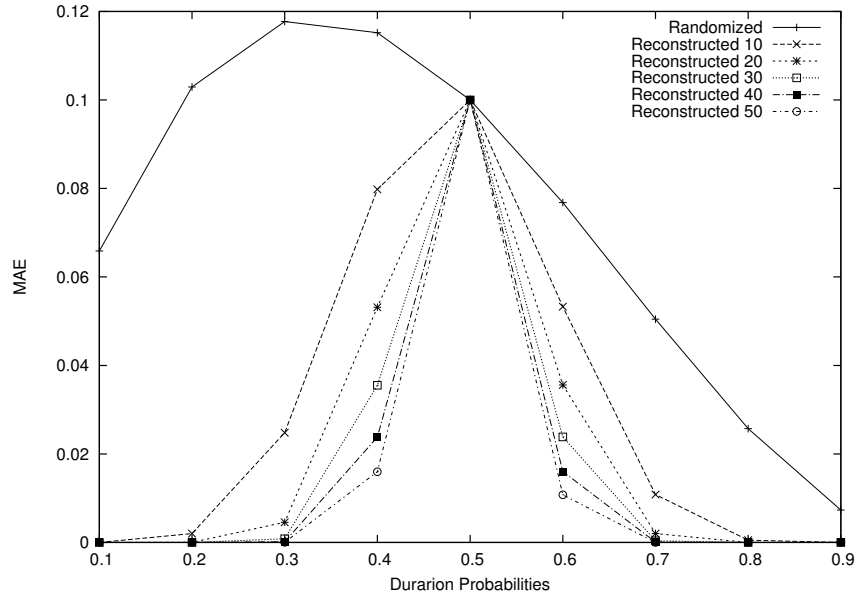


図 A.6: 維持確率 p に対しての誤差の分布 (オリジナルの分布 4:6)

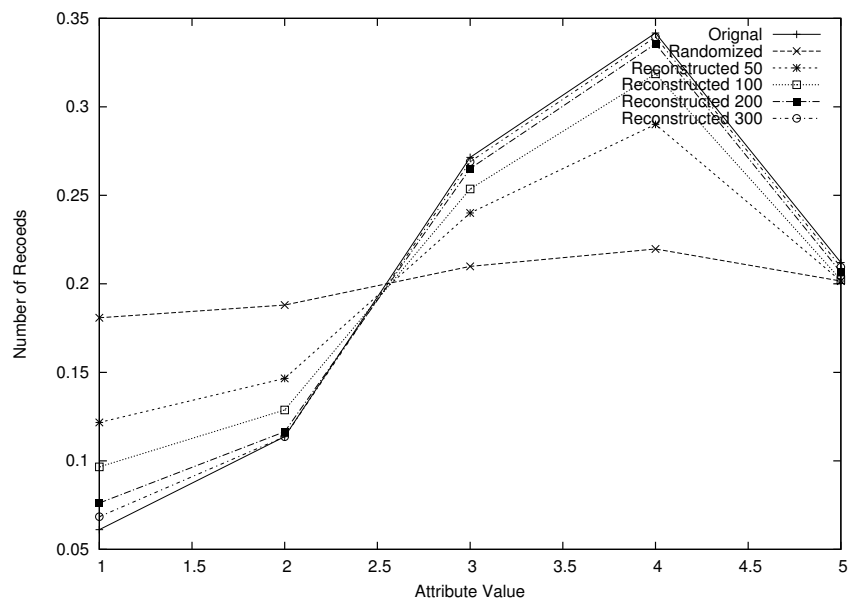


図 A.7: MovieLens Dataset に摂動化した分布と再構築結果

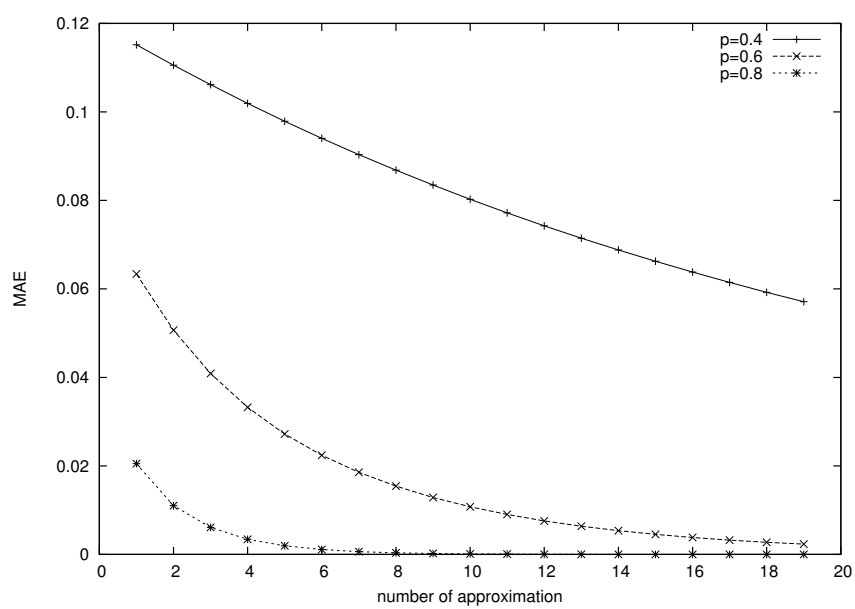


図 A.8: 協調フィルタリングの誤差