
手軽に実現！！ビッグデータ分析

～新製品Big Data Module～

(株)数理システム
bigdata-info@msi.co.jp



目次

デモンストレーション

製品紹介

Big Data Module 応用事例紹介

質疑応答



－ デモンストレーション －

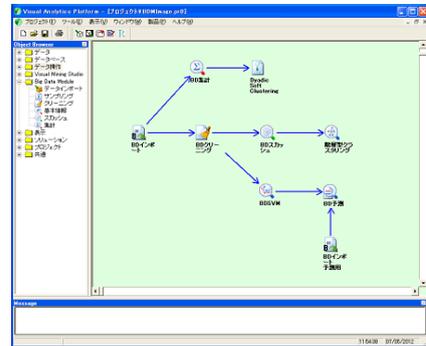
ECサイトのログを例に

- ・ インポート
- ・ 集計
- ・ オンラインロジスティック回帰
- ・ VMSアイコンへの接続

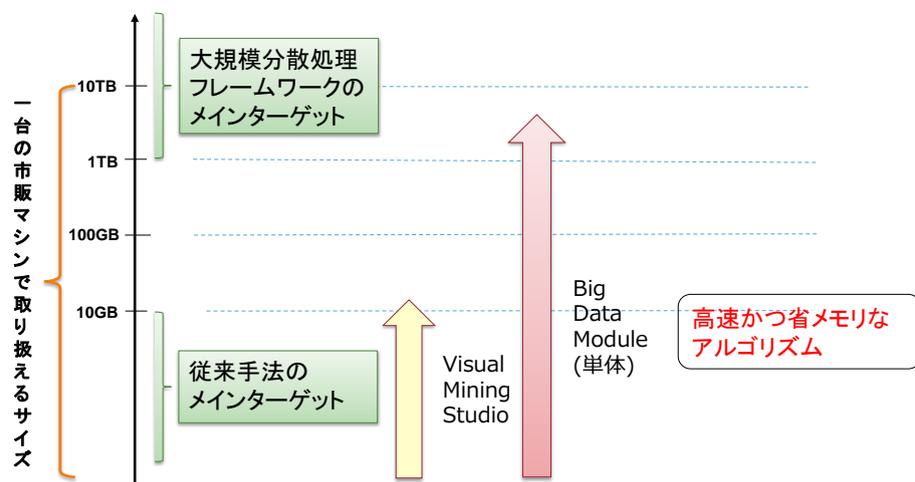
－ 製品紹介 －

Big Data Moduleとは

- 数理システム自社開発
- 簡単な操作で高度なビッグデータ分析が可能
- 特殊な分析専用マシンは不要
市販のマシンを1台用意すれば
それだけで分析が実行可能



データサイズ



データサイズ

アンケートデータ

日本全国 1 億人に、100 の質問
 $100M * 100 * 4B = 40GB (=0.04TB)$

買い物リスト(一年分)

1 億人が、一日 = 1 回、一度 = 10 品
 $100M * 365 * 10 * 4B = 1.5TB$

品質管理センサーデータ (100 台、秒単位、1 年分)

$100 * 365 * 24 * 60 * 60 * 8B = 25GB (=0.025TB)$

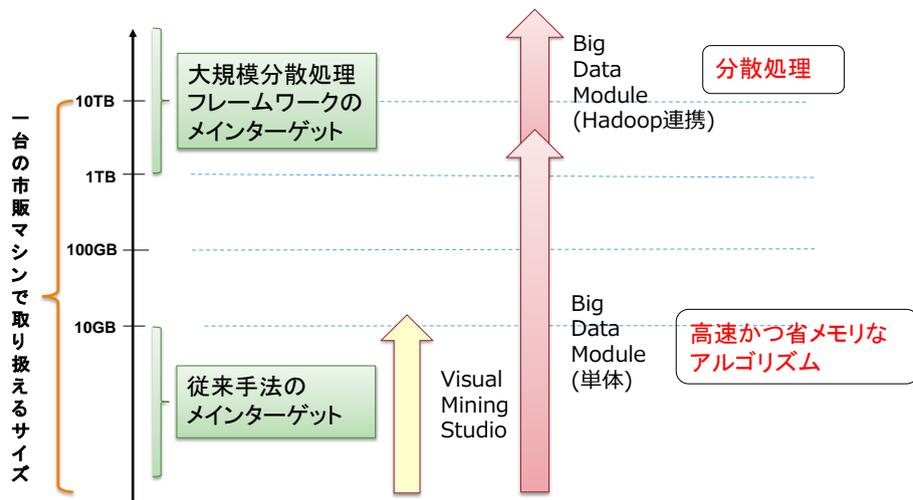
企業評判分析、Web 書込み (1 ヵ月分)

100 万人、1 日 = 10 Tweets、1 Tweet = 140 文字
 $1M * 30 * 10 * 140 * 2B = 84GB$

.....



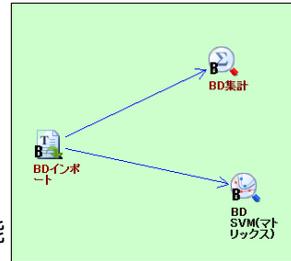
データサイズ



Big Data Module の主な特長

- GUIベースの簡単操作

アイコンを矢印でつなぐだけの簡単操作で、高度な分析機能を使用可能



- 並列処理

並列処理をアイコンベースで簡単に実行可能
マルチコアをフルに生かして分析を実行

- Visual Mining Studio との連携

Visual Mining Studio の機能と組み合わせることで、Big Data Module の分析結果をさらに掘り下げて、多様な分析を行うことが可能



Big Data Module の主な特長

- ビッグデータ分析アルゴリズム

オンラインアルゴリズムなどの
ビッグデータに適した分析アルゴリズムを搭載

- Hadoop との連携

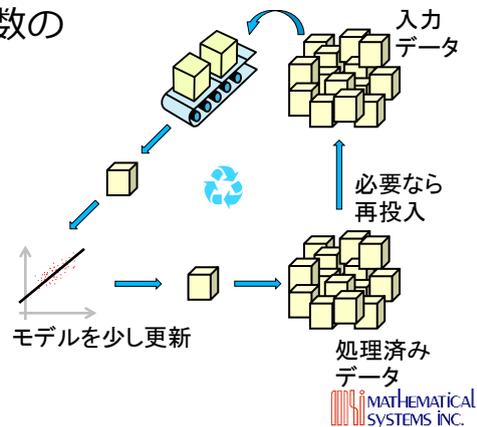


Hadoop上でのデータハンドリングをGUIから簡単に実行
結果を取得してVAPの分析と接続可能

オンラインアルゴリズム

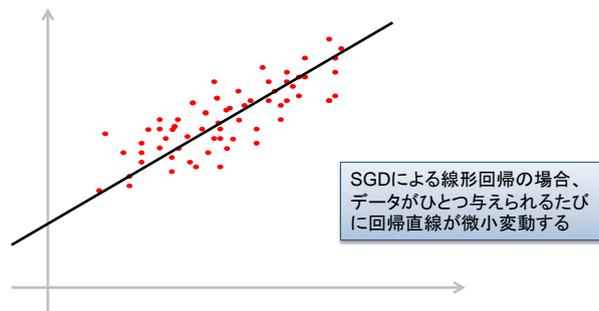
- データを一つずつ読み込み、モデルを逐次更新
- データをためず、**必要最低限のメモリ使用量**
- 計算時間は処理データ数の**線形オーダー**

※「オンライン」はチューリングマシンの理論に由来し、「ネットワークにつなぐ」という意味ではありません。

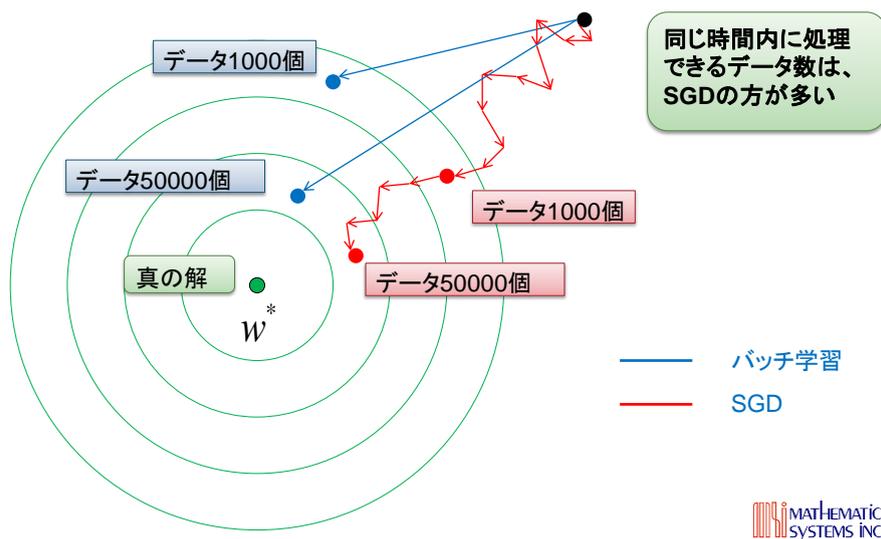


オンラインアルゴリズムの基本エンジンSGD

- **Stochastic Gradient Descent**
 - 確率最適化問題に対するオンラインな計算手法
 - 線形回帰・ロジスティック回帰・SVM・k-means・行列分解...
 - Big Data Module には、**最先端のSGD実装**を搭載

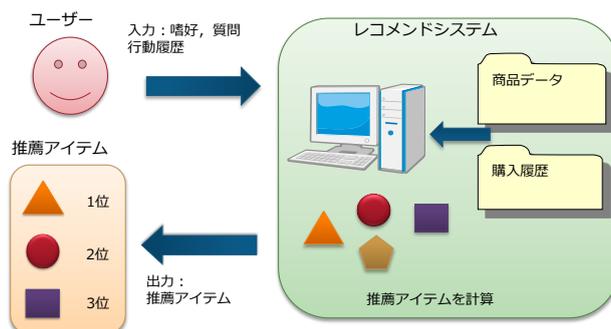


SGDの処理速度



SGDの威力

- Netflix Prize(2006/10~2009/7)
 - レコメンデーションの精度コンテスト
 - 1億件のレンタル履歴データ
 - 100万ドルの賞金
 - SGDを駆使したモデルが優勝



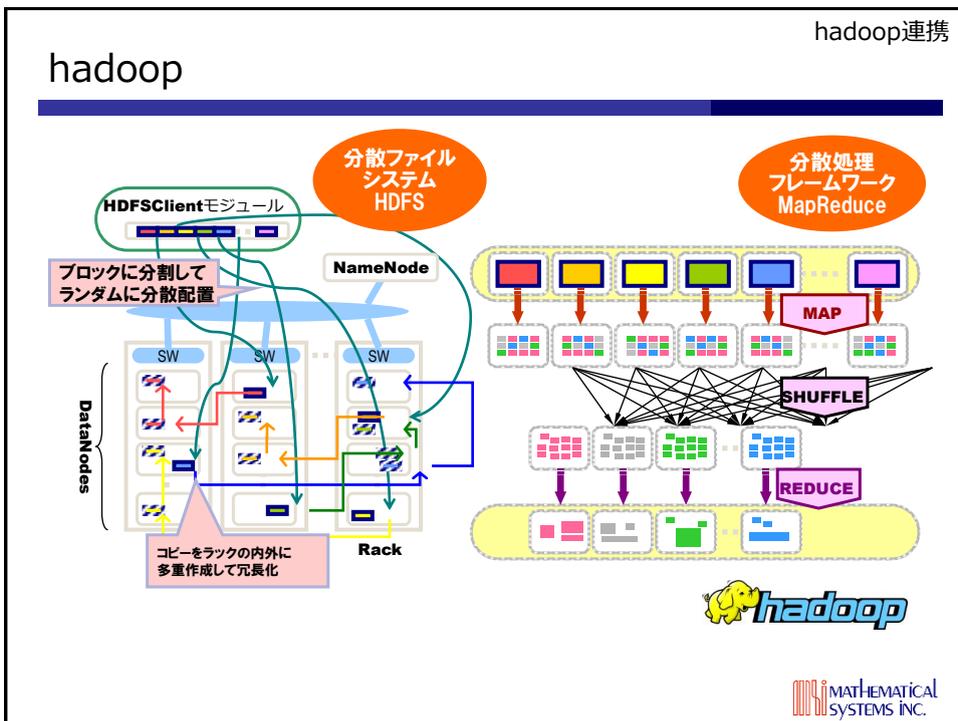
SGDの威力

- Netflix Prize(2006/10~2009/7)
 - レコメンデーションの精度コンテスト
 - 1億件のレンタル履歴データ
 - 100万ト
 - SGDを馬



Big Data Module の主な特長

- ビッグデータ分析アルゴリズム
 - オンラインアルゴリズムなどの
 - ビッグデータに適した分析アルゴリズムを搭載
- Hadoop との連携
 -  Hadoop上でのデータハンドリングをGUIから簡単に実行
 - 結果を取得してVAPの分析と接続可能



hadoop連携

hadoop連携イメージ

The screenshot shows the **Visual Analytics Platform** interface. A red box highlights the text: **バックグラウンドでhadoopが動作** (Hadoop operates in the background). A blue box highlights the text: **hadoopデータを取り込み Big Data Moduleで分析** (Load Hadoop data and analyze with Big Data Module). A yellow box highlights the text: **hadoopの処理結果を確認** (Check Hadoop processing results). The interface includes an **Object Browser**, a **DataGraph Viewer** showing a table of data, and a **DataGraph** with nodes for **Hadoopデータリンク**, **BDアイコン**, **Hadoopデータインポート (Hadoop→BDM)**, and **Hadoopデータエクスポート (BDM→Hadoop)**. A **Hadoopデータエクスポート (Text→Hadoop)** node is also present. The **MATHEMATICAL SYSTEMS INC.** logo is at the bottom right.

Setosa	Versicolour	Versicolour	Versicolour
1	51	35	14
2	49	30	14
3	47	32	13
4	46	31	15
5	50	36	14
6	50	39	17
7	46	34	14
8	50	34	15
9	44	29	14
10	49	31	14

Big Data Module の機能一覧

- 基本機能
 - インポート
 - クリーニング
 - サンプルング
 - 基本情報
 - 集計
 - データハンドリング **新機能**
- マイニング機能
 - オンライン線形回帰
 - オンラインロジスティック回帰
 - サポートベクターマシン(SVM)
 - スカッシング
 - オンライン K-means **新機能**
 - オンライン行列分解 **新機能**
- その他
 - hadoop連携 **新機能**



— Big Data Module 応用事例紹介 —

予測モデル作成
レコメンデーション

予測モデル作成



- 売上予測
- 優良顧客判別



- 株価予測
- 与信管理



- 電力需要予測
- スマートグリッド

- 年々積み重なるデータ
- 技術の進歩で取得できる情報の種類も量も増加
→ ビッグデータ化
- データが豊富になったんだから予測精度も上がるはず

MATHEMATICAL
SYSTEMS INC.

大規模データを用いた未来予測

従来の予測モデル

- 従来の学習アルゴリズムをビッグデータに適用しようとする…
データ数、説明変数の数が多すぎる！
メモリ使用量の爆発的な増加
計算時間の爆発的な増加
→ **計算不可能**
- 従来の学習アルゴリズムをビッグデータに適用することは
現実的に不可能
- 対応方法
アルゴリズム的な工夫

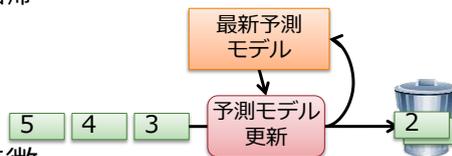


MATHEMATICAL
SYSTEMS INC.

大規模データを用いた未来予測

オンラインマイニングアルゴリズム

- Big Data ModuleはSGDを使用した分析アイコンを搭載
 - オンライン線形回帰
 - オンラインロジスティック回帰
 - オンラインSVM

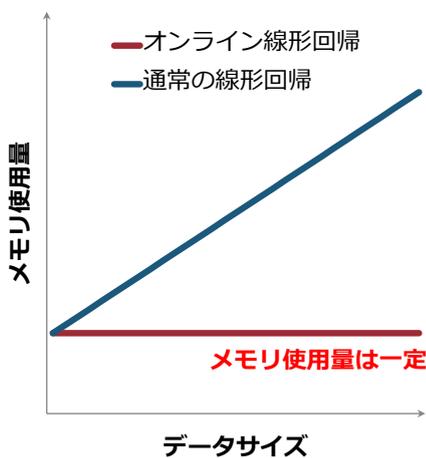


- SGDを使った分析機能の特徴
 - データサイズに依存せず**省メモリ**
 - 計算時間はデータサイズの**線形オーダー**
 - ビッグデータと非常に相性が良い!

大規模データを用いた未来予測

メモリ使用量

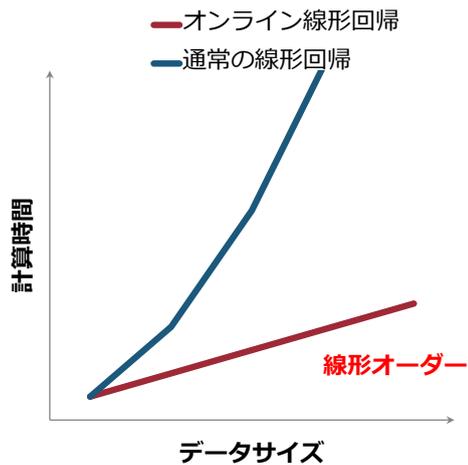
- 通常のアプローチでは、データサイズに応じてメモリ使用量が増加
→ 計算不可能に
- オンライン線形回帰では、データサイズが増加しても、メモリ使用量は一定



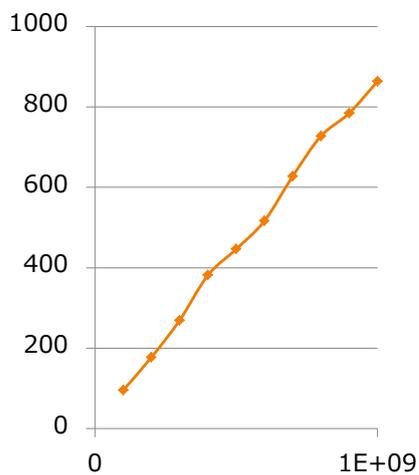
計算時間

- 通常のアプローチでは、データサイズに応じて計算時間が爆発的に増加
- オンライン線形回帰では、データ量の線形オーダーの時間で計算が終わる

$$\text{計算時間} = \alpha \times (\text{データサイズ})$$



オンラインロジスティック回帰の計算時間



行数	計算時間(秒)
1億	95
2億	177
3億	269
4億	381
5億	446
6億	516
7億	627
8億	727
9億	784
10億	863

レコメンデーション

商品、記事、サービスなどのアイテムを推薦する機能。

- ECサイト
amazon、iTunes、
アパレル、オンライン書店等
のおすすめ商品

あなたへのお勧め商品は、...



- SNS サイト
知り合い候補を列挙
- ニュースサイト
関連記事

その他、さまざまな箇所で導入されている技術

MATHEMATICAL
SYSTEMS INC.

協調フィルタリング

レコメンデーション

ユーザーの嗜好データを用いた推薦方法。

1. 嗜好が似たユーザーが高得点を付けたアイテムを推薦。

	映画A	映画B	映画C	映画D	...
ユーザー1	-	5	4	5	
ユーザー4	3	5	4	-	
...					...



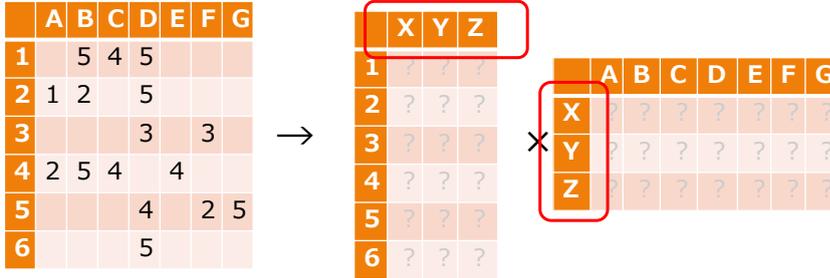
おすすめ!!

- 必要な技術
 1. ユーザーの類似度計算
 2. 類似ユーザーの探索
(コンテンツベースの方法と同様の技術。)

MATHEMATICAL
SYSTEMS INC.

行列分解

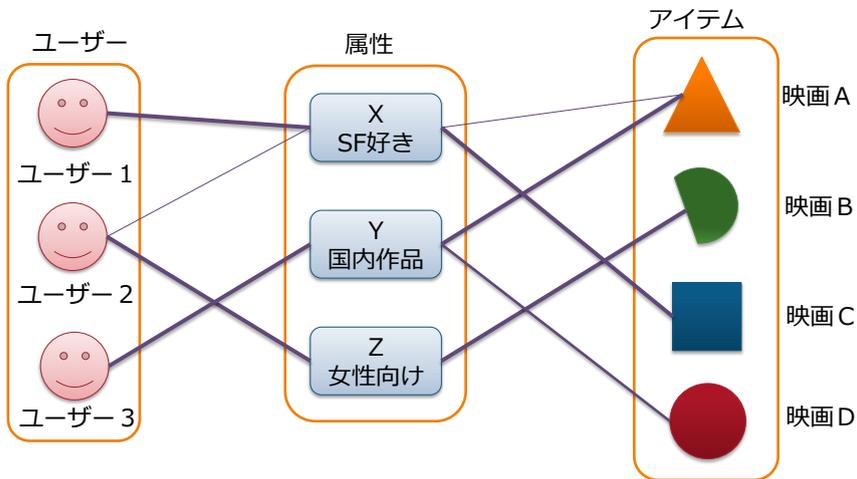
ユーザーと映画の関係を表す行列を
2つの行列の積に分解

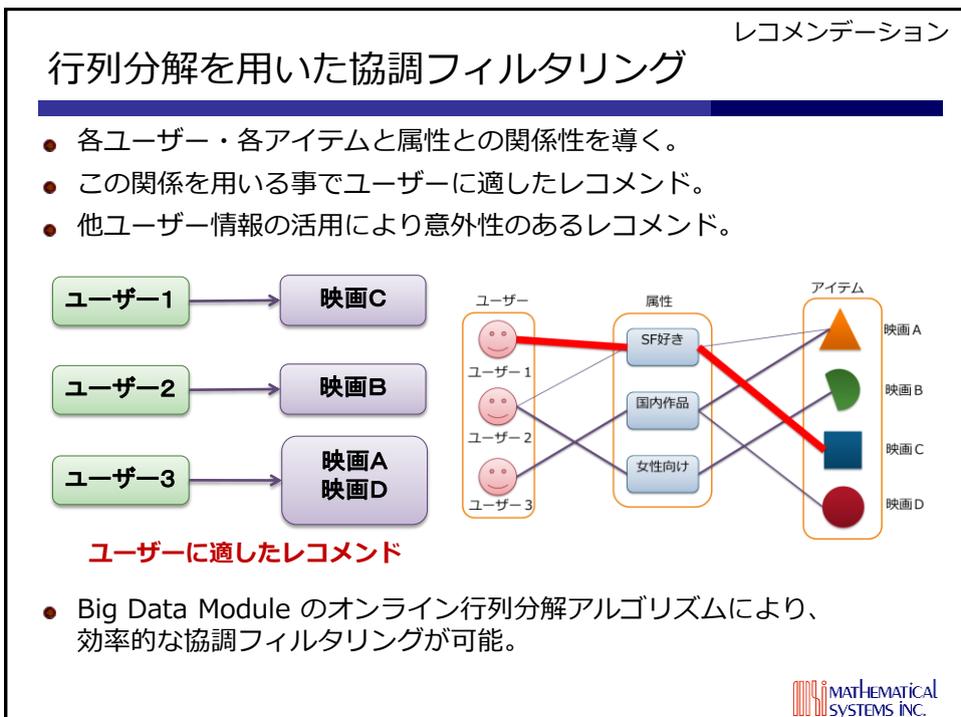
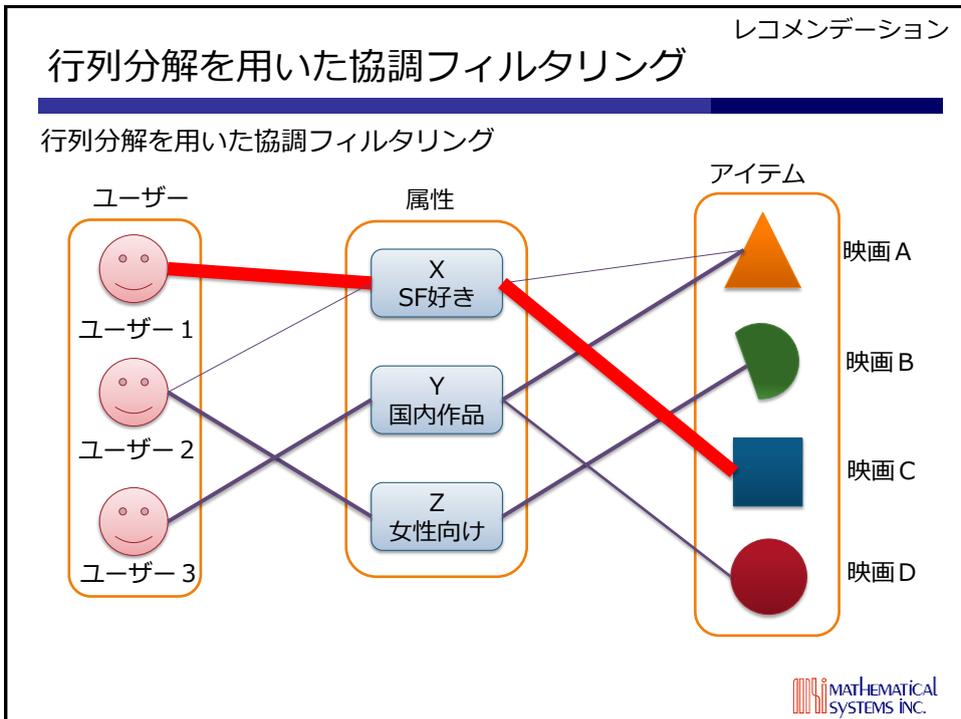


新たに表れたX、Y、Zが属性を表す

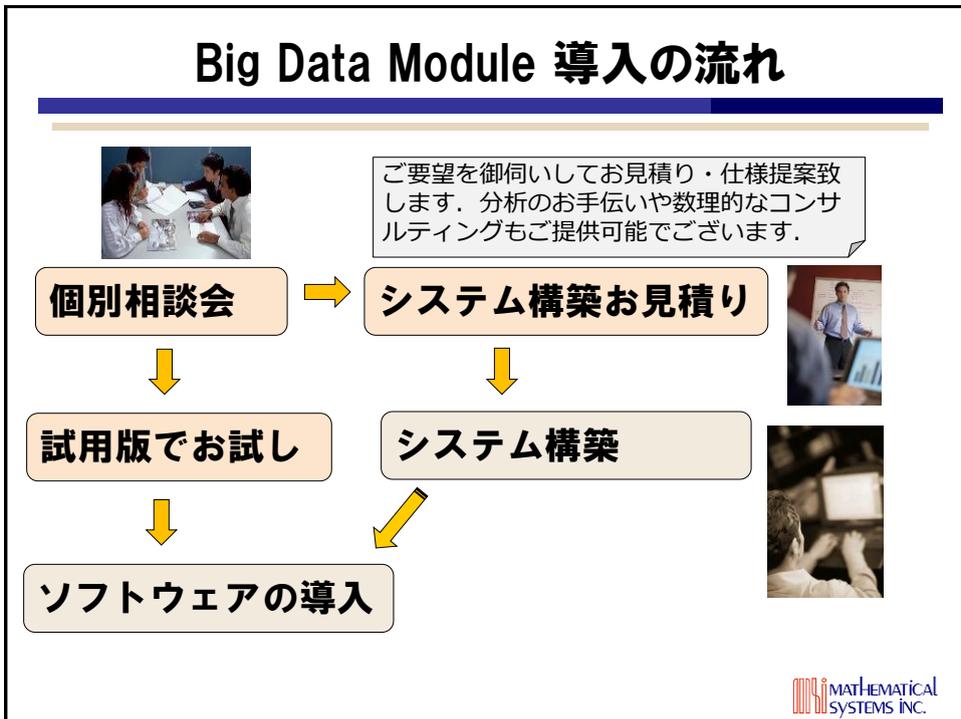
行列分解を用いた協調フィルタリング

行列分解を用いた協調フィルタリング





Big Data Module 導入の流れ



本日はありがとうございました

★数理システムがお手伝いできること

※自社展開／他社へ向けてサービス展開は不問

ご相談・ヒアリング
解決のためのご提案

既存のシステムに分析モジュールの
組み込みご相談

分析のコンサルテーション
のご提供

他社に向けての共同でのご提案

ソフトウェアのご提供

bigdata-info@msi.co.jp

何でもご相談下さい!

MATHEMATICAL SYSTEMS INC.