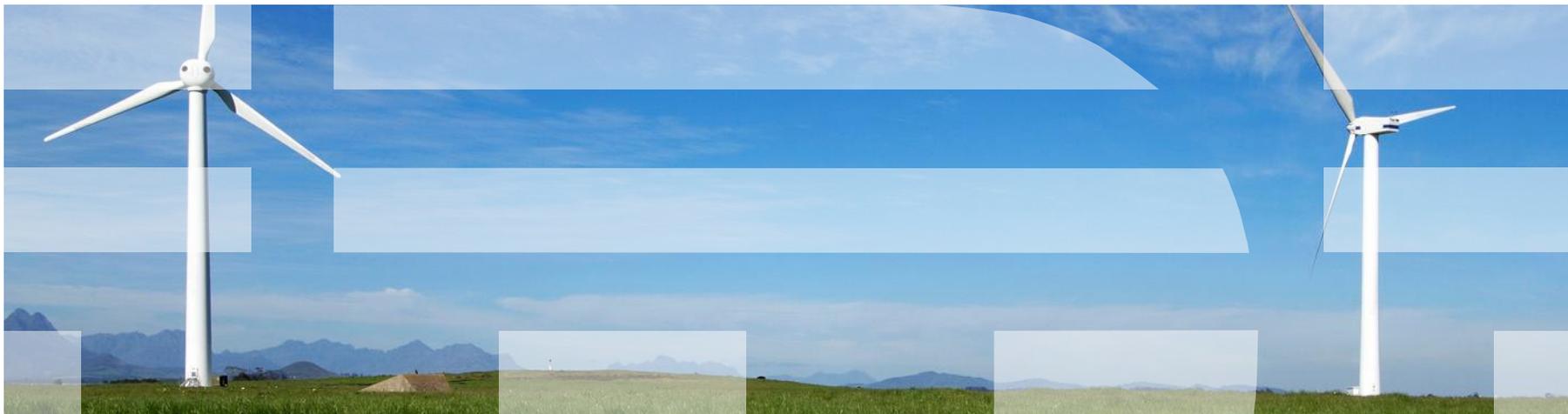


スパース構造学習による異常検知技術

IBM東京基礎研究所 数理科学担当
担当部長 井手 剛



目次

- **イントロダクション**
 - 機械学習/データマイニングと制御工学
 - データマイニングをめぐるIBMのビジネス戦略
 - 私たちが解いている問題の例
- **スパース構造学習による異常検知**
- **まとめ**

目次

- **イントロダクション**
 - 機械学習・データマイニングと制御工学
 - データマイニングをめぐるIBMのビジネス戦略
 - 私たちが解いている問題の例
- **スパース構造学習による異常検知**
- **まとめ**

自己紹介

機械工学 → 物性物理 → 液晶工学 → データマイニング

▪ 学歴

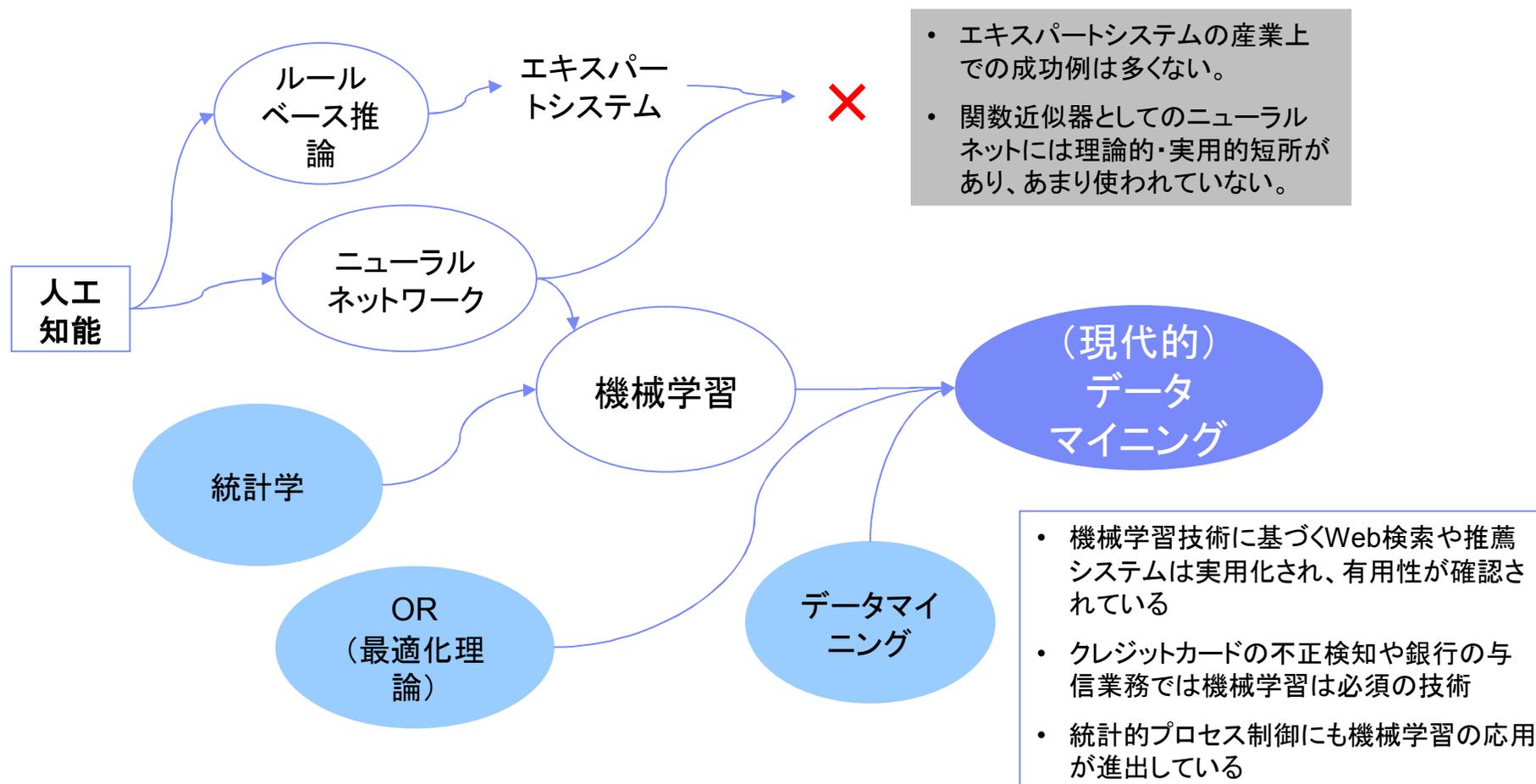
- 北海道の田舎の高専(機械工学科)
- 東北大学工学部(機械工学科、編入)
- 東京大学大学院・理学系研究科(物理学専攻)
 - 物性理論で博士号

▪ 入社後(2000年～)

- アカデミックポストが見つからず、やむを得ずIBM東京基礎研究所入所
- 液晶ディスプレイ
- 自律型コンピューティング
- データ解析(2004～)

▪ 現在は、IBM東京基礎研究の数理科学グループを担当しています

ここ10年の機械学習技術の急速な発展に伴い、その産業応用が進展しています



制御工学との融合が直近の重要なテーマとなると予想されます

計測自動制御学会誌・特集企画

— プラントモデリングの新展開 — 効率的な制御対象モデリングと制御システムの開発を目指して

— 解説 No.9

— 「機械学習技術の最近の発展とシステムモデリングへの応用」

- IBM東京基礎研究所 井手剛,
- 東京大学先端科学技術研究センター 矢入健久

✓ 機械学習は、最近の情報技術の発展を象徴する研究分野のひとつである。本稿では、この10年の機械学習技術の急速な進歩を象徴するキーワードとして、「非線形性」と「スパース性」を取り上げ、これらを軸に関連する手法を概観する。同時に、機械学習技術のシステムモデリングへの応用について最近の研究事例を紹介する。



(ご参考) 機械学習の主要国際会議のひとつICMLでは、部分空間同定法に関する論文がbest paperに選ばれました

▪ **Hilbert Space Embeddings of Hidden Markov Models**

– Le Song, Byron Boots, Sajid Siddiqi, Geoffrey Gordon, Alex Smola

- Hidden Markov Models (HMMs) are important tools for modeling sequence data. However, they are restricted to discrete latent states, and are largely restricted to Gaussian and discrete observations. And, learning algorithms for HMMs have predominantly relied on local search heuristics, with the exception of spectral methods such as those described below. We propose a non parametric HMM that extends traditional HMMs to structured and non-Gaussian continuous distributions. Furthermore, we derive a local-minimum-free kernel spectral algorithm for learning these HMMs. We apply our method to robot vision data, slot car inertial sensor data and audio event classification data, and show that in these applications, embedded HMMs exceed the previous state-of-the-art performance.



International
Conference on
Machine
Learning

Haifa, Israel
June 21 - 24



系のモデリングは制御のための第一歩ですが、モデリング自体に非常に手間がかかる場合があります

- たとえば自動車業界では Matlab/SimuLinkを用いたモデルベース開発が盛んです
- この手法は、比較的単純な物理系では大変有用ですが、たとえば排ガス中のNO_xの制御など、複雑な物理現象を扱うのは困難です
- したがって、物理モデルを、手元の実験データと合わせるための新しい手法が必要です

物理モデルを立てるのが難しいような分野では、機械学習の手法の併用が効果的だと考えられます

「ホワイトボックス」の方法

物理学の原理に基づいてモデルを立てる

$$\frac{\partial \mathcal{L}}{\partial f} - \frac{d}{dx} \left(\frac{\partial \mathcal{L}}{\partial f'} \right) + \frac{d^2}{dx^2} \left(\frac{\partial \mathcal{L}}{\partial f''} \right) - \dots + (-1)^n \frac{d^n}{dx^n} \left(\frac{\partial \mathcal{L}}{\partial f^{(n)}} \right) = 0$$

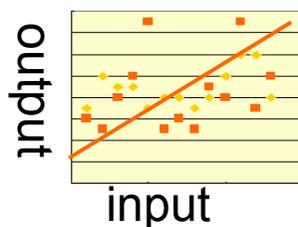
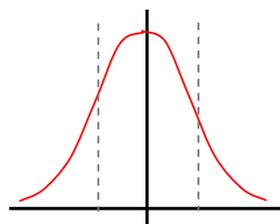
$$\rho \left(\frac{\partial \mathbf{v}}{\partial t} + \mathbf{v} \cdot \nabla \mathbf{v} \right) = -\nabla p + \mu \nabla^2 \mathbf{v} + \left(\frac{1}{3} \mu + \mu^v \right) \nabla (\nabla \cdot \mathbf{v}) + \mathbf{f}$$

「灰色ボックス」モデル

- 実験データをよく説明する
- 物理モデルと矛盾しない

ブラックボックスの方法

実験データに基づいて入力と出力を結ぶ関数関係を学習



機械学習の手法のシステムモデリングへの適用例: ボルテラ級数の同定

- 例: Volterra級数の同定

- 問題: 入力 (\mathbf{u}) と出力 (y) の間の関数関係を、Volterra級数の形で与えること

$$y = h^0 + \sum_{n=1}^{\infty} \sum_{i_1=1}^m \cdots \sum_{i_n=1}^m h_{i_1, \dots, i_n}^n u_{i_1} \cdots u_{i_n}$$

- 観測データ: (入力、出力)のN個の組

$$\mathcal{D} \equiv \left\{ (\mathbf{u}^{(t)}, y^{(t)}) \mid \mathbf{u}^{(t)} \in \mathbb{R}^m, y^{(t)} \in \mathbb{R}, t = 1, 2, \dots, N \right\}$$

- 従来手法は、級数の有限次の打ち切りに基づいていた
- しかしカーネル法を使えばそのような人為的な近似なしにシステムを同定できる

M. O. Franz and B. Schoelkopf. "A unifying view of Wiener and Volterra theory and polynomial kernel regression". *Neural Computation*, 18(12):3097–3118, 2006.

(ご参考)現代の機械学習の2つのキーワード: 非線形性とスパース性

今日の話の中心

非線形性

- 入力と出力の間の非線形的な関係を柔軟かつ簡単に取り込む
- 技術的キーワード
 - カーネル法

スパース性

- 入力と出力の間の関係を、より少ないパラメータや変数で表す
- 技術的キーワード
 - L1正則化、Lasso
 - 関連度自動決定 (ARD: automated relevance determination)

このパートのまとめ

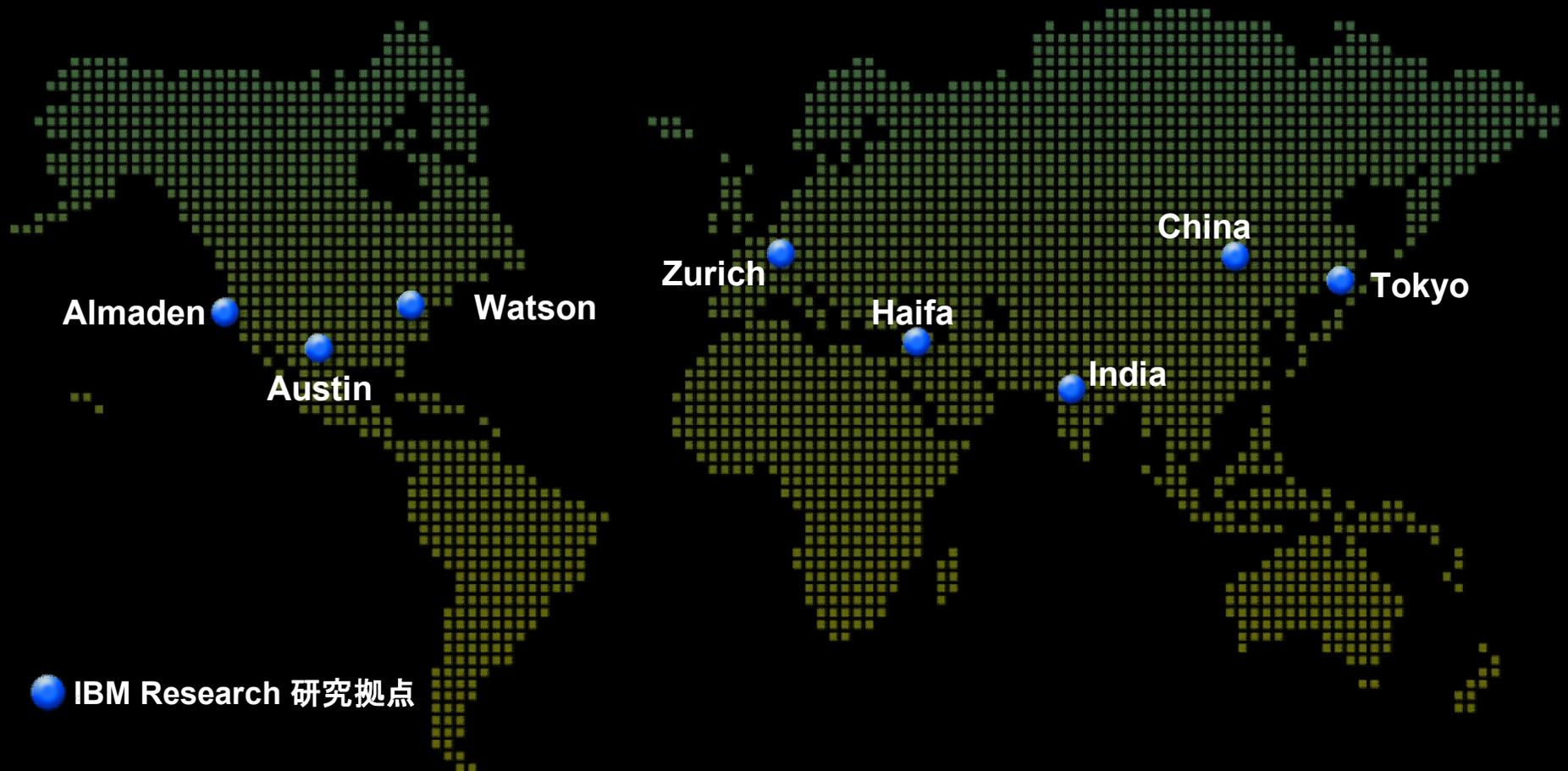
- 機械学習はここ10年で目覚ましい発展を遂げ、いまや社会基盤に直接インパクトを与えうる産業技術になっています
- 機械学習はこれまで多くの他分野と相互作用しながら発展してきましたが、制御工学との連携は、直近の最重要テーマのひとつです

目次

- イン트로ダクション
 - 機械学習/データマイニングと制御工学
 - データマイニングをめぐるIBMのビジネス戦略
 - 私たちが解いている問題の例
- スパース構造学習による異常検知
- まとめ

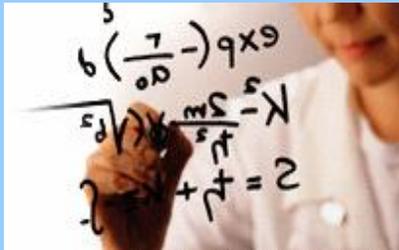
IBM東京基礎研究所は世界に8箇所ある研究拠点のひとつです

全世界に8拠点、約3000名が研究に従事



IBM基礎研究部門のストラテジー・エリア

数理科学



インダストリー・ソリューション



サービス



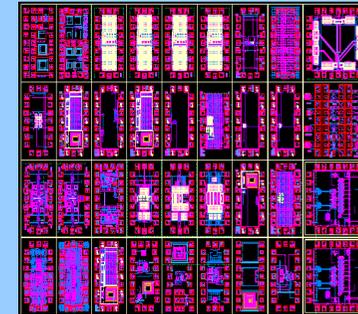
ソフトウェア



システム



テクノロジー



基礎科学研究

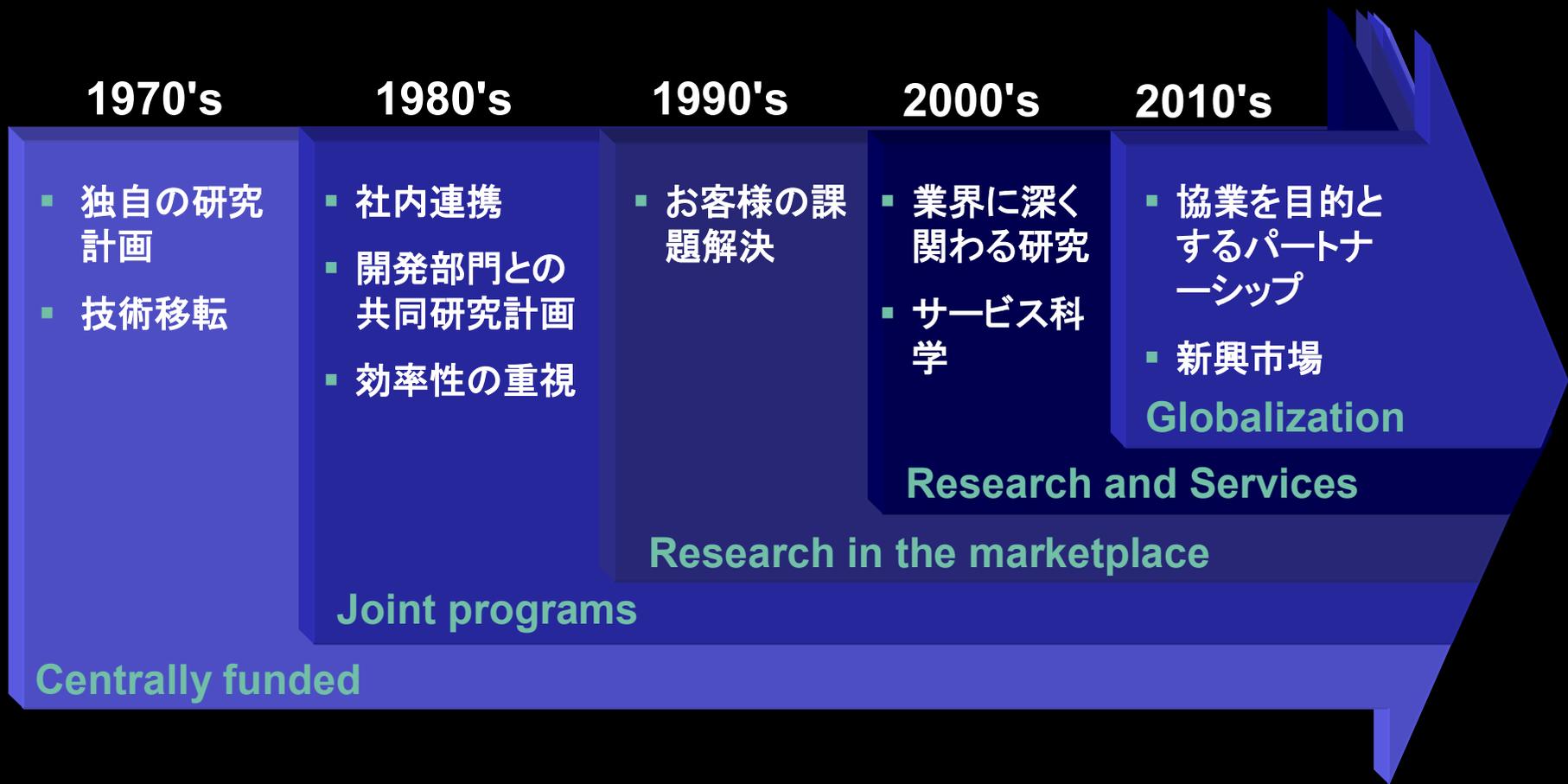


IBM東京基礎研究所の数理科学チームは、「つぶしがきく」チームとして大いに活躍しています



- 実問題を解く上ではデータ解析に関する幅広い知識が必要
 - 機械学習
 - 最適化理論
 - アルゴリズム
 - データ工学
 - ソフトウェア工学
 - 数値計算
- 研究員は、最低何かひとつの分野の専門家であることが期待されている

IBM研究部門は現在、ビジネスと研究の結合を世界で最も成功させている企業研究所です



データマイニングを基盤とするコンサルティングビジネスは IBMの最重点投資領域です

Japan [変更]

検索

ホーム ソリューション サービス 製品 サポート & ダウンロード My IBM

ようこそ [ログイン] [登録]

「BAO」が導く新たな知見。 膨大な情報を経営に生かす秘訣とは？

散在する多様な情報からビジネス・チャンスを発掘。
予測的分析に基づく業務最適化で、さらなる成長を牽引します。

はじめに BAOでビジネス・チャンスを ソリューション お客様導入事例

BAOソリューション体系
ビジネスの課題にきめ細かく対応させ、BAOサービスを利用しやすくしました。
[ソリューション体系図を表示する](#)

情報活用診断
膨大なデータをもっと有効活用したいが、どこから手をつけたいのかというお客様に最適です。情報活用に関するお客様の現状を調査・診断して、あるべき姿と現実のギャップを明らかにし、優先着手領域を特定します。
・情報活用の成熟度診断
・BIアプリケーション診断

クイック・アセスメント・サービス
「ソリューションの本格導入の前に効果を検証したい」というお客様のご要望に対して、短期間に効果を検証します。
・KPI定義による見える化構想支援
・テキストマイニングを活用したお客様の声活用
・異常値の効率的な検知
・統計アプローチに基づく在庫最適
・価格最適化

BAOソリューション

- スマートな経営の可視化**
・ビジネス・パフォーマンスを高める経営の可視化
・企業の活動評価指標の最適化
- スマートな顧客分析**
・お客様の声活用
・インターネット上の顧客行動の分析
・販促費用の最適化
・販売価格の最適化
- スマートな生産と物流**
・配送経路の最適化
・生産計画と製造ラインのスケジューリング
・動的在庫の最適化
・製造現場と経営をつなぐマネージメント・コックピット
- スマートな品質管理・異常検知**
・異常値の検知
・映像・画像の自動分類と検索
- スマートなコンテンツ管理**
・ワンソース・マルチユースを実現するドキュメント作成

フォームでお問い合わせ
まずはお気軽にご相談ください。
[入力フォーム](#)

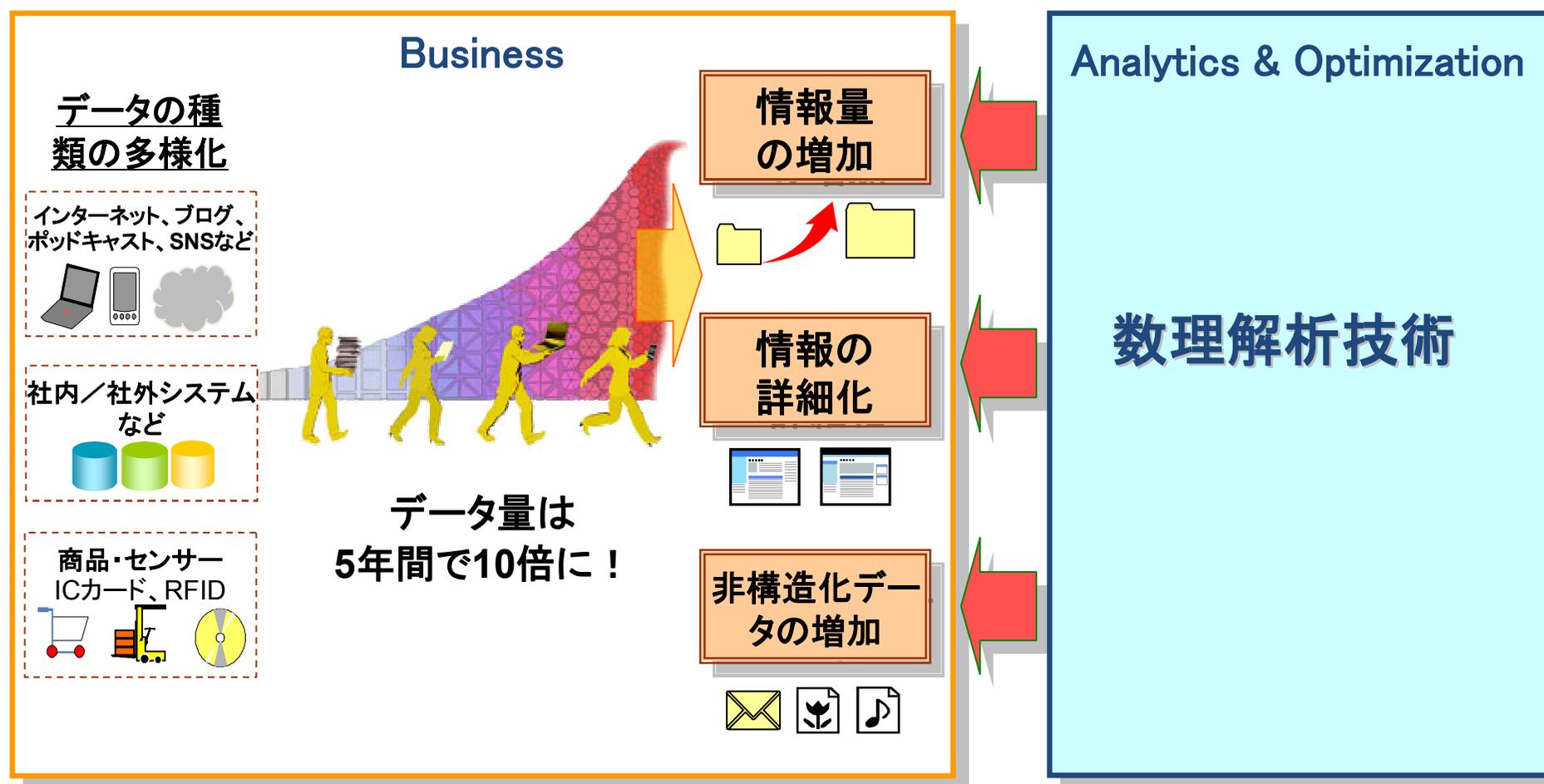
IBM Smart Analytics System

スマートなビジネスを加速する分析アプライアンス
[詳細はこちら](#)

最先端のビジネス・アナリティクス

IBMのBAO (Business Analytics & Optimization) サービスは、数理解析技術を中心としたコンサルティングサービスです

- 誰でもできるような業務の自動化はすでにコモディティ化している
- 高度な解析技術に基づいて、ビジネスに有用な知見を与えることが求められている



IBM東京基礎研究所の数理科学チームは、BAOチームと連携して主にコンサルティングビジネスを展開しています

(ご参考) 数理科学の技術が経営判断に役立つ実例

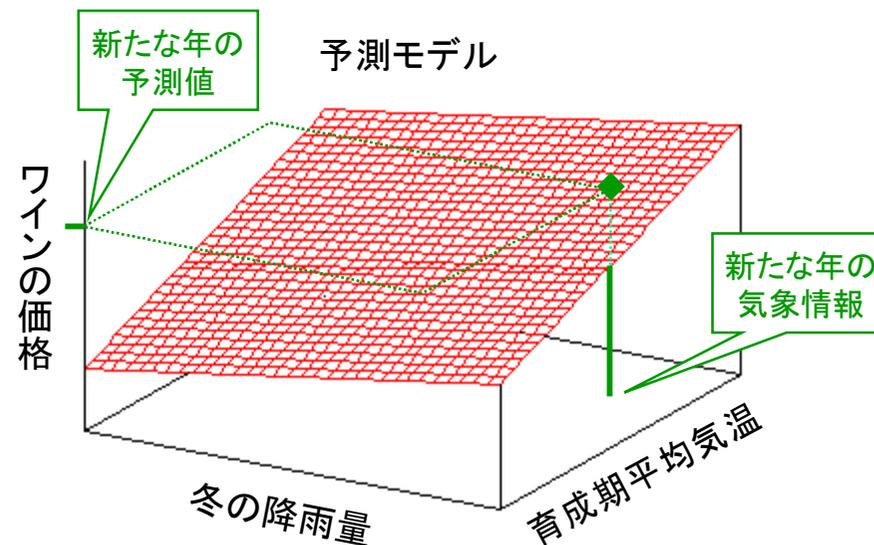
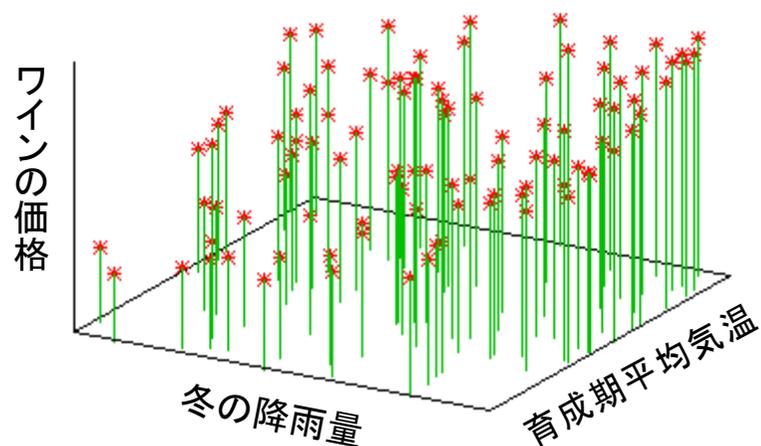


- 将来のワイン価値を製造年の気象情報から予測するモデル

$$\text{ワインの質} = 12.145 + 0.00117 \times \text{冬の降雨量} + 0.0614 \times \text{育成期平均気温} - 0.00386 \times \text{収穫期降雨量}$$

→「1989年ものは今世紀最高のヴィンテージで、翌年1990年ものはそれ以上」といった予測が、製造時点で可能に。

回帰分析

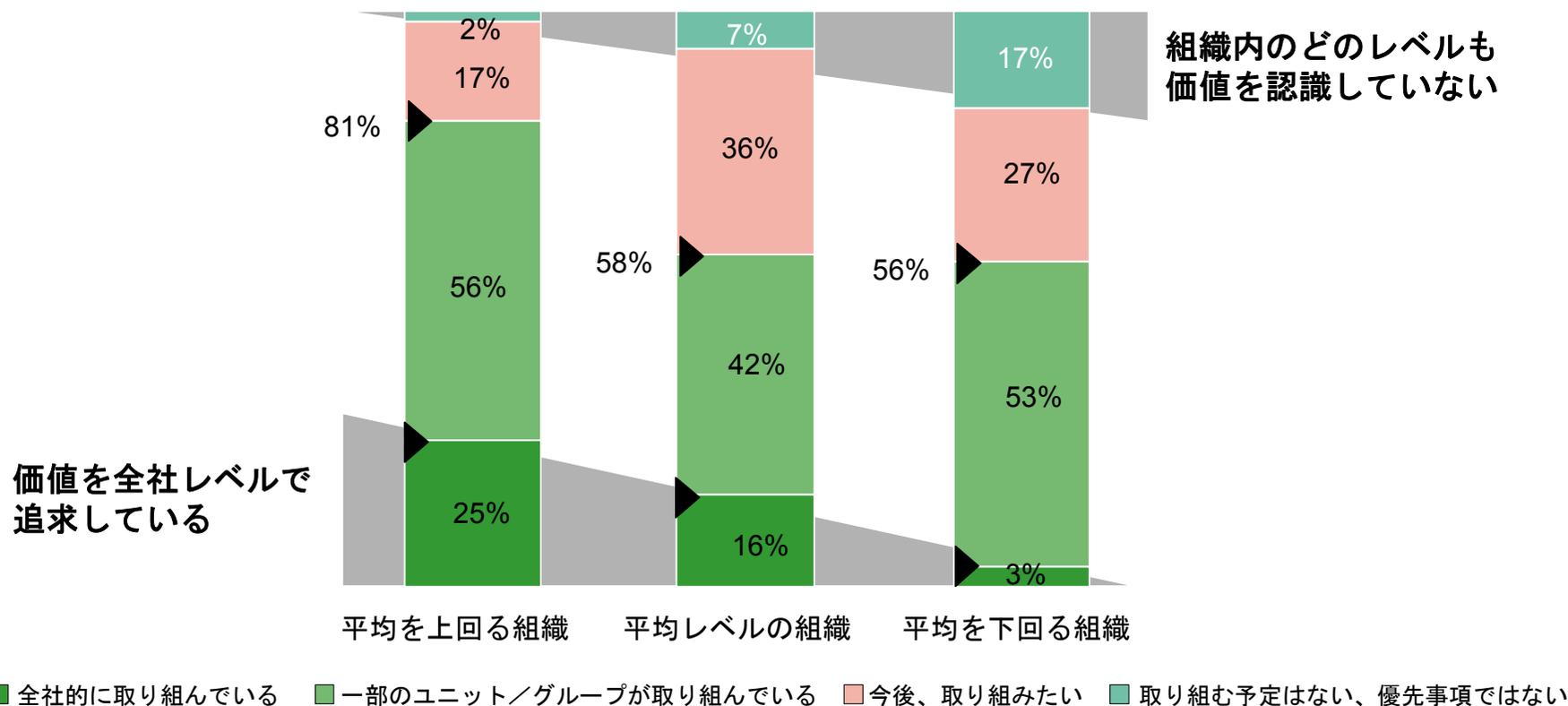


情報活用が競争力になる

■比較的業績の優れている企業は情報の価値を認識し積極的に追求しています。

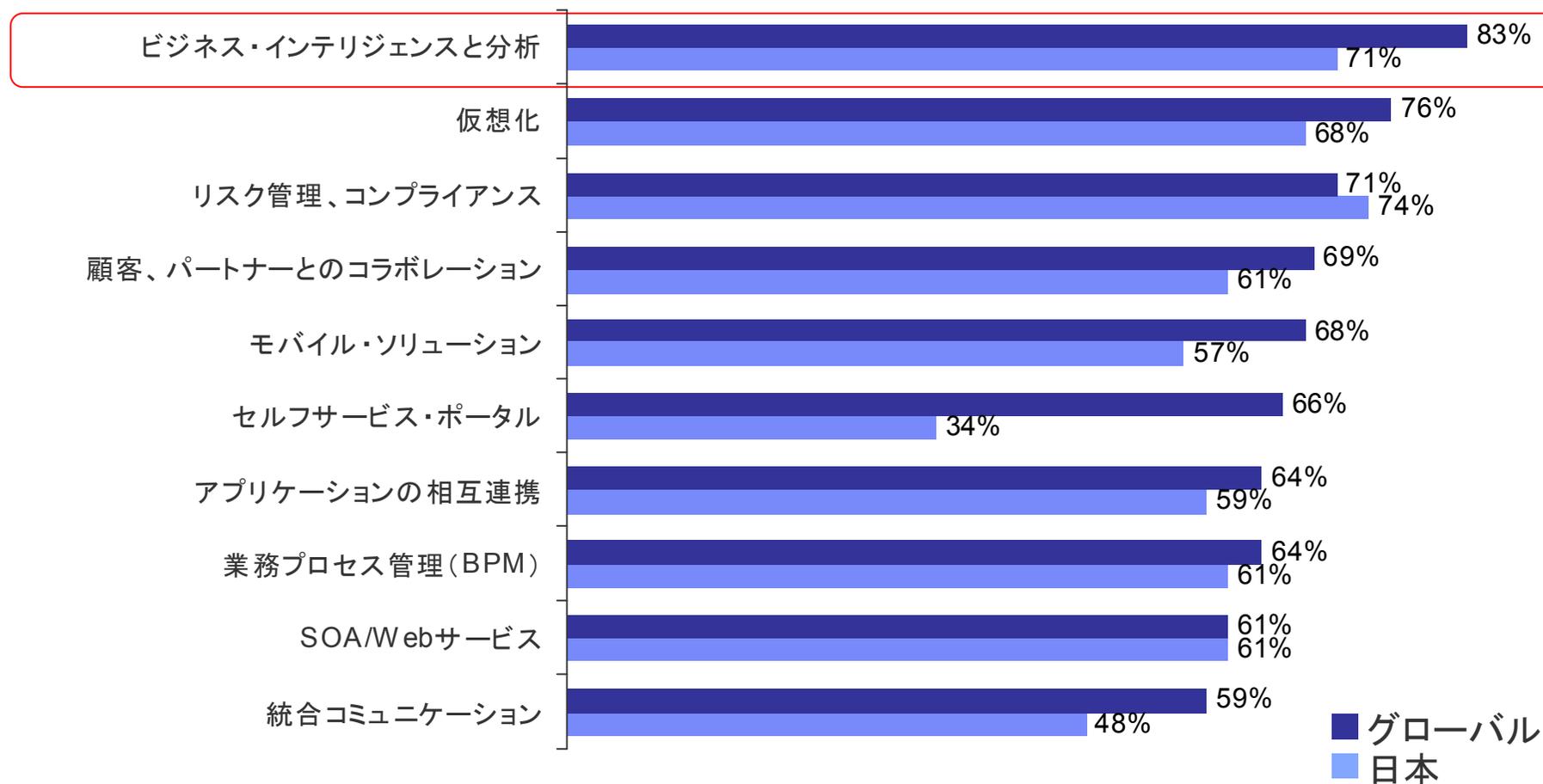
Q. あなたの組織には、次の分野における改善計画がありますか？

- ①データの収集および分析能力
- ②関連情報の提供
- ③従業員が情報に基づいて行動する権限の付与



情報活用に対する需要が高まっています

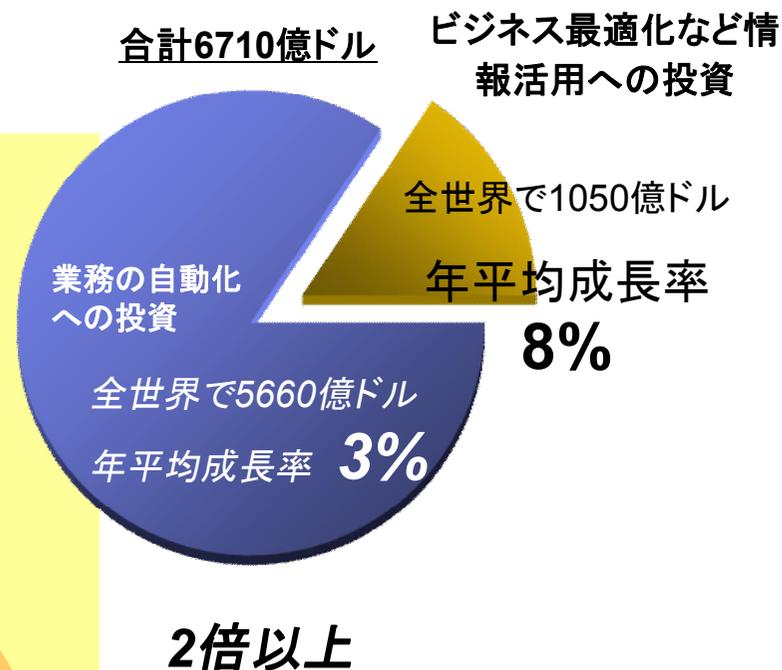
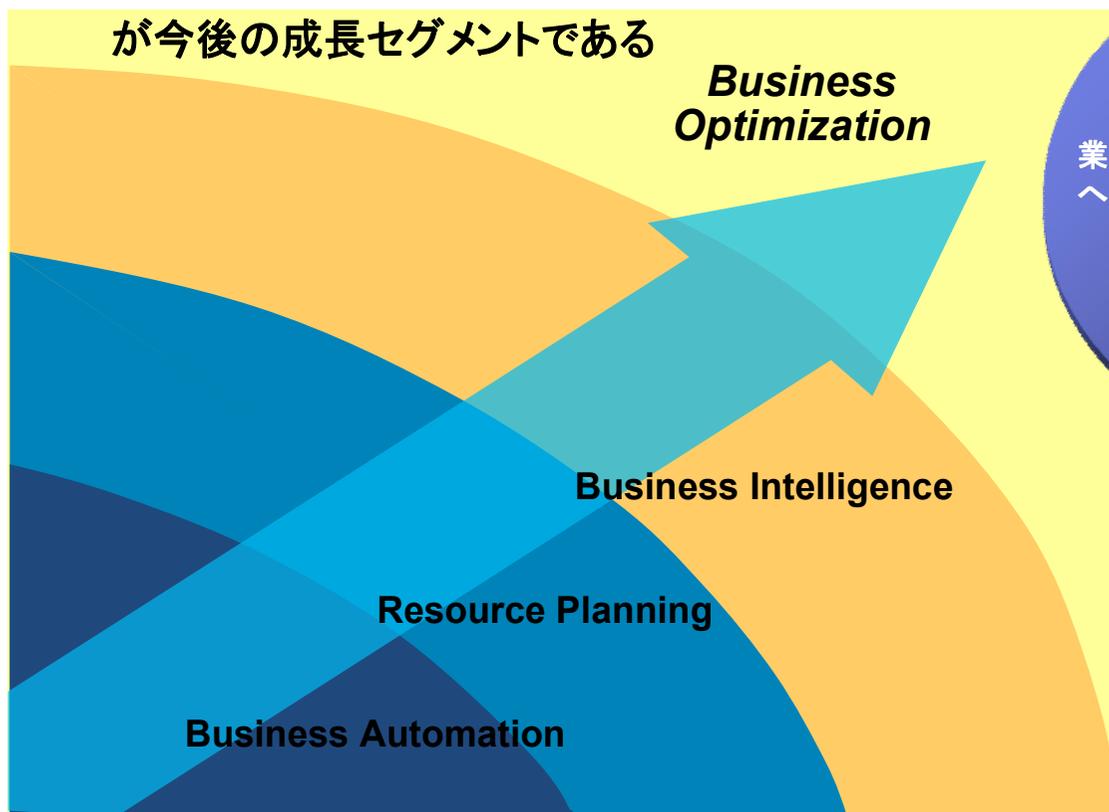
Q. 競争力強化のために、どのような分野での取り組みを検討されていますか？



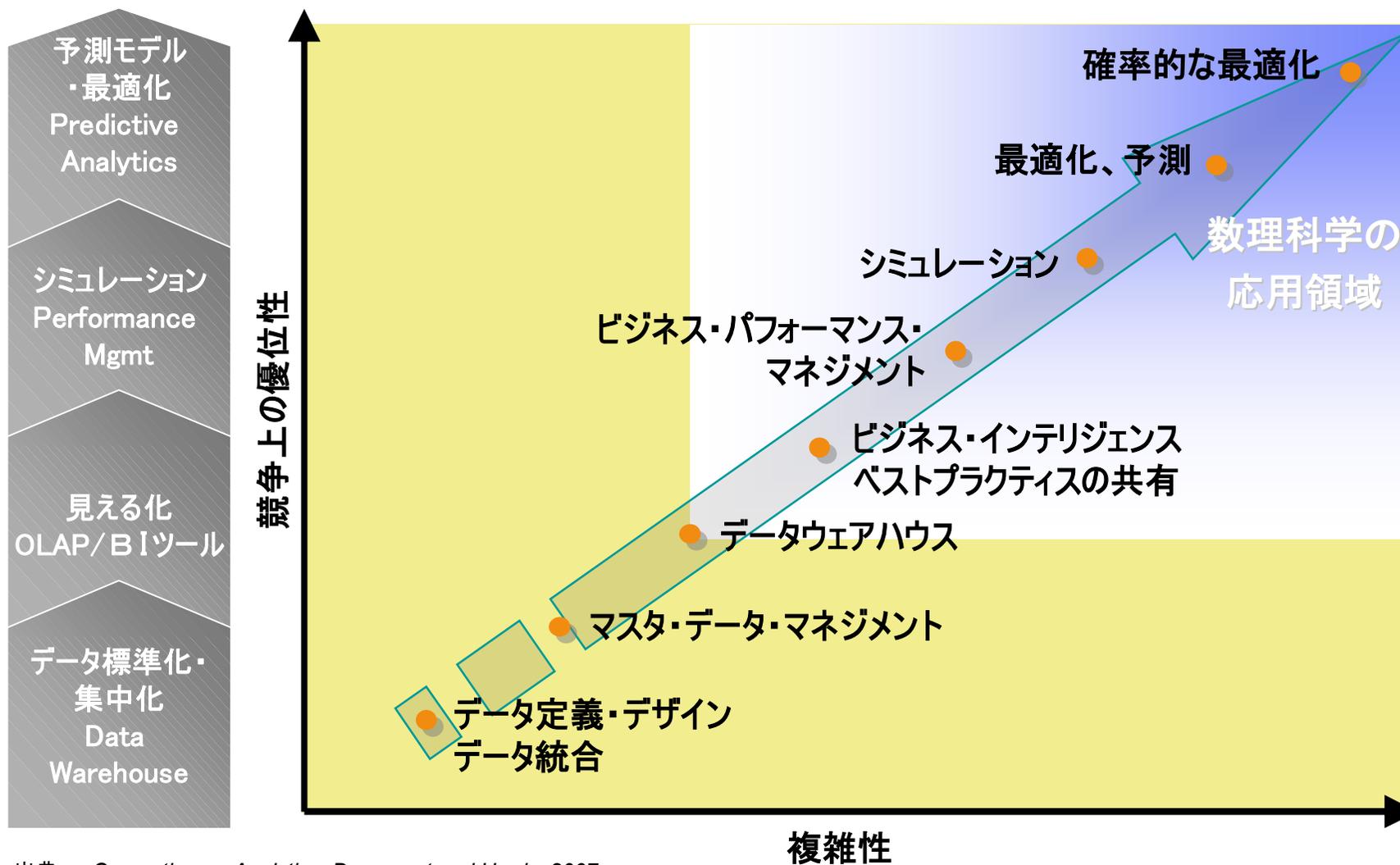
出典：IBM Global CIO Study 2009; 世界78カ国、19業種、2,345名（日本からは162名）の企業ならびに公共機関のCIOの方々に対して実施した、役割や行動様式に関するインタビュー調査による。2009年1月～4月にかけて調査。

データマイニングへの投資は、IT技術の成熟化に伴う歴史的必然です

- 比較的単純な業務の自動化を目的としたIT投資は漸減傾向
- 情報活用を目的とした積極的なIT投資が今後の成長セグメントである

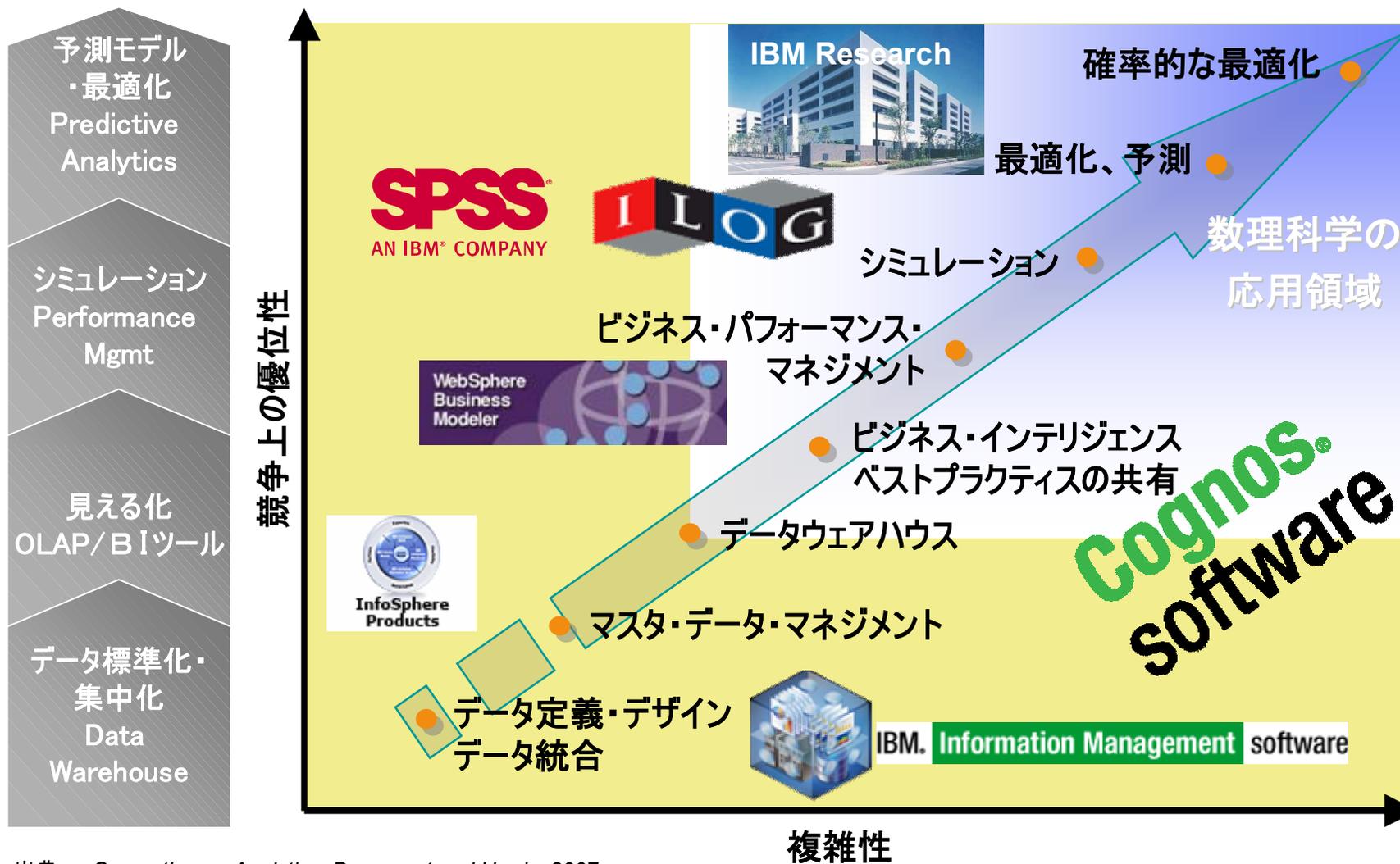


BAOサービスの実行のため、IBMは製品の品揃えを系統的に充実させています



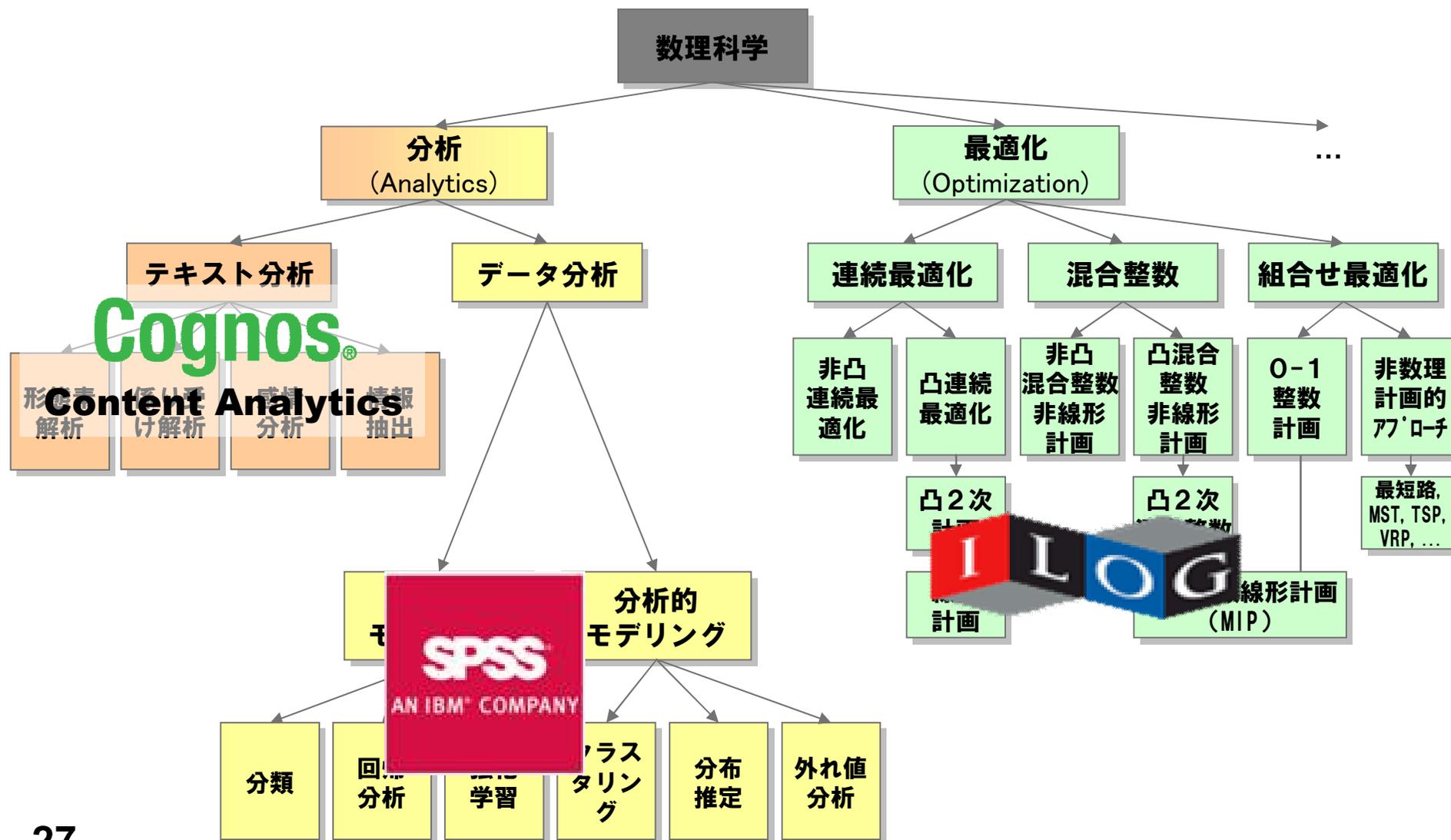
出典： *Competing on Analytics, Davenport and Harris, 2007*

BAOサービスの実行のため、IBMは製品の品揃えを系統的に充実させています



出典： *Competing on Analytics, Davenport and Harris, 2007*

(ご参考) データマイニング技術に関するIBMのソフトウェア
製品: テキストマイニング、データマイニング、数理最適化



このパートのまとめ

- データマイニングはIBMの最重点投資領域であり、ビジネス拡大のための施策を戦略的に行っています
- IBMのBAOサービスは、データマイニングの技術を活用して合理的な判断を支援するサービスです

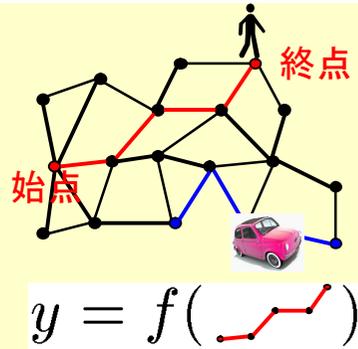
目次

- イン트로ダクション
 - 機械学習/データマイニングと制御工学
 - データマイニングをめぐるIBMのビジネス戦略
 - 私たちが解いている問題の例
- スパース構造学習による異常検知
- まとめ

私たちが解いている問題の例(1/2)

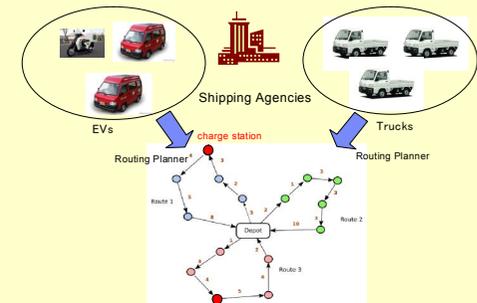
トラジェクトリ解析

- 自動車や人間の移動の軌跡を解析
- マーケティング等に有用な情報を抽出



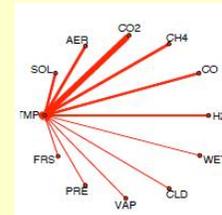
配送経路最適化

- 電気自動車のバッテリーに関する制約を正しく取り込んで配送経路を最適化



私たちが解いている問題の例 (2/2)

事例: 地球温暖化の原因推定



- 異常な温度上昇と、他の環境要因との依存関係を気象データから自動学習
- CO₂との因果関係が見出された

Deep QA – クイズ番組でチャンピオンを目指す

2009年、IBMは、「ジョパディー」というアメリカのクイズ番組に挑戦することを宣言しました。最先端の自然言語処理技術と高度な数理解析技術に基づいて、任意に与えられた質問文に高速に答えを返します



Watson takes on Jeopardy!

Advanced computing system has potential to take business intelligence to a new level

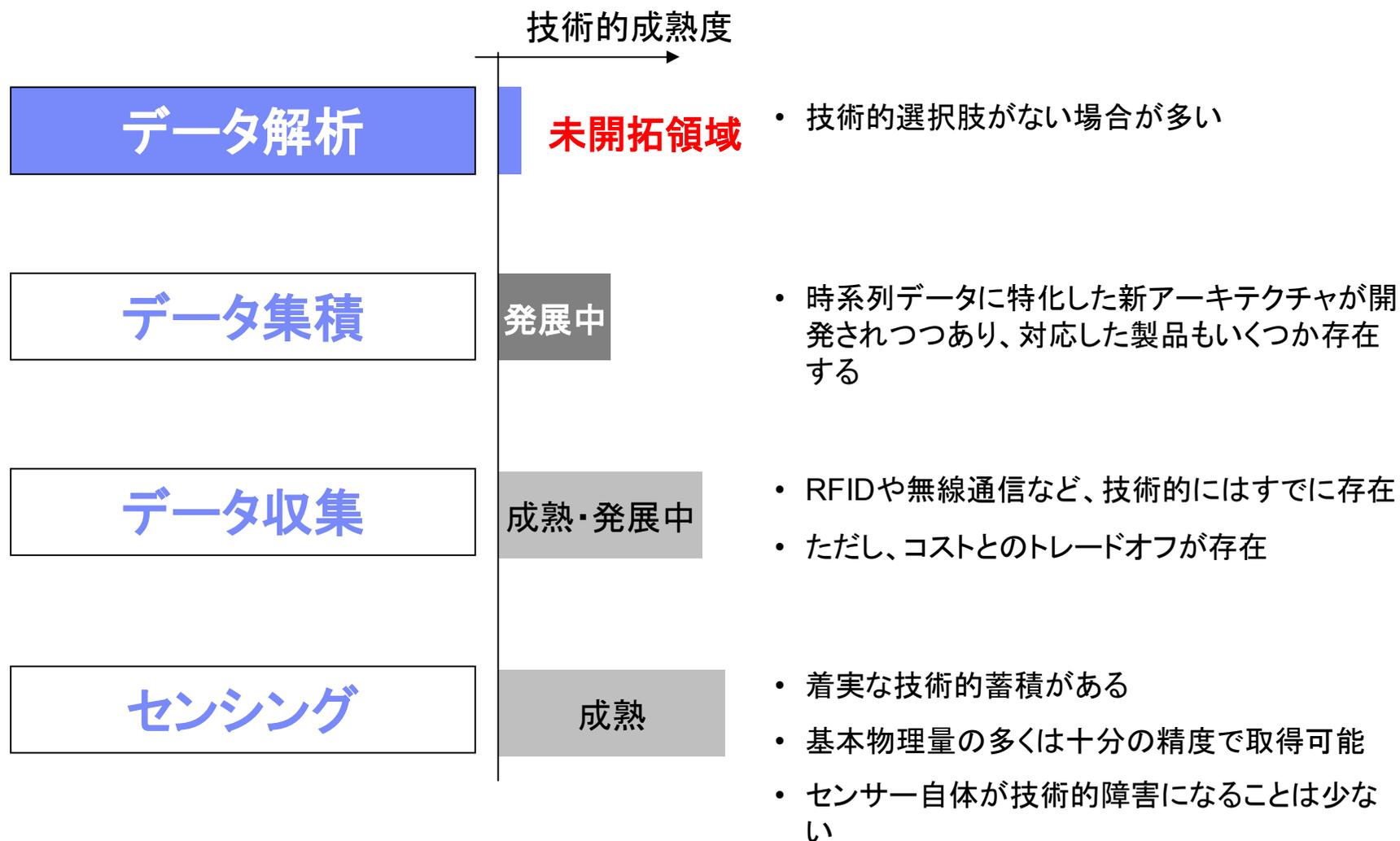
東京の数理科学チームのひとつのフォーカスは、センサーデータの活用のための技術です

- 比較的安価にセンサーデータを取得・蓄積できるようになった今、いくつかの業界では、保守・保全に関する業務の自動化と収益化に興味をもたれています。
- しかしそのためには、データ管理インフラに加えて、センサーデータの解析技術の進歩が不可欠です

建設機械の早期異常検
知とメンテナンスサービス

ある建機メーカーでは、世界中の建機の稼動状況をモニタリングし、保守業務に役立てるシステムを構築している

センサーデータの活用はまだ進んでいないのが現状です



なぜセンサーデータは活用されていないのでしょうか

巨大なデータサイズ

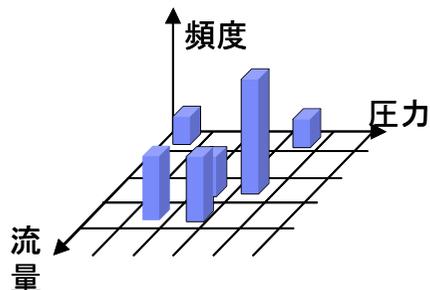
伝統的な関係DBは、物理センサーからの巨大な時系列データの扱いに最適化されてはいるわけではありません

数百のセンサーからの時系列データを数日間蓄積するだけで、そのサイズは容易にテラバイト程度となります。

クエリ処理技術の欠如

簡単な集計をしようと思っても、素朴な方法で作られたDBでは、処理に時間がかかりすぎるのが通例です

例えば、ある圧力と流量の範囲に入る頻度を知りたい、というような素朴なクエリであっても、素朴に設計されたDBでは、処理に相当の時間がかかってしまいます。



解析技術の欠如

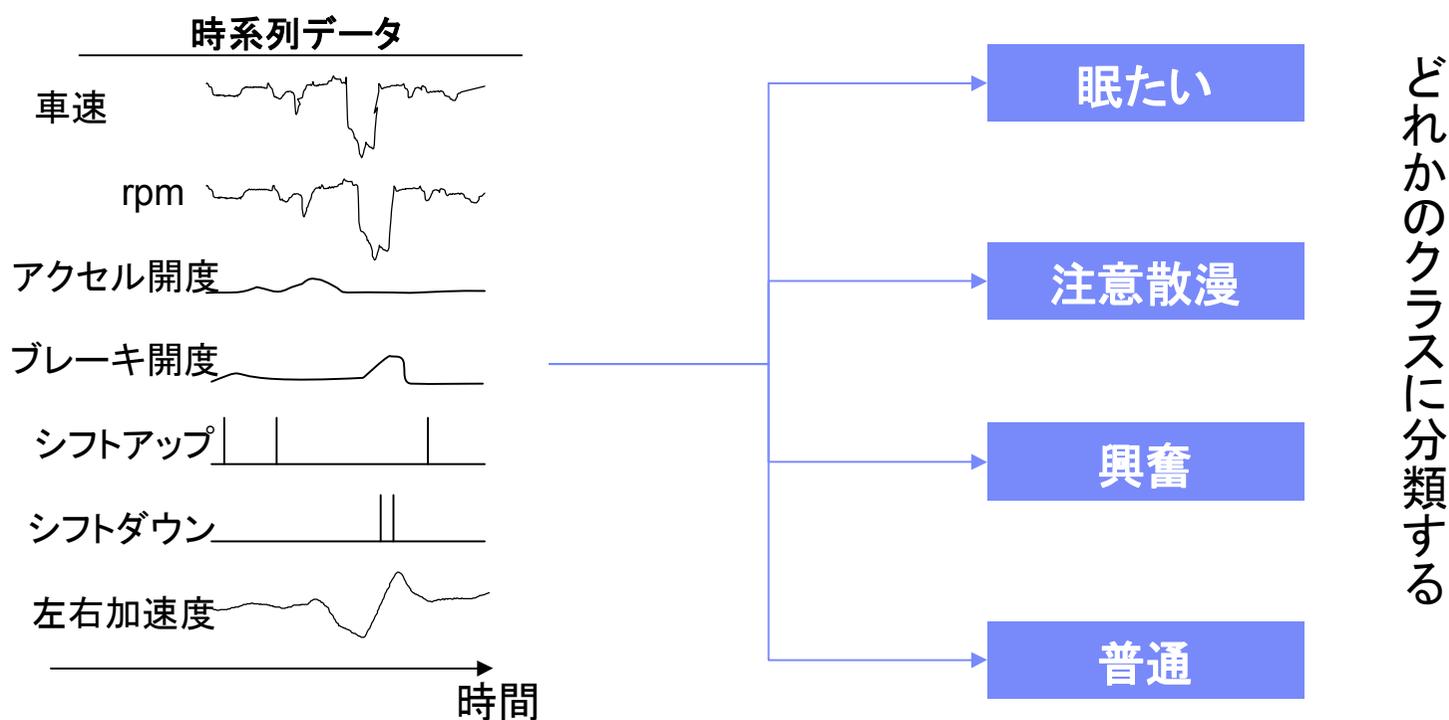
売り上げデータなどと異なり、そもそもどういう「クエリ」をすれば不具合解析が可能なのか自明で花ありません

不具合解析のためには、素朴な集計を越えた高度な解析技術が必要です。現状ではそのような実用的機能を提供している時系列DB製品は存在しません。

(ご参考)脳科学の分野ではいわば「念力」が実用化されつつあり、センサーデータ解析に進歩をもたらしています

- 脳波の微弱な電気信号を時系列データとして取得
- それを元に思考を認識
 - 「上」と念じるとカーソルが上、など
- 手を使わずに脳で念じるだけでテレビゲームができる

(ご参考) 私たちはこの技術を、センサーデータからの運転手の状態判別などに適用しています



(ご参考) 運転手の状態判別

このパートのまとめ

- IBMの数理解析チームは、機械学習と最適化技術の専門知識を背景に、多くの実問題を解いてきました
- センサーデータの解析は技術の進歩が待たれる領域であり、私たちのフォーカスのひとつです。

目次

- イン트로ダクション
 - 機械学習/データマイニングと制御工学
 - データマイニングをめぐるIBMのビジネス戦略
 - 私たちが解いている問題の例
- **スパース構造学習による異常検知**
- まとめ

内容

- やりたいこと
- グラフィカル・ガウシアン・モデルと関連研究
- 疎構造学習の方法
- 相関異常度の定義
- 実験結果
- まとめ

▶ Acknowledgement

- This is a joint work with Aurelie C. Lozano, Naoki Abe, and Yan Liu (IBM T. J. Watson Research Center).

資料は下記にあります。

- http://latent-dynamics.net/01/2010_LD_Ide.pdf

[Japanese / [English](#)]

Latent Dynamics 研究会

— 潜在世界の中でどのような変化が起こっているか —

関連イベント

第1回 Latent Dynamics ワークショップ

- 日時：2010年6月16日(水) 10:00 – 17:00
- 東京大学工学部6号館3階 セミナー室A・D
- [第1回 IBISML研究会](#)と連続開催

詳細は[こちら](#)をご覧ください。

発起人

- [山西健司](#) (東京大学情報理工学系研究科数理情報学専攻)
- [太澤幸生](#) (東京大学工学系研究科システム創成学専攻)
- [井手剛](#) (IBM東京基礎研究所)

目的・背景

大量データが溢れる現在、データの表層をみただけではわからない潜在世界に注目し、その変化を捉えることが重要性を増している。

第1回Latent Dynamic
Workshop での話

目次

- **イントロダクション**
 - 機械学習/データマイニングと制御工学
 - データマイニングをめぐるIBMのビジネス戦略
 - 私たちが解いている問題の例
- **スパース構造学習による異常検知**
- **まとめ**