

# Hadoop Recommendation Machine Learning

本文中の会社名・製品名・サービスネームについて  
Amazon Web Services は Amazon.com, Inc.の商標または登録商標です。  
Apache Hadoop は Apache Software Foundationの商標または登録商標です。  
hybris は hybris AGの商標または登録商標です。  
その他、すべての会社名・製品名・サービスネームは、それぞれ各社の商標または登録商標もしくはサービスマークです。

2014/7/16  
EVF2014

スマートプレイス開発部 吉田 匠

- **吉田 匠**

- スマートプレイス開発部
- 11年目



- **検索エンジンを苦節(?) 8年**

- 社内、コーポレートサイト、ECサイト
- Apache Solrが主力、Elasticsearch は次世代のホープ！
- 昔は、FAST ESP / Autonomy IDOL / Verity K2 など

- **今年のテーマはクラウドと機械学習**

- AWS ソリューションアーキテクト取りました！
- Softlayer は勉強中？
- Hadoop / Mahout を中心に調査・検証



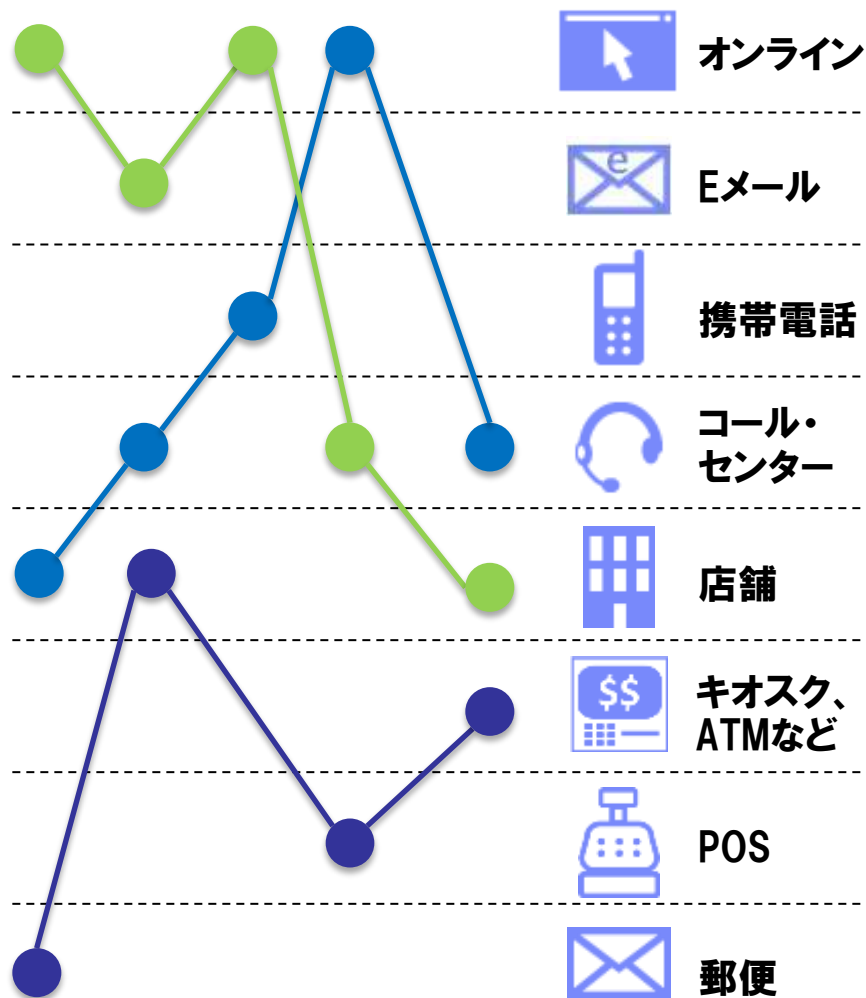
- **ことの経緯**
  - スマートプレイスについて
  - Hadoop についてのおさらいと最新の動向
- **レコメンドとは**
  - レコメンドの基本的な仕組みと課題
  - オムニチャネル・レコメンド
- **機械学習への応用**
  - 機械学習とは
  - 現実的な課題と皆様へのお願い

# ことの経緯

## 「マルチチャネル」から 「オムニチャネル」へ

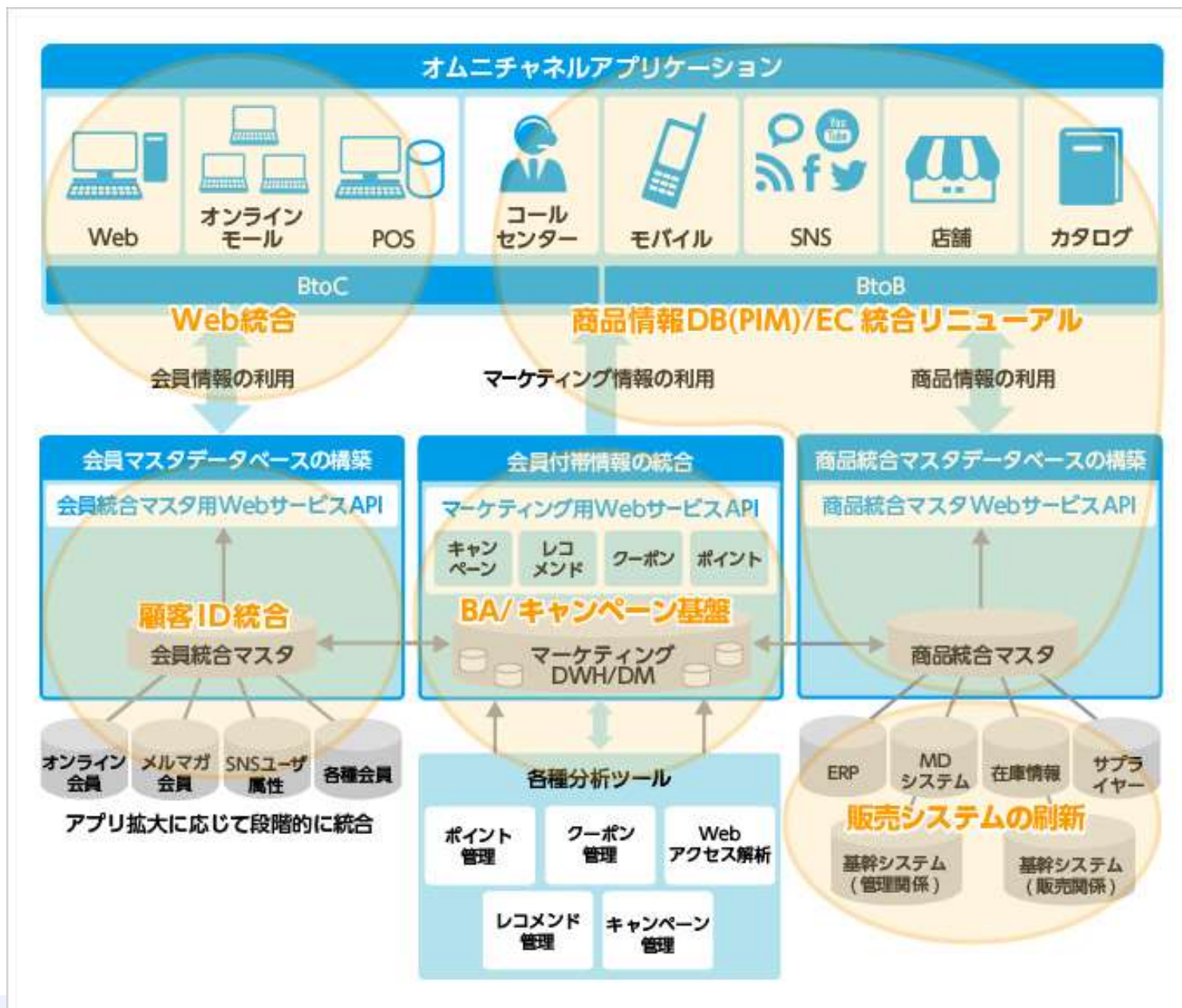


顧客

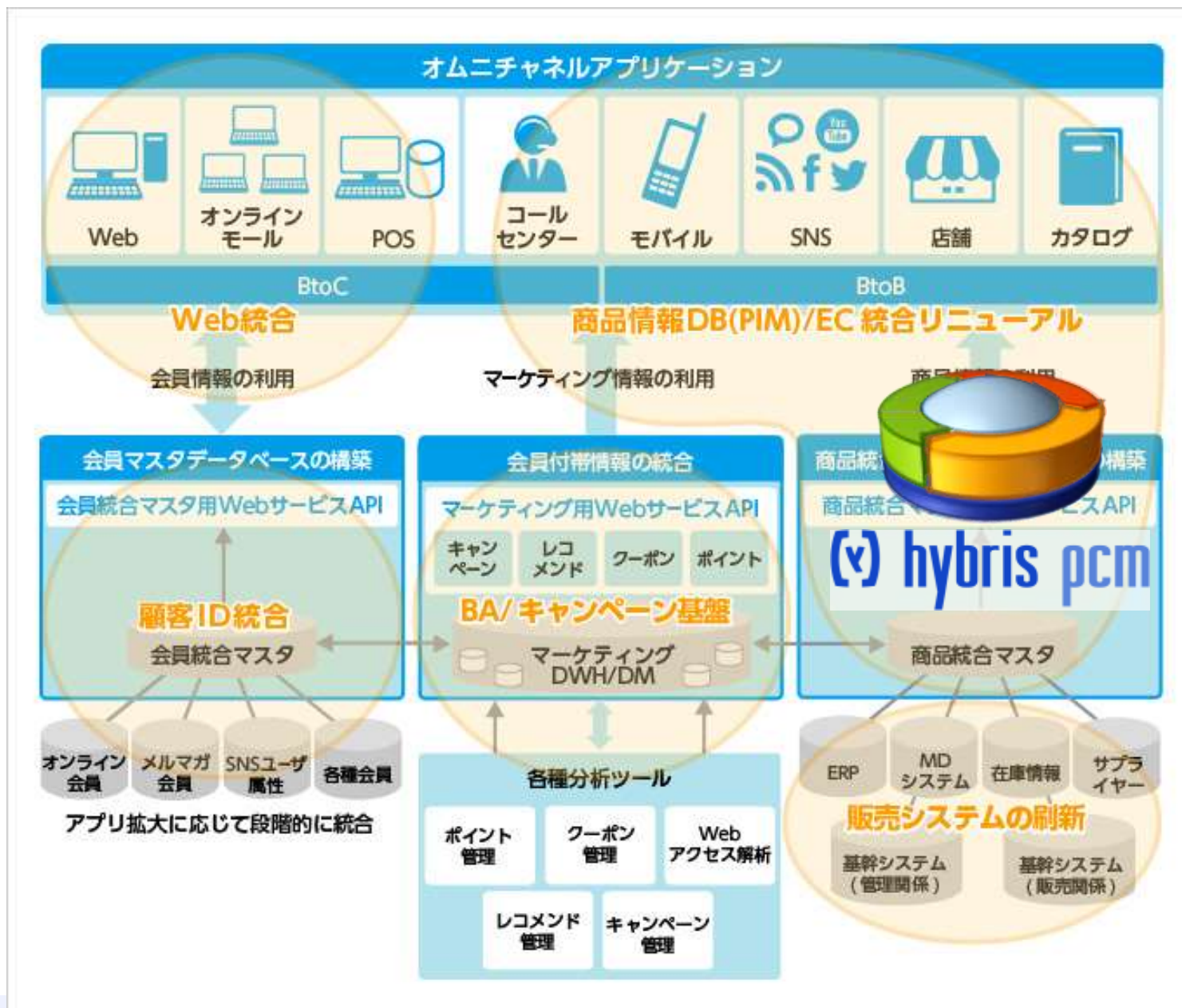


企業

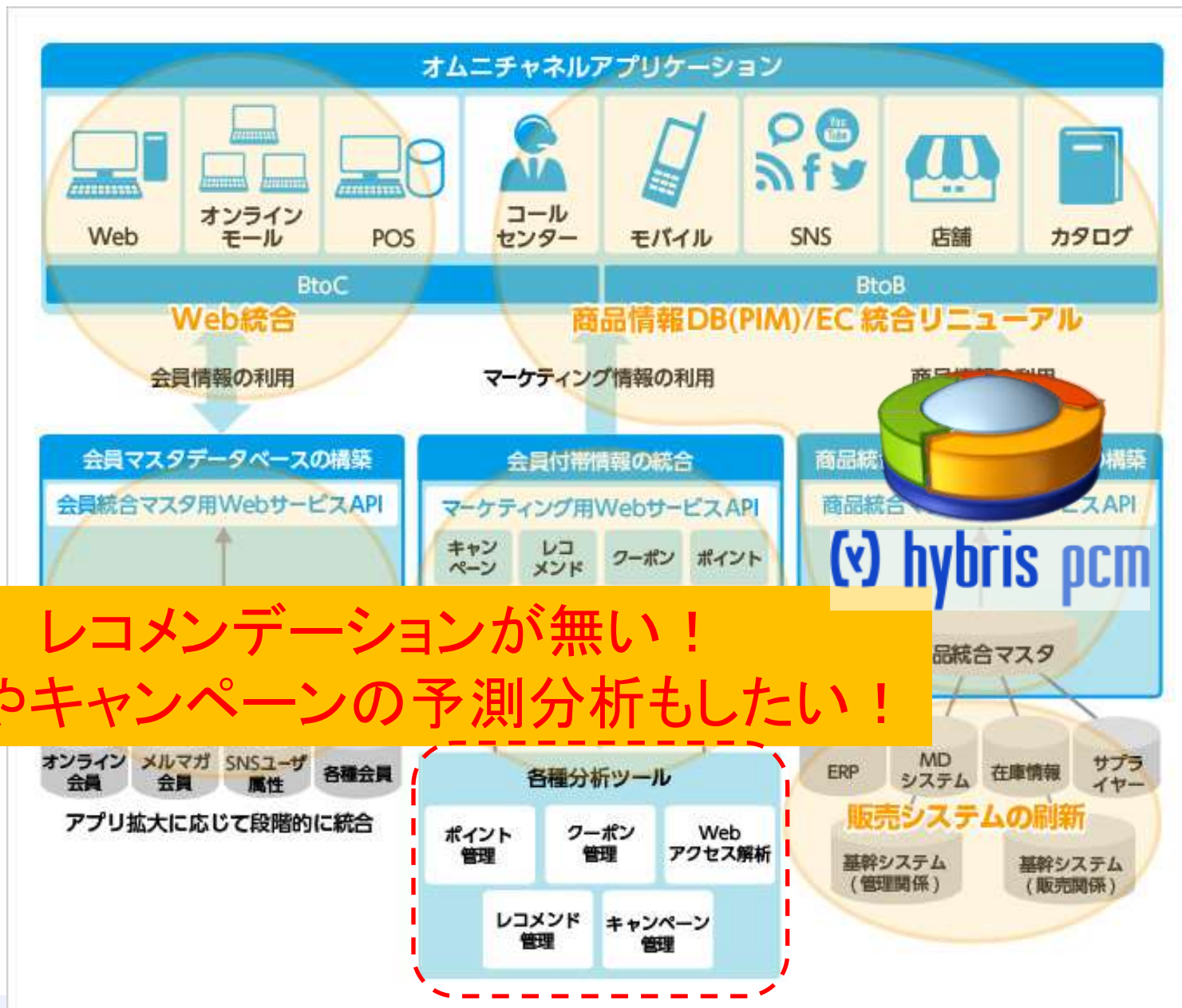
# スマートプレイスとは？



# スマートプレイスとは？



# スマートプレイスとは？



レコメンデーションが無い！  
売上やキャンペーンの予測分析もしたい！



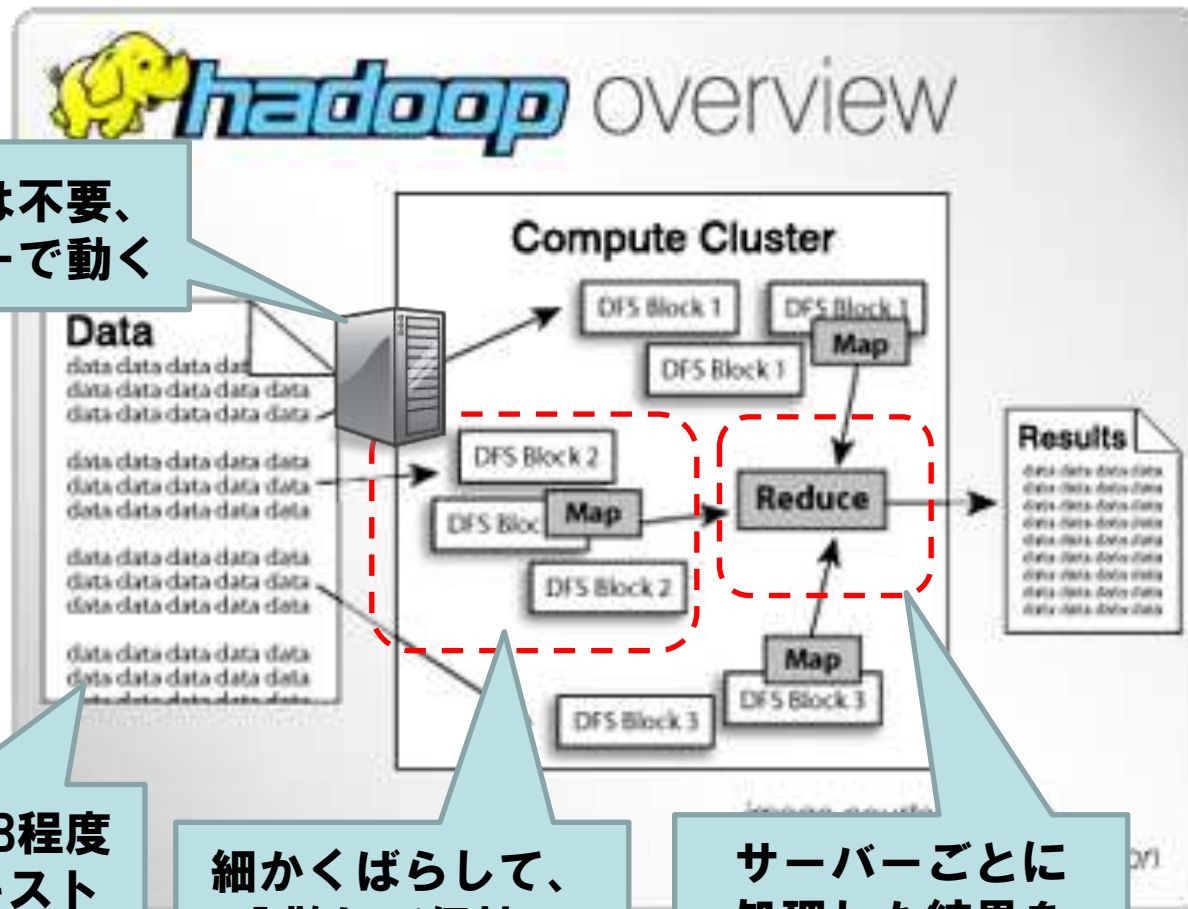


って、何だっけ??

- Hadoop とは、分散ファイルシステム (HDFS) を使用した、高速な**バッチ処理基盤**です。
- 大量のデータを高速に**バッチ処理**できるという効果があります。
  - 多数のマシンに分散することで、高スループットを実現し処理時間を短縮
- MapReduce というフレームワークに沿って処理を記述する (プログラミングする) 必要がありました。
  - Hive (SQL)
  - Pig (DSL)
  - Hbase (KVS、DWH)
  - Mahout (機械学習、レコメンド)



# Hadoop とは？



特別なハードは不要、  
普通のサーバーで動く

数十GBからTB程度  
の巨大なテキスト  
ファイル

細かくばらして、  
分散して保持。  
行単位で処理

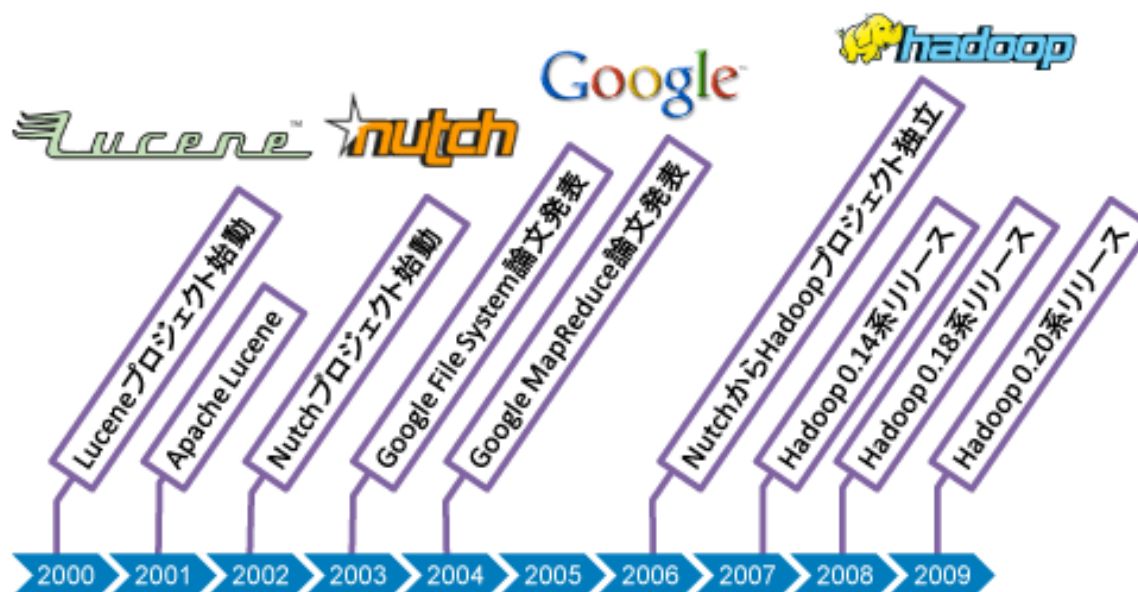
サーバーごとに  
処理した結果を  
集計、結合。

出典: How jStart is leveraging distributed computing for business  
<http://www-01.ibm.com/software/ebusiness/jstart/hadoop/>



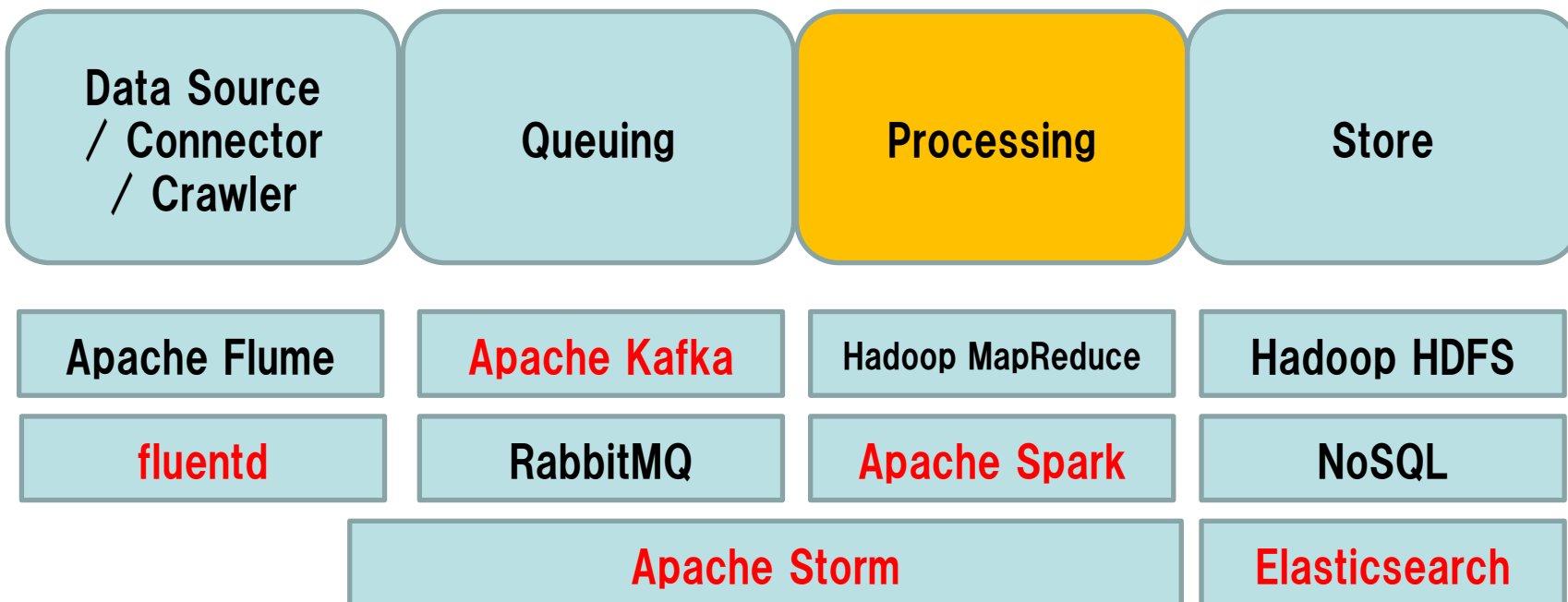
- 7月8日（火）  
@ベルサール汐留
- 参加者: 1299人
- 私服組（Web系企業）が多いが、メッセージボードを見ると **Sler** も多数参加。
- YARN（MRv2）、Spark および関連アプリ、機械学習のセッションが中心。

- ・ 開発開始からおよそ10年。安定、**普及期へ**。
  - Web系だけでなく、一般企業も検討を始めている
- ・ MapReduce からの脱却、パフォーマンスが劇的に向上。**リアルタイム性**をうたいはじめる。
  - 2013年10月 - Hadoop 2 リリース (MRv2 / YARN)
  - 2014年 5月 - Apache Spark リリース



出典: [第2回]Hadoopの生い立ち

<http://itpro.nikkeibp.co.jp/article/COLUMN/20120215/381721/>

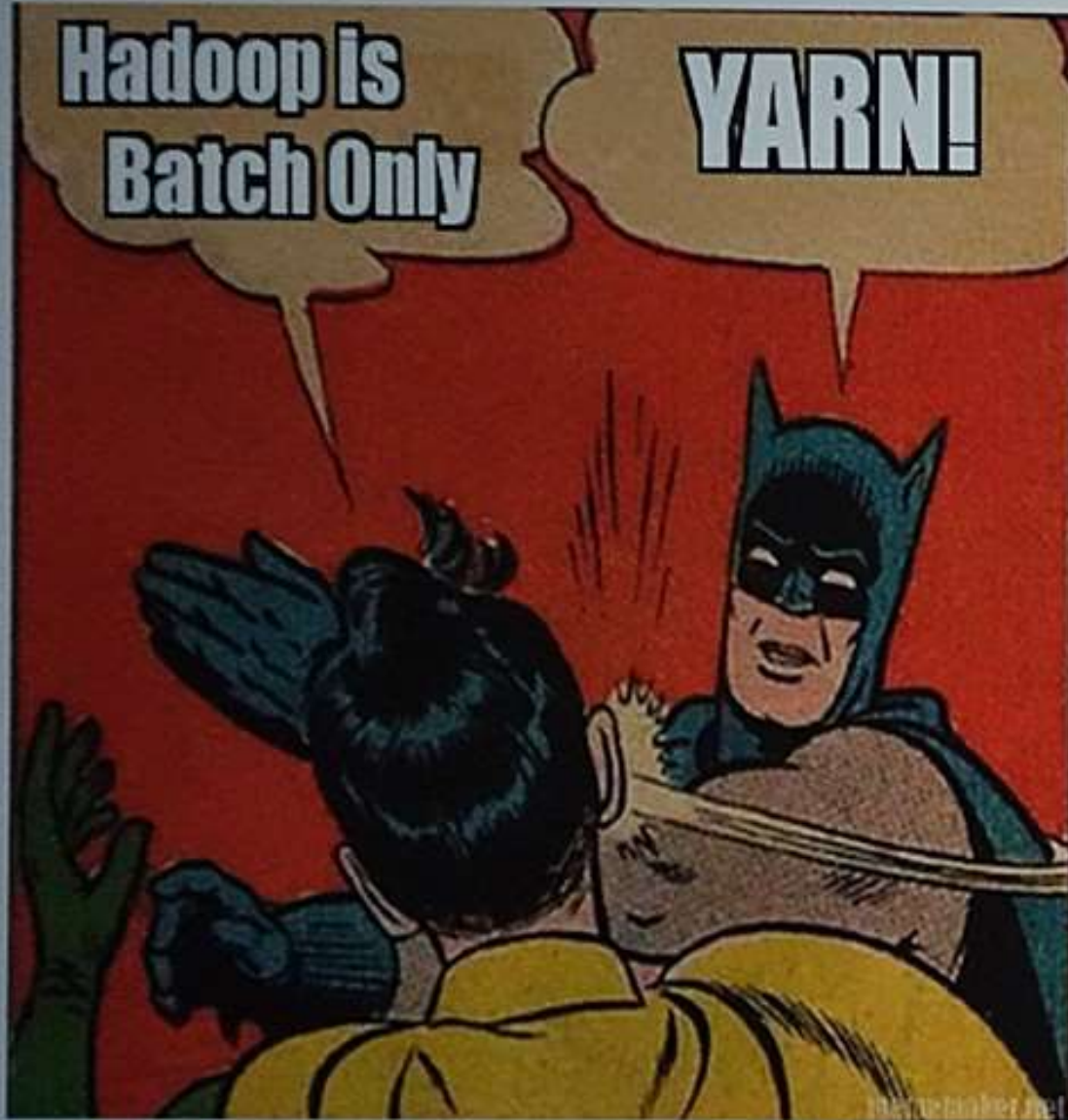


クラウド



## 速いは正義！

- 大量のデータに、繰り返し集計や統計処理を行い、業務上、優位なデータを生成する
- Hadoop の MapReduce を使ったソリューションがバッチ処理であった点に対し、処理を分散して高速化する優位性を、リアルタイムに実現する
- Applications
  - 売上やアクセス数の一時間あたりの速報値（KPI ダッシュボード）
  - センサーデータ分析
  - リアルタイム・レコメンド（生放送のコメントに連動したレコメンド）



## YARN Has Fundamentally Changed Hadoop

It Enables...

- **More Workloads**  
From batch to interactive & real-time
- **More Data**  
Multiple data sets of varying types and structures
- **More Value**  
Hosting multiple business cases in a single Hadoop cluster



出典 : Hadoop Summit 2014でホートンワークスが示したスライドより  
<http://itpro.nikkeibp.co.jp/article/COLUMN/20140619/565326/>



- ・ 検索ライブラリ Apache Lucene と生みの親が同じ
- ・ 検索エンジンも、Hadoopと同じく大量のデータを扱う



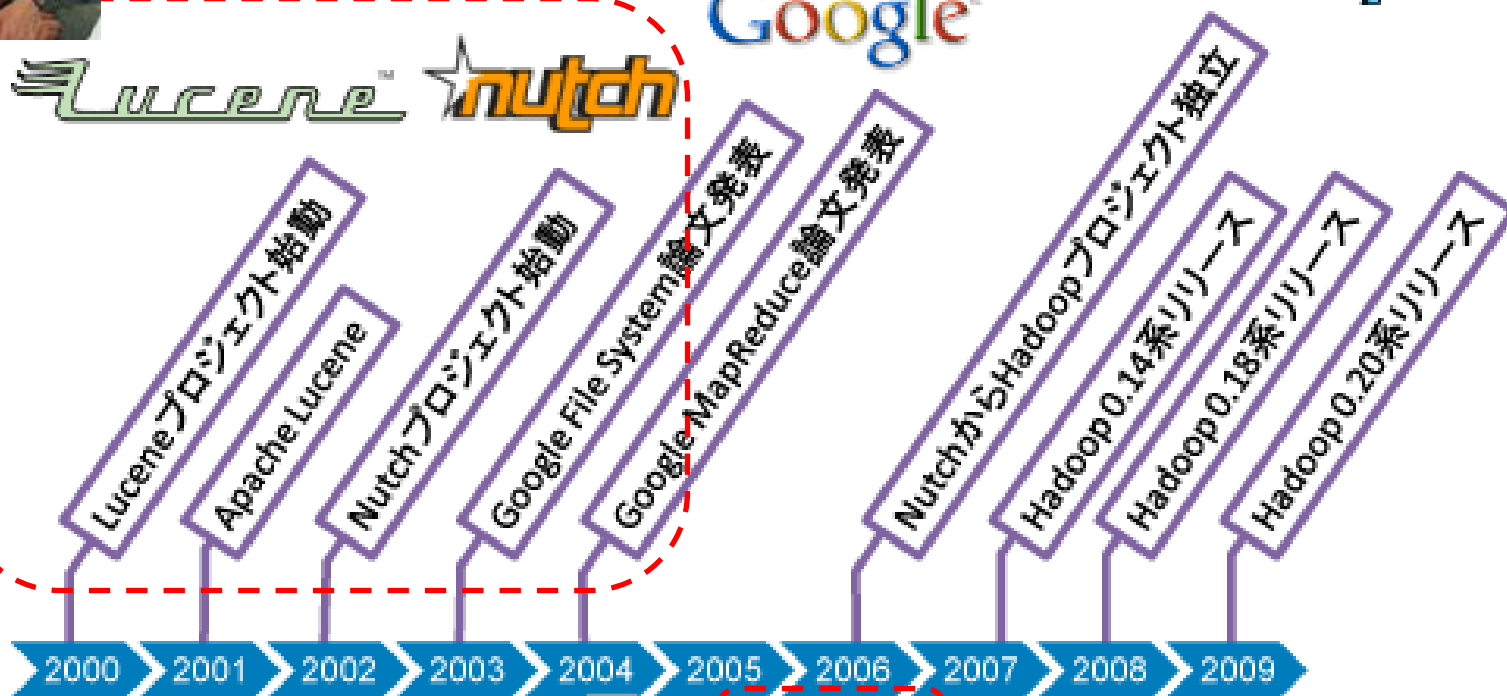
Dug Cutting

- Lucene
- Nutch
- Hadoop
- Solr は別の人。(Yonick Seely @ CNET Networks)



Google

Lucene Nutch



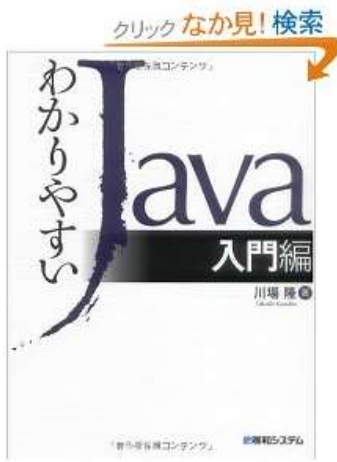
# レコメンドとは

amazon.co.jp マイストア | Amazonポイント | ギフト券 | タイムセール | 出品サービス | ヘルプ JCBのOki Doki ポイントが Amazon.co.jp で使える [広げて見る](#)

カテゴリからさがす    こんにちは、**吉田匠さん** アカウントサービス 今すぐ登録 プライム カート ほしい物リスト

本 | 詳細検索 | ジャンル一覧 | 新刊・予約 | Amazonランキング | コミック・ラノベ | 雑誌 | 文庫・新書 | Amazon Student | 本のお買い得情報

Amazon Kindleでは、[わかりやすいJava入門編](#)をはじめとする200万冊以上の本をご利用いただけます。 [詳細はこちら](#)



[自分のイメージを掲載する](#)  
[この本の中身を閲覧する](#)

## わかりやすいJava入門編 [単行本]

川場 隆 (著)   
★★★★☆ (30件のカスタマーレビュー)

価格: **¥ 3,024** 通常配送無料 [詳細](#)

**在庫あり。** 在庫状況について  
この商品は、Amazon.co.jp が販売、発送します。ギフトラッピングを利用できます。

住所からお届け予定日を確認  [詳細](#)

7/5 土曜日にお届けするには、今から **19時間 22分**以内に「お急ぎ便」または「当日お急ぎ便」を選択して注文を確定してください（有料オプション。Amazonプライム会員は無料）

**中古品の出品: 15¥ 2,025より**

フォーマット	Amazon 価格	新品	中古品
Kindle版	¥ 2,520	--	--
単行本	¥ 3,024	¥ 3,024	¥ 2,025

Would you like to see this page in English? [Click here.](#)

数量:

**新品を購入**

または  
1-Clickで注文する場合は、サインインをしてください。

**中古品を購入**  
[中古品 - 良い 詳細を見る](#)

価格: **¥ 2,221**

Amazon.co.jp が発送

または

amazon.co.jp

マイストア | Amazonポイント | ギフト券 | タイムセール | 出品サービス | ヘルプ

JCBのOki Doki ポイントが Amazon.co.jp で使える [広げて見る](#)

カテゴリ からさがす

こんにちは、吉田匠さん [アカウントサービス](#) | [今すぐ登録 プライム](#) | [カート](#) | [ほしい物 リスト](#)

本 | [詳細検索](#) | [ジャンル一覧](#) | [新刊・予約](#) | [Amazonランキング](#) | [コミック・ラノベ](#) | [雑誌](#) | [文庫・新書](#) | [Amazon Student](#) | [本のお買い得情報](#)

Amazon Kindleが... [詳しくはこちら](#)

## この商品を買った人はこんな商品も買っています

ページ: 1 / 19

- 

わかりやすい Java オブジェクト指向入門編  
川場 隆  
★★★★☆ (9)  
単行本  
¥ 2,808
- 

わかりやすい Java オブジェクト指向編  
川場 隆  
★★★★☆ (18)  
単行本  
¥ 3,024
- 

スッキリわかる Java 入門  
中山 清彦  
★★★★★ (50)  
単行本 (ソフトカバー)  
¥ 2,808
- 

なぜ、あなたは Java でオブジェクト指向開発が...  
小森 裕介  
★★★★☆ (25)  
単行本  
¥ 2,354
- 

ゲーム作りで学ぶ Java プログラミング入門 J...  
中島 省吾  
★★★★☆ (4)  
単行本  
¥ 2,376





[自分のイメージを掲載する](#)  
[この本の中身を閲覧する](#)

単行本	¥3,024	¥3,024	¥2,025
-----	--------	--------	--------

Amazon.co.jp が発送

または

- ・ 利用者の動線に基づいた共起分析  
- 次に何をおすすめしますか？

			
安藤さん	<input type="radio"/>	<input type="radio"/>	
佐藤さん			<input type="radio"/>
	<input type="radio"/>		

- ・ **さまざまなチャンネルでの利用履歴を元にしたレコメンド**
  - オンラインでの購入履歴を元に、店舗でおすすめ（タブレットと合わせてコンシェルジュ）
- ・ **パーソナライゼーション**
- ・ **速さ（早さ）は正義！**
  - 今日、今、おすすめすべき商品は何か？

**次世代オムニチャンネル・レコメンド**

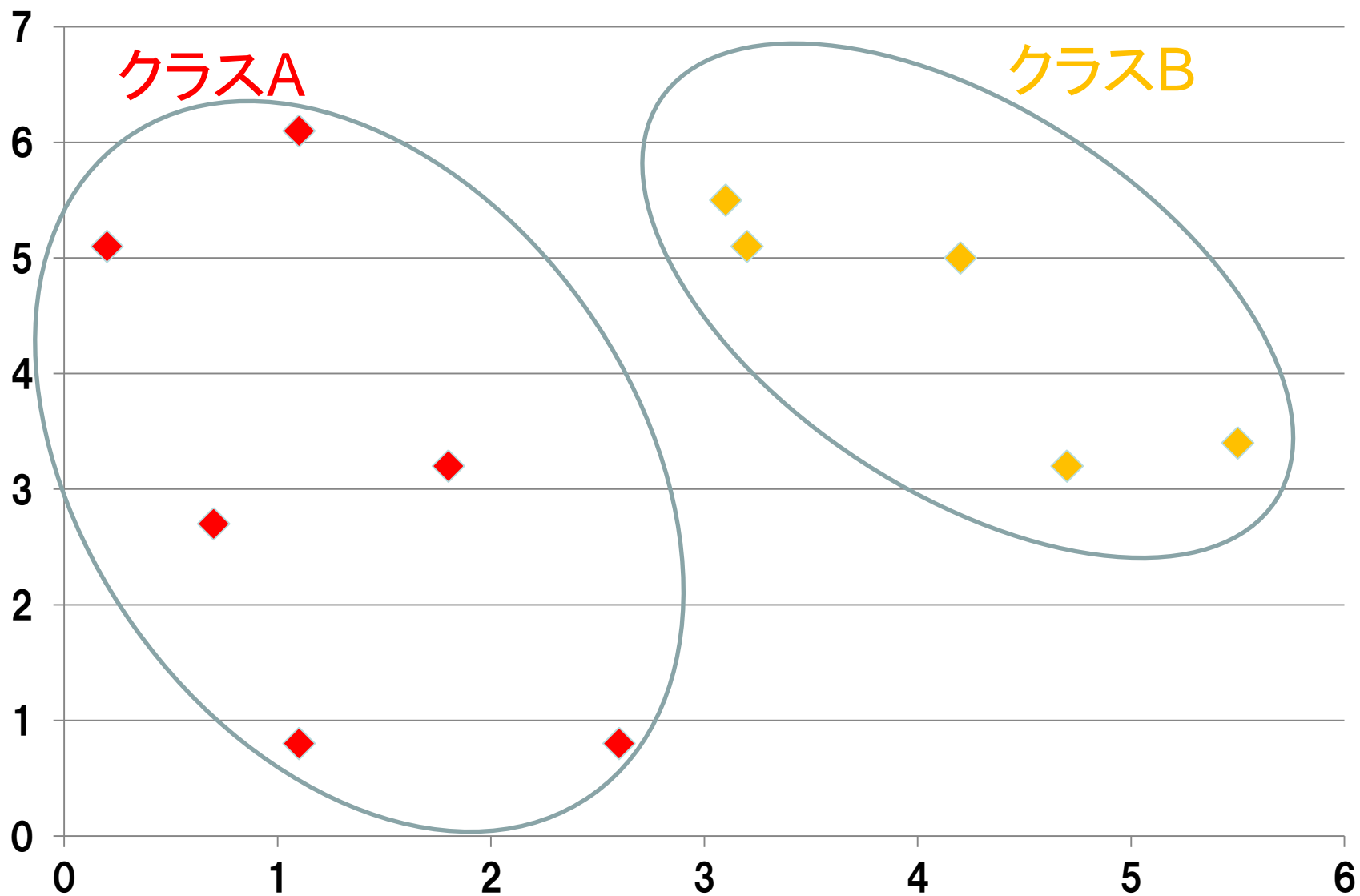
# 機械学習への応用

- ・ **機械学習とは「明示的にプログラミングすることなく、コンピュータに行動させるようにする」**
  - **コンピュータでアルゴリズムを構築し、学習データを読み込ませることで、自動的に今あるデータを分類&まだ見ぬデータを予測できるようにする**
  - **ヒトの情報処理能力を超えて、複雑なデータを分類&予測できるようにする**

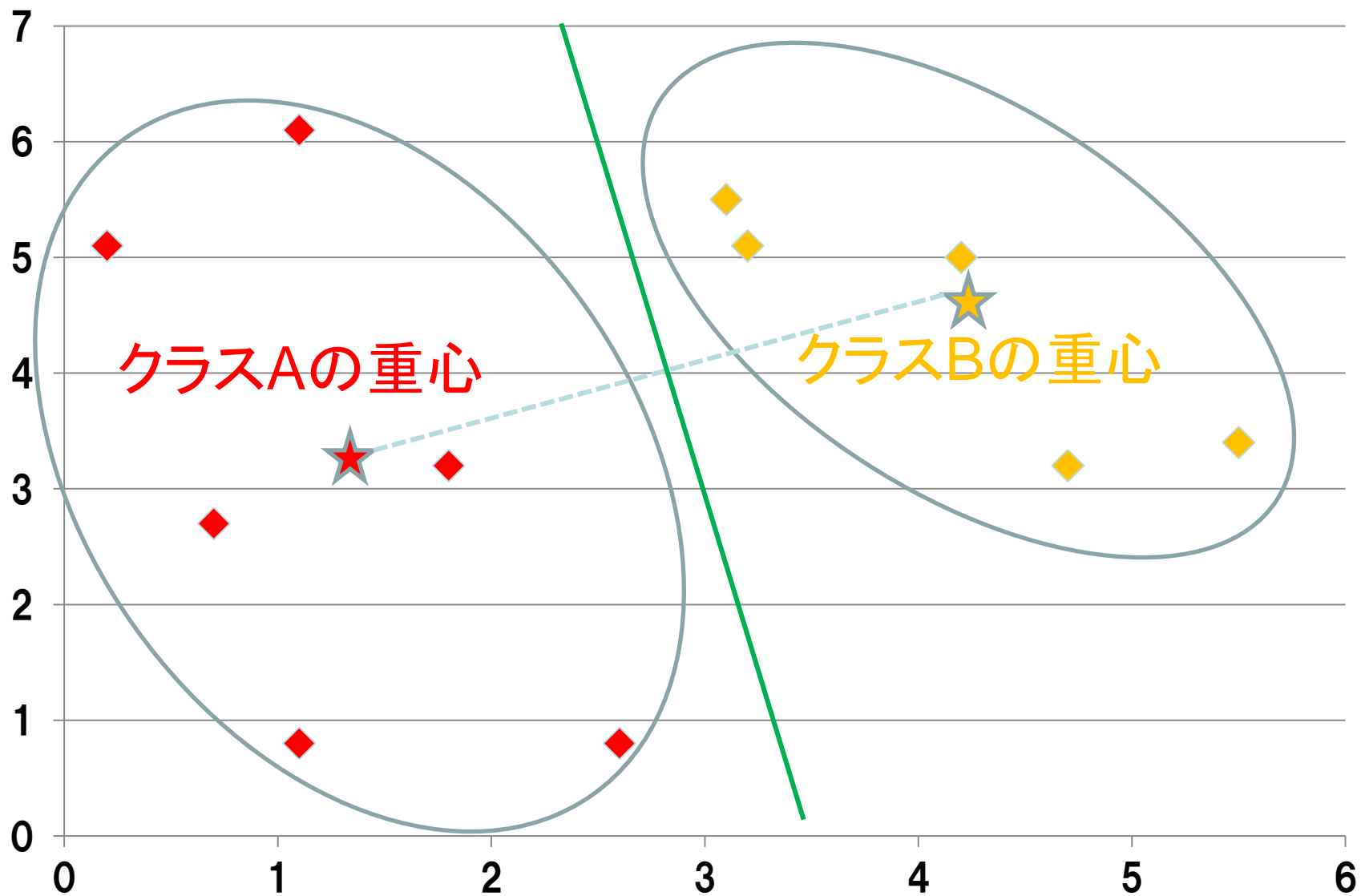
**レコメンドも機械学習の一種  
ノウハウが貯まれば、Hadoopと合わせて横展開が可能**

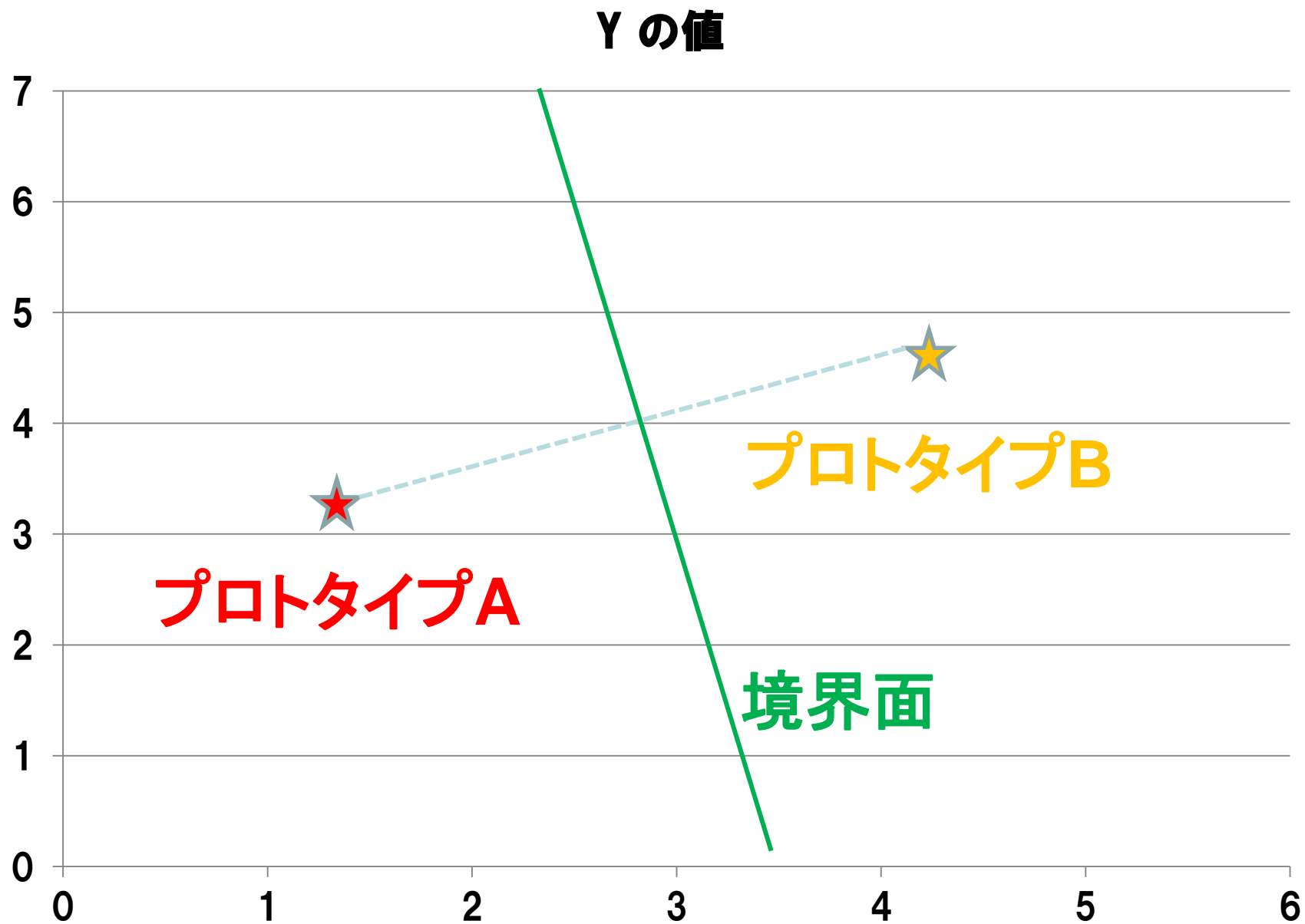


## Y の値

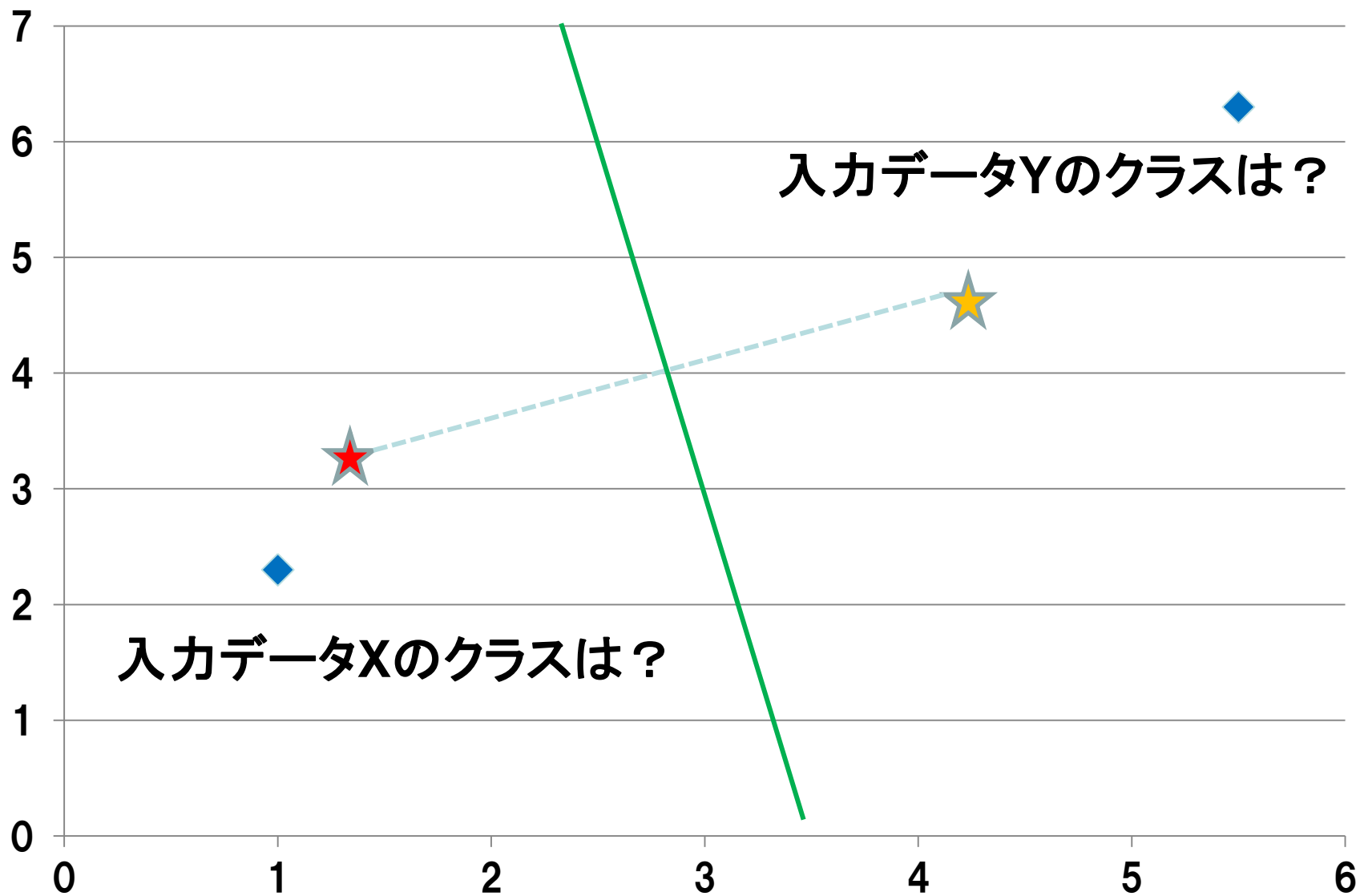


Y の値

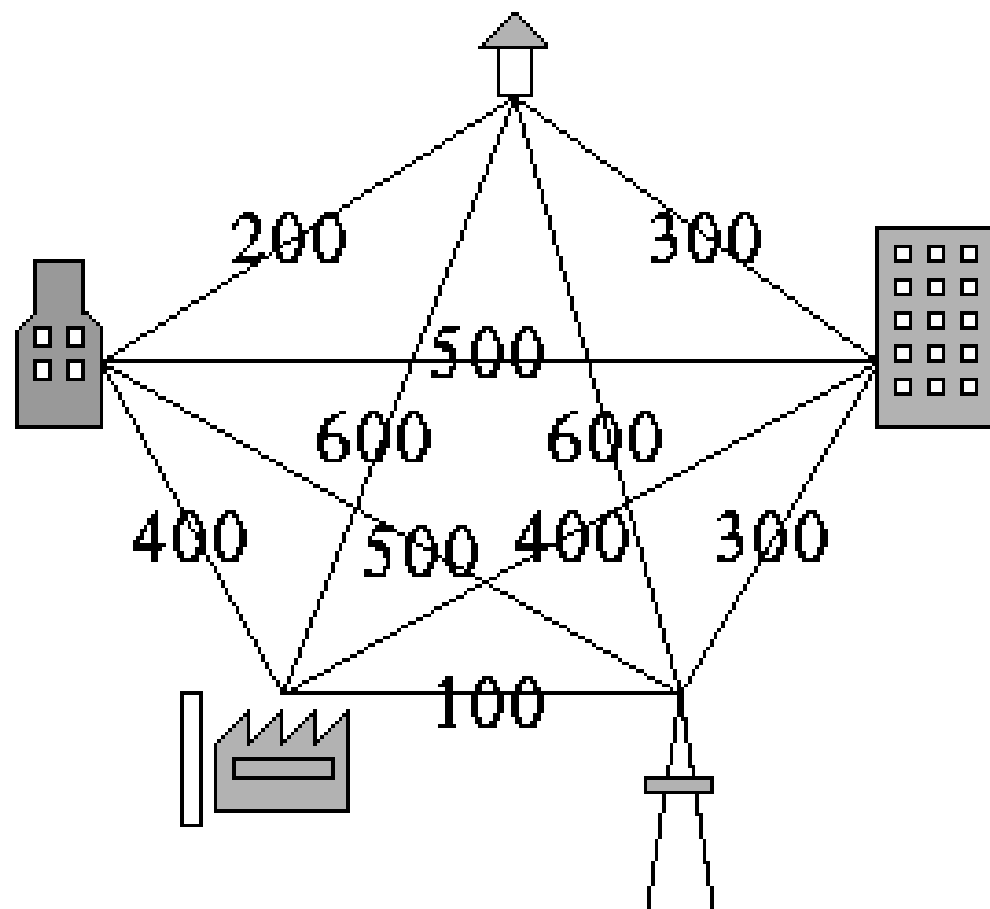




Y の値



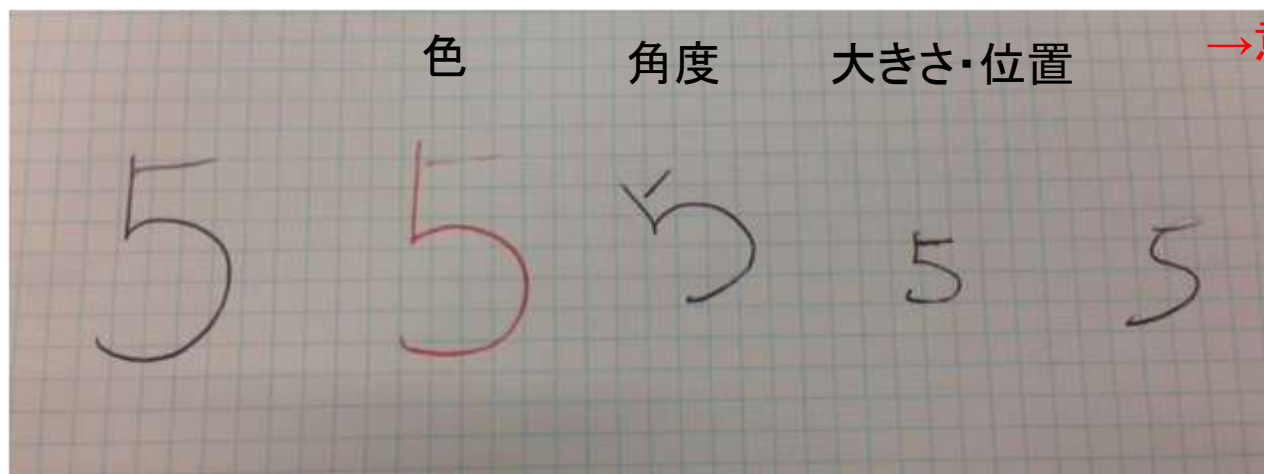
- ・ 膨大な組み合わせの中から最適な解を見つけ出す。
  - 全ての街を訪れること。
  - 同じルートは2度通れない。
  - 最も安いルートを通ること。



わずか30箇所、  
10の30乗！！！！

全ての組み合わせ数  
=  $n! / 2n$  通り  
5箇所 = 12通り  
10箇所 = 181440  
30箇所 = 10の30乗

- 分析したい事象の特徴を示すパラメータが必要
  - ・ 手書き文字の OCR の例
- 故障に関連するパラメータは？
  - ・ ワット数、電流、電圧、長さや重さ、気温、天気、湿度、温度、使用時間、利用者、使用方法・・・
  - ・ 手法や繰り返し分析することで、関連するパラメータを推定することが可能だが“あたり”をつける必要がある。



→意味の無い  
パラメータ

## ・ 異常検知

- 機器から出力されるパラメータを元に、異常な挙動を検知し、事故や故障を未然に防止する。
- 例) クレジットカードの不正利用

## ・ 予防保守

- いつ故障しそうか、過去のデータから推測する。
- 部品の交換の頻度や、生産ラインを休止する頻度を下げることができる。



# まとめ



- Hadoop は、安定化・高速化し、Web系から一般企業へ浸透し始めている
- オムニチャネル・コマースにあった、オムニチャネル・レコメンドを提案します。
- 機械学習の適用には業務知識が必須  
エクサには実現するノウハウがあります。

# Q&A

**ご静聴ありがとうございました。**

---