

2000年9月1日

SUJ J2000

# マイニング・ツールの比較評価と 選択のポイント

ベンダー 8社のマイニング・ツール評価

株式会社三和銀行

リテール統括部

小 野 潔

kiyosi\_ono@sanwabank.co.jp

03-5252-1518

(本報告は個人的見解である)

## 目 次

- ・ 現状のマイニング・ツール
  - ツールの分析手法と機能一覧
  - ツールの新しい手法と機能紹介
- ・ ツールの新たな方向性
  - マイニング理論&技法の最前線
  - 最強モデル AdaBoost とは？
- ・ まとめ
  - マイニング・ツールの選択の要

## ツールはマイニングに必要不可欠

- ・ 作業効率 は 3 ~ 5 倍以上
- ・ 日本では 20 種類、米国では 70 種類以上
- ・ 問題点 :
  - 大規模ソフトウェアのため 価格が高い
  - ベンダーごとに 機能や分析手法が異なる
  - 選択を誤れば無用の時間と費用の浪費となる
  - すべてを兼ね備えたツールは存在しない
  - ユーザー本位の「ツール比較レポート」がない

**結局、ユーザーは「価格」と「分析手法」のみで選択**

## ベンダーと製品名

販売ベンダー	本 社	会 社	製 品 名	Ver	発売日	機 能
東 芝	日本	メーカー	KINOsuite-P/R	1.0	99.04	ルール抽出機能付きキューロ
			KINOsuite-IDTF	1.0	99.04	ファジィ決定木
日本エシス	米国	メーカー	MiningPro21 (開発は日本)	1.0	99.10	マイニングエンジン
			分析テンプレート	1.0	99.10	マイニングナビ
日本IBM	米国	メーカー	Intelligent Miner	6.1	99.10	マイニングエンジン
			Relationship Marketing	6.1	99.10	マイニングナビ
日本SGI	米国	メーカー	MineSet	3.0	99.07	マイニングエンジン
日立製作所	日本	メーカー	DATAFRONT	2.01	00.03	マイニングエンジン
富士通	日本	メーカー	SymfoWARE Parallel Mining Server	4.0	99.12	マイニングエンジン
			SymfoWARE VisualMiner	4.0	99.12	平行座標ソフト
S A S	米国	統計ソフト	Enterprise Miner	2.02	98.07	マイニングエンジン
S P S S	米国	統計ソフト	Clementine	5.2	99.05	マイニングエンジン
			AnswerTree	2.1	97	決定木
			Neural Connection	2.0	95	ニューラルネットワーク
			SmartScore	1.0	99.06	開発支援ツール

## 分析手法

分類	分析手法	東芝	日本エニックス	日本IBM	日本SGI	日立	富士通	SAS	SPSS
決定木	C4.5	-	-	-	-	-	-	-	-
決定木	C5.0	-	-	-	-	-	-	-	-
決定木	CART	-	-	-	-	-	-	-	-
決定木	CHAID	-	-	-	-	-	-	-	-
決定木	QUEST	-	-	-	-	-	-	-	-
決定木	Pseudo Decision Tree	-	-	-	-	-	-	-	-
決定木	Option Tree	-	-	-	-	-	-	-	-
決定木	ファジィ決定木	-	-	-	-	-	-	-	-
決定木	領域分割決定木	-	-	未定	-	-	-	-	-
決定木	その他機能拡張決定木	-	-	-	-	-	-	-	-
ニューロ	Back Propagation	-	次期	-	-	-	-	-	-
ニューロ	Radical Basis Function	-	-	-	-	-	-	-	-
ニューロ	ベイジアン・ネットワーク	-	-	-	-	-	-	-	-
クラスタリング	コホーネン・ネットワーク	-	次期	-	-	次期	-	次期	-
クラスタリング	K-means法	-	-	-	-	-	-	-	-
クラスタリング	Ward法	-	次期	-	-	-	-	-	-
クラスタリング	コンドルセの手法	-	-	-	-	-	-	-	-
クラスタリング	概念クラスタリング	次期	-	-	-	-	-	-	-
アソシエーション	Apriori	次期	-	-	-	-	-	-	-
アソシエーション	Generalized Rule Induction	-	-	-	-	-	-	-	-
アソシエーション	順序アソシエーション	-	-	-	-	-	-	-	-
アソシエーション	類似時系列パターン	-	-	-	-	-	-	-	-
ルール抽出	ニューラルネットワーク、MLP	-	-	-	-	-	-	-	-
ルール抽出	ルールインダクション	-	-	-	-	-	-	-	-
最小近傍法	MBR	-	-	-	-	-	-	次期	-
回帰分析	ロジスティック回帰分析	-	-	-	-	-	-	-	-
テキストマイニング	テキストマイニング	-	-	-	-	未定	-	次期	-
テキストマイニング	Concept Base Search	-	-	-	-	-	-	-	-

## 機能一覧

	東芝	日本エニックス	日本IBM	日本SGI	日立	富士通	SAS	SPSS
並列処理		-					-	-
UNIX版	-							
NT版								
LINUX版	-	-	-	次期	-	-	未定	-
メインフレーム版	-	-		-	-	-	未定	-
リフト図（モデル的中率の比較）	-	次期		-	未定	-		
説明属性の選択機能	-	-	-					
高度な多次元散布図	-	-	-					
KDDプロセスの可視化								
対費用効果を含めた収益計算	-	-			-	-		
初心者向けのナビゲータ機能	未定			-	次期	-	-	-
プロセスのC or XLM言語変換	-	-	-	-	-	-	次期	
レポート機能	-		-	-	-	-		
日本語のマニュアル本	-			-	-	-		
定期セミナーの開催	-			-	-	-		
専門の質問応答セクション								
価格（Aは300万円未満、Cは1000万円以上）	B	A	C	B	A	A	C	B

## ツールの新しい機能

### プロセスのC or XML言語変換機能

XMLはインターネットの標準言語になる

### 運用形態

クライアント・サーバー型

シングル・アローン型

### インターフェース用開発支援ツール

モデル更新に伴うインターフェース画面の変更を自動的にこなす

### マイニング・ナビゲーター機能

初心者でもある程度のマイニングが可能

仕組はシステム部門が前もって加工データを用意

## 注目の分析手法と技法

- MBR(Memory Based Reasoning,記憶に基づく推論)
  - 最小近傍(Nearest Neighbor)法
  - モデルの更新が必要ない
  - 欠点 :ブラックデータが数千個必要
- Rough Sets
  - 欧米ではファジィ集合より使われている
- SVM(Support Vector Machines)
  - key word: カーネル法    マージン
  - AdaBoosting との関係
- Cross Validatin(=Jack knife)    精度を高める技法
  - k-fold cross-validation (leave-one-out,)
  - multiple cross-validation

# マイニング理論の最前線

- 分類器の組み合わせ

## Hybrid Model

Cataract Model (連続型)

Cascade Model (直列型)

Revise Model (補正型)

Flag Model (CARTを利用してカテゴリフラグ作成)

## Committee Model

Bagging (並列型)

AdaBoosting (モデルの重み付け結合型)

ECOC (Error Correcting Output Codes)

- SAS の Ensemble Model の構成

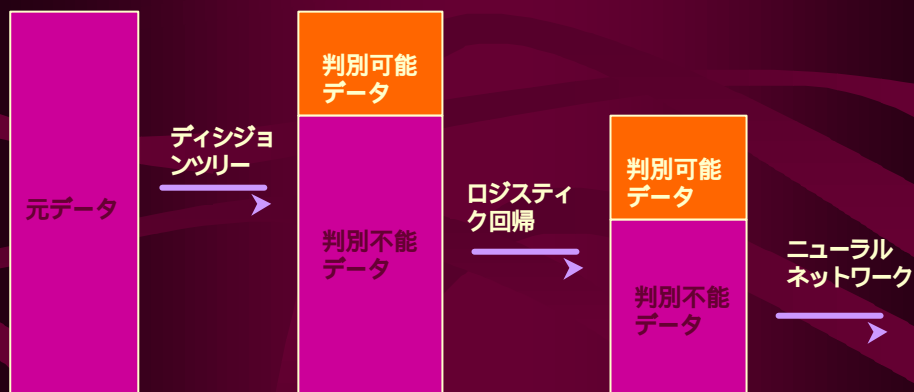
Combined Model

Stratified Model

Bagging Model

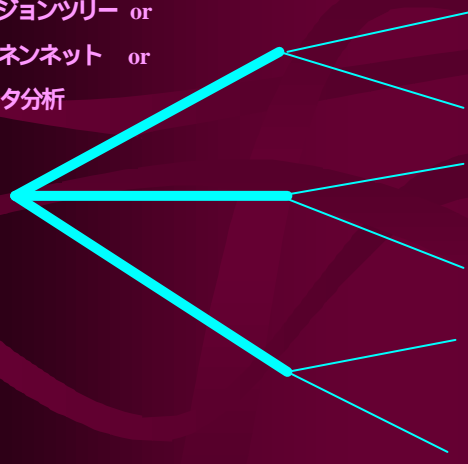
Boosting Model

## Cataract Model



## Cascade Model

First Step  
ディビジョンツリー or  
コホーネンネット or  
クラスタ分析



セグメント1  
セグメント2  
セグメント3  
セグメント4  
セグメント5  
セグメント6

Second Step  
ロジスティック回帰モデル  
ロジスティック回帰モデル  
ニューラルネットワーク  
ニューラルネットワーク  
ラフ集合  
OK

## Revise Model

FIRST STEP

決定木 (n個説明変数)



セグメント毎の確信度

SECOND STEP

ロジスティック回帰

$$y = a_1x_1 + a_2x_2 + \dots + a_nx_n + a_{n+1}x_{n+1}$$

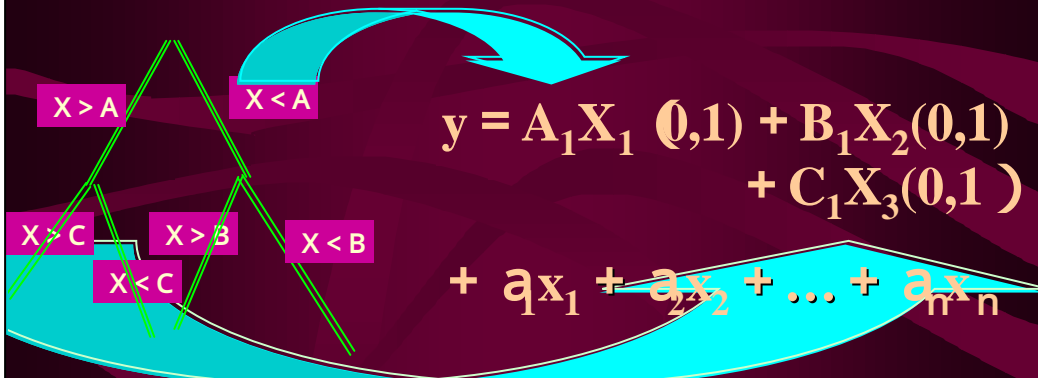




# Flag Model

FIRST STEP

SECOND STEP



1 説明変数の決定木

但し  $X$  は 2 値 (0 or 1)

# Committee Model

	ディジニツリー	ニューラルネットワーク	ロジスティック回帰
顧客 A	90%	95%	70%
顧客 B	15%	10%	11%
顧客 C	70%	50%	40%
顧客 D	60%	65%	55%
顧客 E	10%	15%	20%



多数決、線形結合、最小値、最大値、平均 等

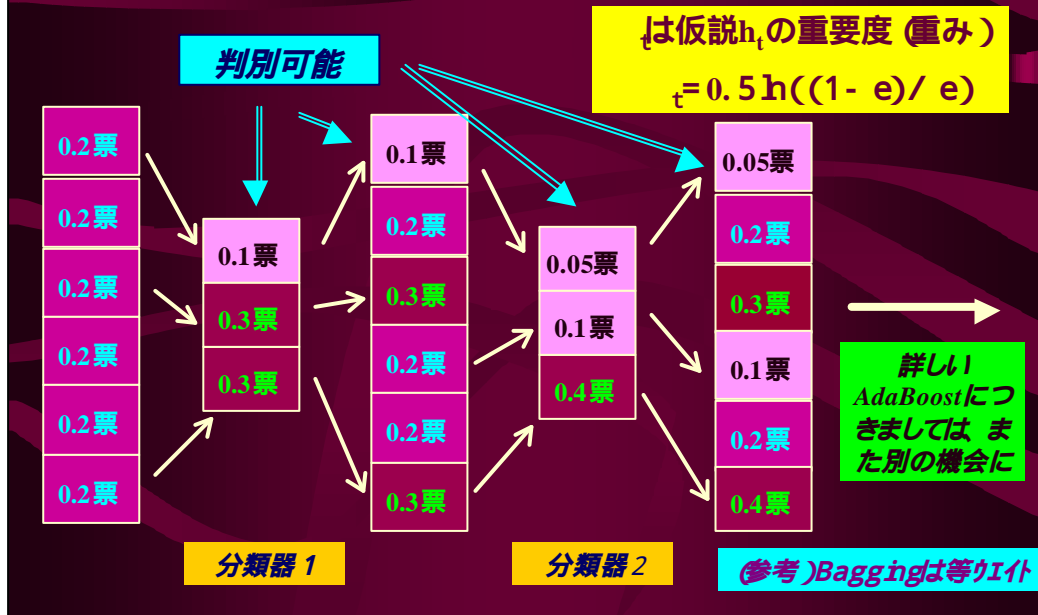
	総合点
顧客 A	70%
顧客 B	12%
顧客 C	50%
顧客 D	55%
顧客 E	17%

同じ分類器



Bagging, Boosting

## AdaBoostの仕組み



## 最強の AdaBoost

- 精度の低い弱学習分類器を精度の高い強学習分類器にするためのアルゴリズム
  - 例えば、10 ~ 100個の複数の分類器の多数決
  - 簡単に言うと、AdaBoostは分類機のターボエンジン
- 今、最も米国の学会ではホットな話題
  - 数学的な裏付け?
  - 過学習が起きない?
- エンジン (分類器) によって最適ターボが違う
  - 現在、多値、Local Classification、雑音、データ不足に対応した新しい AdaBoost モデルが現われた



## ツールの選定基準 (最後に)

### 1.マイニングの成功は、良いデータ分析者・ツールが必要不可欠

マイニング・ツールは未発達である  
最初はベンダーによるコンサルトが望ましい

### 2.ツールの決定要素は？

使用目的、ベンダーのコンサルタント能力、運用形態、  
分析者の素質等

### 3.ユーザーは新しい分析手法や機能を絶えずチェック

これからツールには、Bagging, AdaBoostが必要  
SAS/EM ver4.0 は頭一つリードしているようだが  
他社ツールの追従確実

## Question & Answer

比較一覧



AdaBoost



最新理論

