

2003年度 統計関連学会 連合大会
企画セッション「データとその属性情報の記述」

DandDインスタンス



慶應大・理工 島津秀康



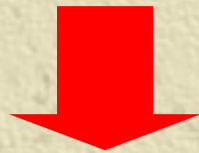
DandDインスタンスとは

✧ XML (eXtensible Markup Language) 文書

- ✦ タグを使った記述

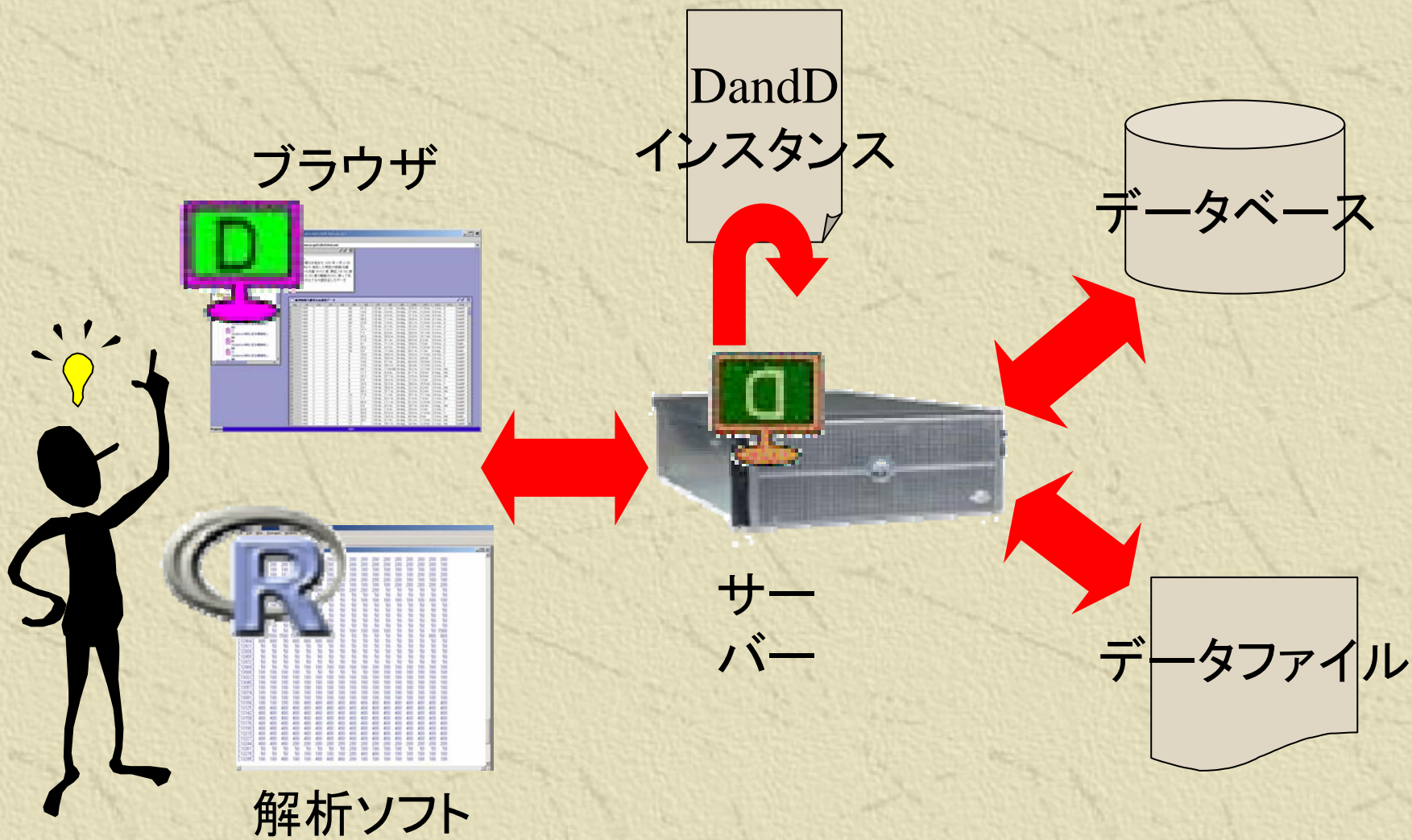
✧ インスタンスは単に「記述」

- ✦ データの所在, データベクトル, 構造...



記述を解釈するのは支援システム

DandDの全体像



従来の記述

✦ Relational (関係形式)

- ◆ ベクトル(変量)を並べたもの

✦ 正規形

◆ 第1正規形

- すべての要素がアトミック

◆ 第2正規形

- いくつかの項目でレコードを一意に定められる

◆ 第3正規形

- レコードを一意に定める項目以外の項目は互いに独立

ベクトル間の関係は不明

DandDでのデータ表現

✦ Relational (関係形式)

データベクトルの並び

◆ TimeSeries (時系列)

- 時間, 値を示すベクトルの並び

◆ PointProcess (点過程)

- 時刻, 場所, 値を示すベクトルの並び

✦ Array (配列形式)

軸とデータベクトルの並び

✦ 基数系

時間 → 時, 分, 秒を示すベクトルの並び

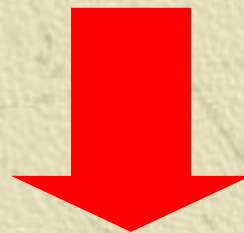
データベクトル間の関係でデータを表現

従来の記述との違い

✦ ベクトル間の**関係**に基づく組織化

✦ 組織化に人間の意図が反映される

View



データの解析に直結した表現

解析におけるデータの流れ

- ✦ 取得
- ✦ クリーニング, 記述
- ✦ ブラウジング
- ✦ 解析
- ✦ モデリング
- ✦ モデルの評価

試行錯誤を重ねるうち
データの見方 (View) が
変化

変化をインスタンスに記述

水泳競技データ

✦ (財)日本水泳連盟医・科学委員会

✦ 1996年～2003年

✦ 13329レコード

✦ 約100項目

- ◆ 大会名, 種目

- ◆ 選手名

- ◆ ストロークタイム

- ◆ ラップタイム など

取得

- ✦ 会場に設置されたビデオカメラより取得
- ✦ 手入力

<理由>

- ✦ 1選手1競技の結果を1行で管理
- ✦ 様々な競技を1つのテーブルで管理
- ✦ 更新が容易

クリーニング

✦ 2次データの混在

◆ ストローク頻度 = 60/ストロークタイム

➡ 1次データのみを記述

✦ 個人識別

➡ 新たなIDベクトルを作成

記述

✦ 意味を捉えた分かりやすい形へ (テーブルの分割)

◆ 競技情報

- 競技ID, 年, 月, 日, 大会名, 性別, 種目, 距離...

◆ 個人情報

- 個人ID, 氏名

◆ 時間情報

- 競技ID, 個人ID, 観測ポイント, ラップタイム

◆ ストローク情報

- 競技ID, 個人ID, 観測ポイント, ストロークタイム

解析へ向けて

✦ 注目する変量

- ◆ ストロークタイム

- ◆ ラップタイム

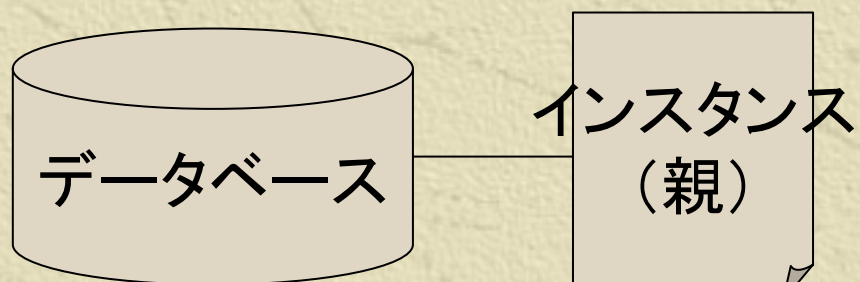
✦ 注目する水準

- ◆ 男子, 200m, 自由形

変化の記述 (Relatives)

-
- ✦ 変化するのはデータの見方 (**View**)
 - ✦ データベクトルの継承の関係を記述

インスタンスのネットワーク



<Relatives>

<Child Id=" " URL=" ">

<Inherit Id=" ">

継承させるベクトルの作成手順

</Inherit>

</Relatives>



<Relatives>

<Parent Id=" " URL=" ">

<Inherit Id="継承するベクトル" />

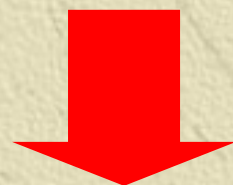
</Relatives>

まとめ

✦ ベクトル間の関係に注目する組織化により、
多様なデータを表現

- 元のデータを書き換えることなく、インスタンスの書き換だけでデータの見方(**View**)を自由に換えられる

✦ 特定のソフトに依存しない利用環境



高い汎用性