# リレーショナルデータベースの 高度利用環境DandD

慶應義塾大学理工学部 横内大介 慶應義塾大学理工学部 柴田里程

## データの高度利用

- データの公開
  - -情報公開
  - インターネットの普及、発展
  - データの配布形式、組織化の多様性
- 具体例
  - インターネットによる配布
    - LIBOR (London InterBank Offered Rate)
    - 地震データ
  - CD-ROMによる頒布
    - 気象データ(気象業務支援センター)
    - 民力(朝日新聞社)

## データの高度利用を阻む障壁

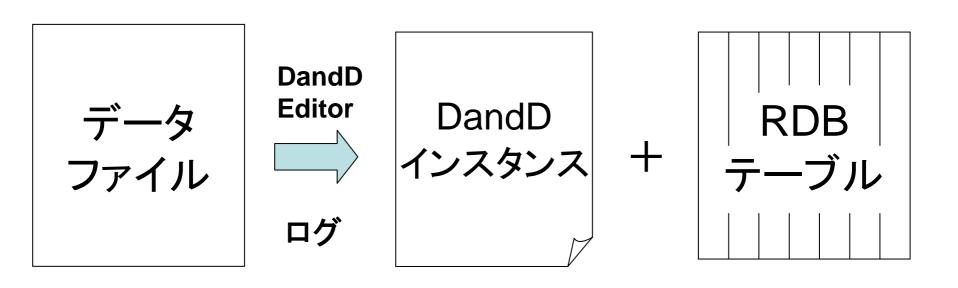
- 分割,分散
  - 公開側の都合
- 形式の多様性
  - 統一基準がない
- 規模
  - 情報技術の向上
- 信頼性
  - 保証する仕掛けがない
- 記述の曖昧さ
  - 客観的に評価するシステム
- 公開の原則の欠如
  - 風土



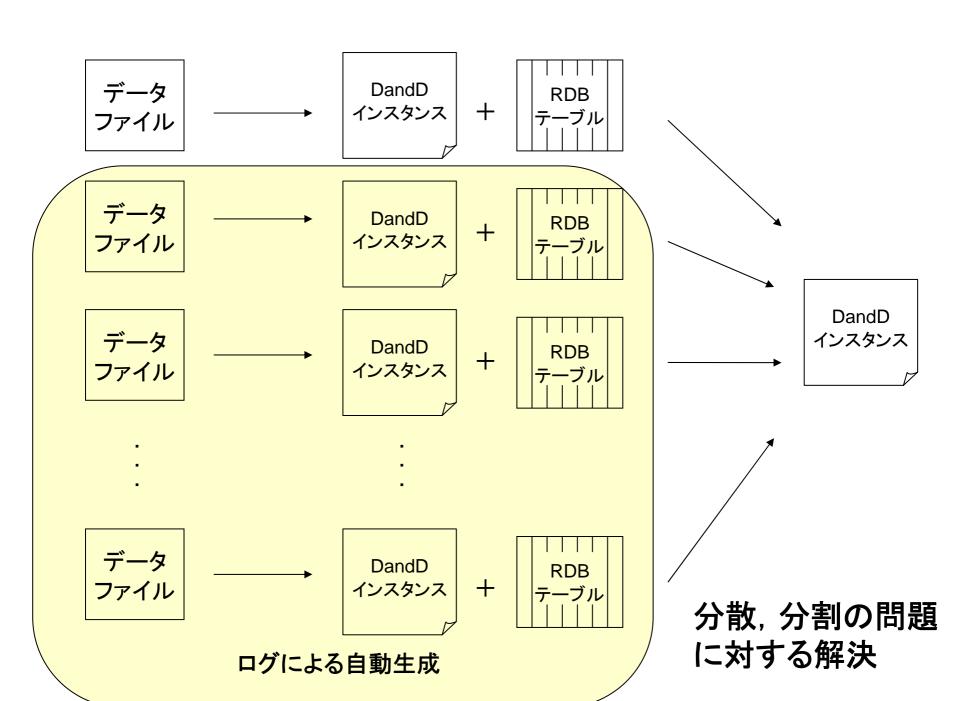
### DandD環境による 5つの障壁の除去

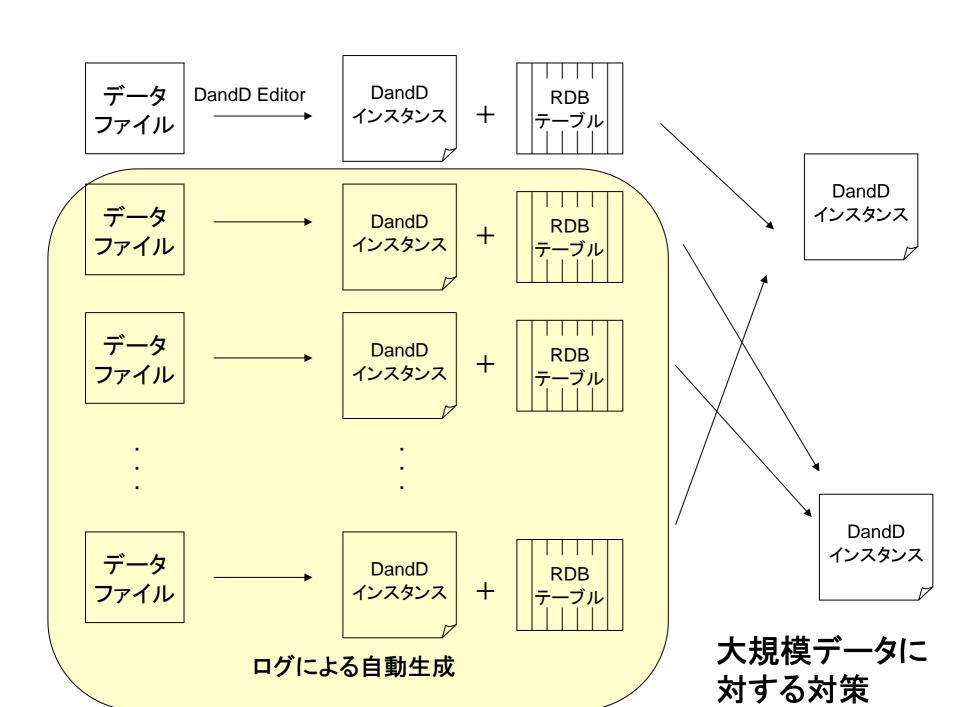
- DandDインスタンス
- リレーショナルデータベース

## DandDインスタンスとRDB



- 形式の多様性の吸収
- 信頼性の保証
- 形式的な記述

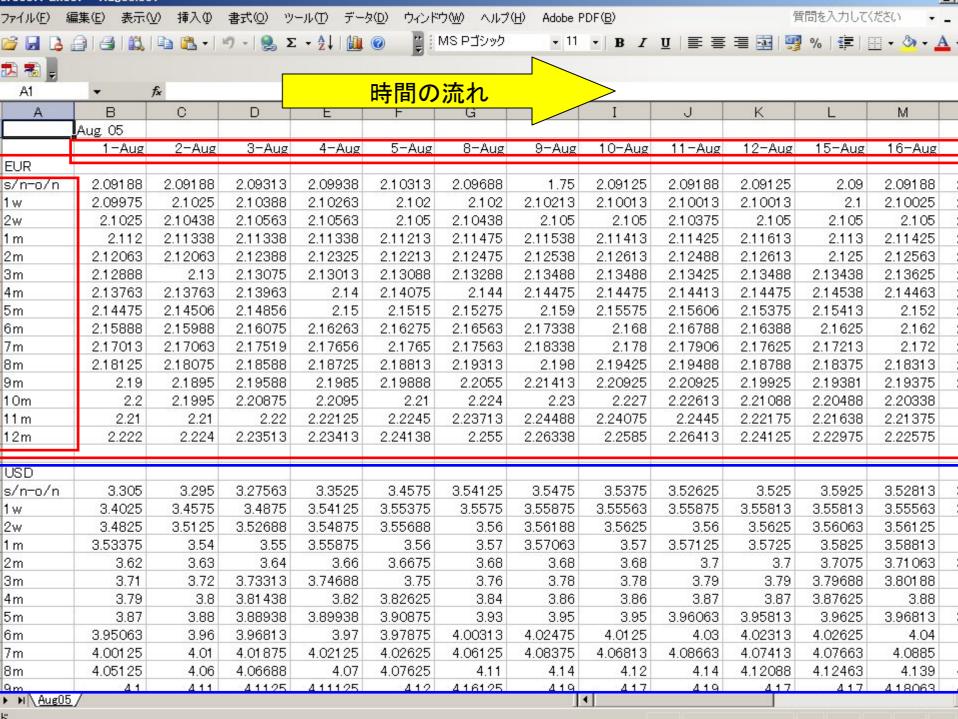




## 実際例

#### LIBOR

- ロンドン市場での銀行間平均貸し手金利
- BBA(The British Banker's Association)が日に 一度発表
- 1ヶ月単位のExcelシート
- シート 1 枚に9種類の金利データ
- 構造が複雑
- モデルケース
  - 2005年1月から8月までのEURO金利を併合しR 上で利用する



# 表計算ソフトウェアとRを用いた一般的な作業工程

#### 表計算ソフト

- データの切り出し
- データの転置
- 複数のシートを結合
- 適切なヘッダーの付与
- R への取り込み

#### DandD環境

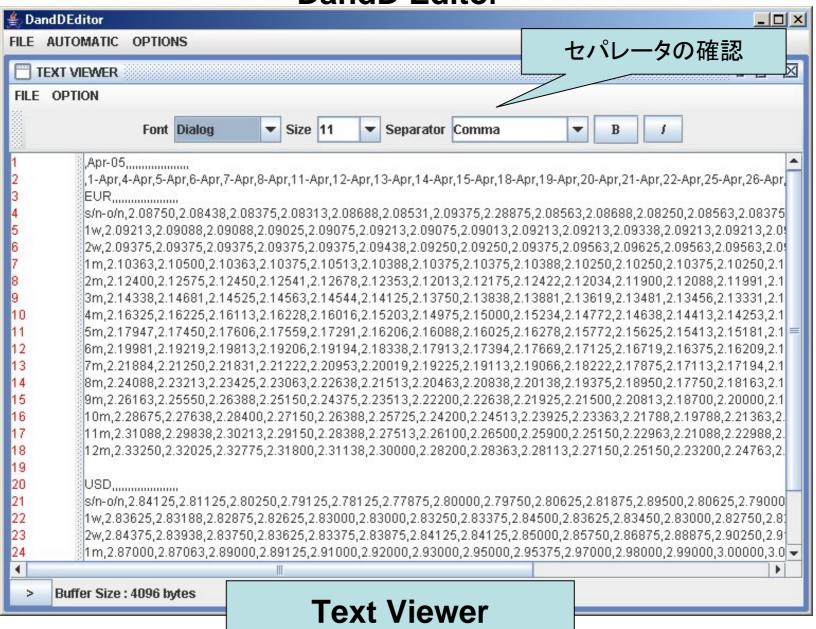
#### **DandD**

- データの切り出し
- データの転置
- DandDインスタンスの生成
- R への取り込み

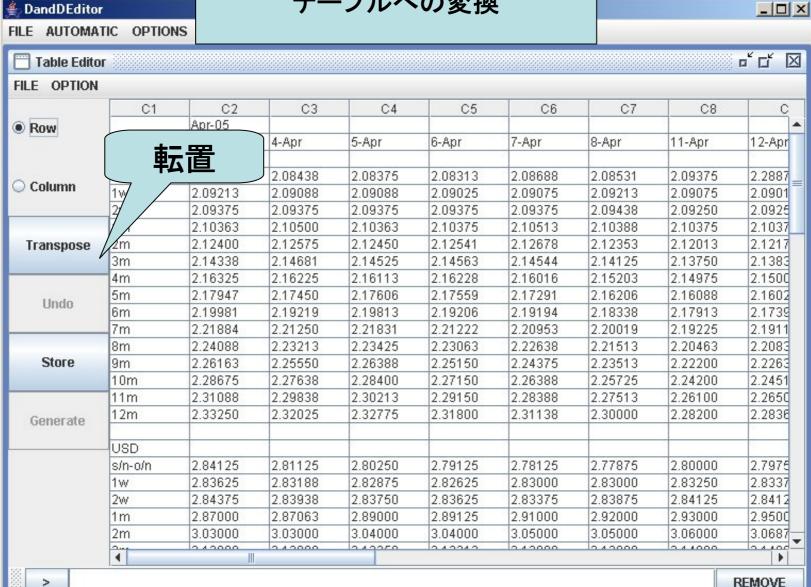


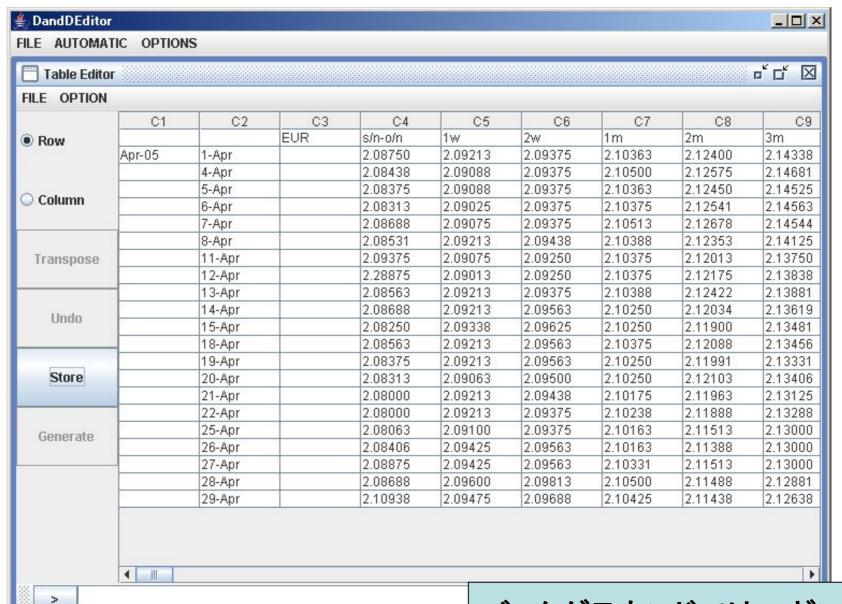
DandDEditor ≥ DandDR

**DandD Editor** 



#### テーブルへの変換

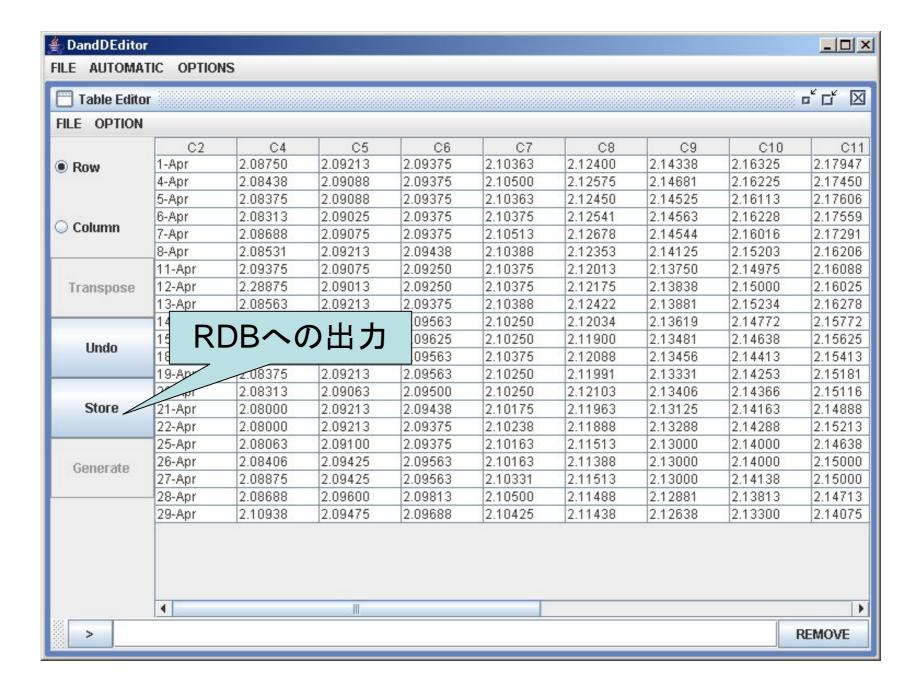


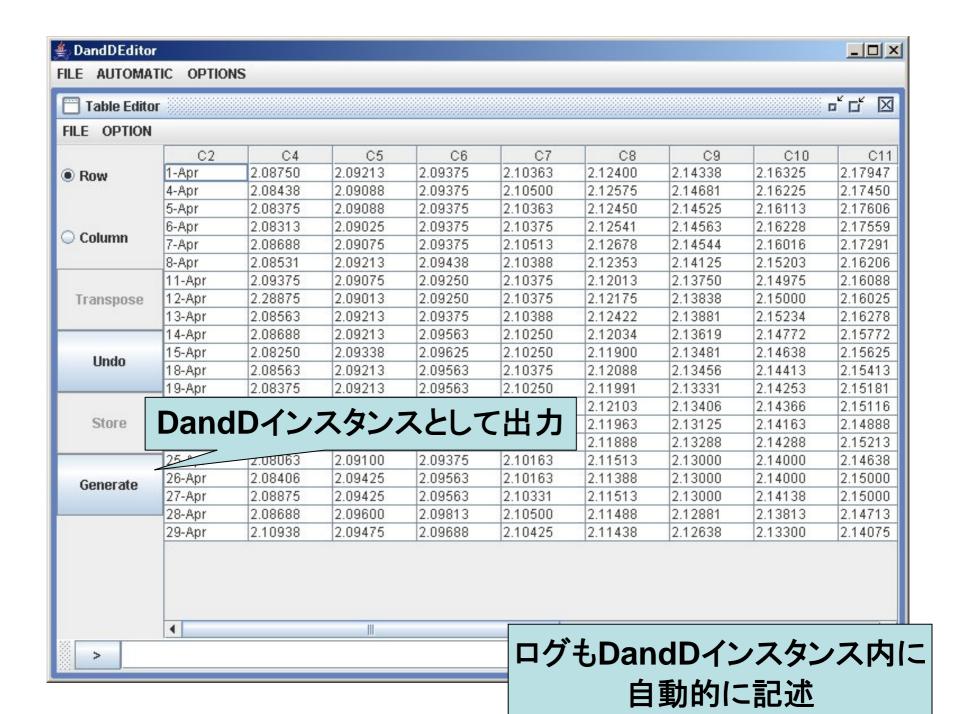


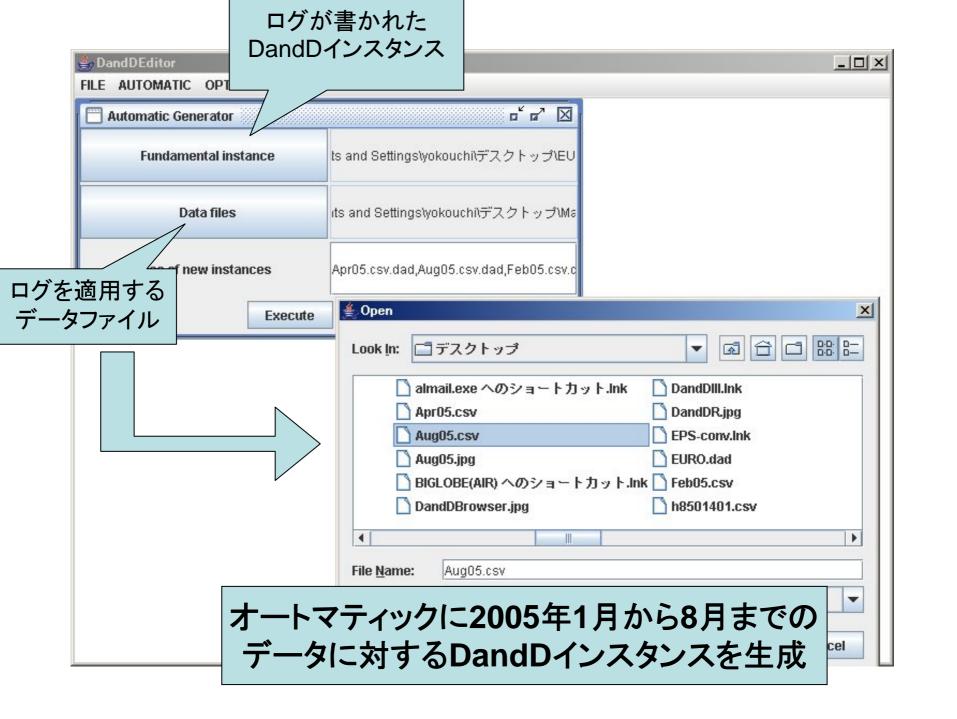
バックグラウンドではロギング

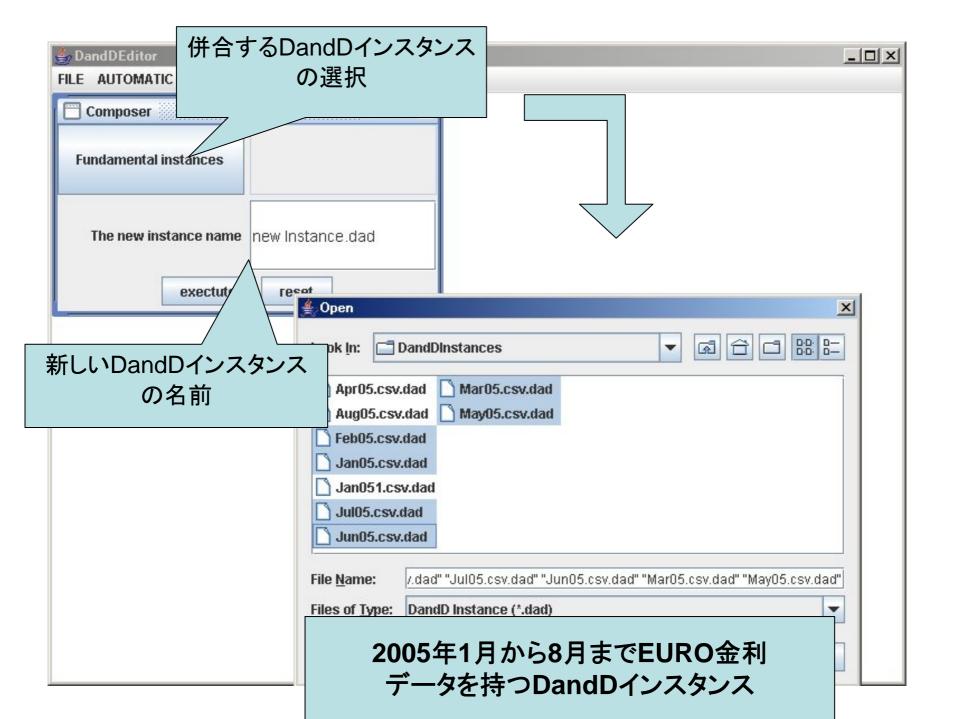


バックグラウンドではロギング

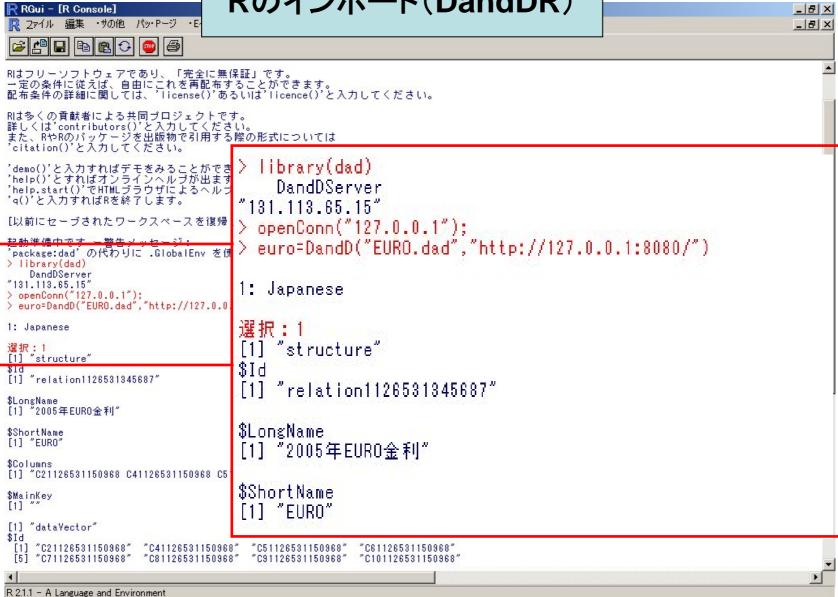




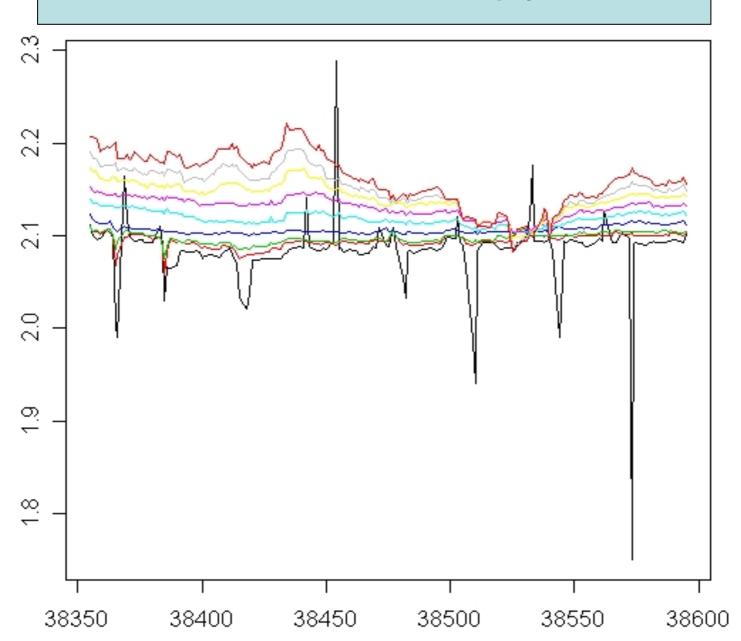




#### Rのインポート(DandDR)



### RによるEURO金利の時系列図



# 表計算ソフトウェアとRを用いた一般的な作業工程

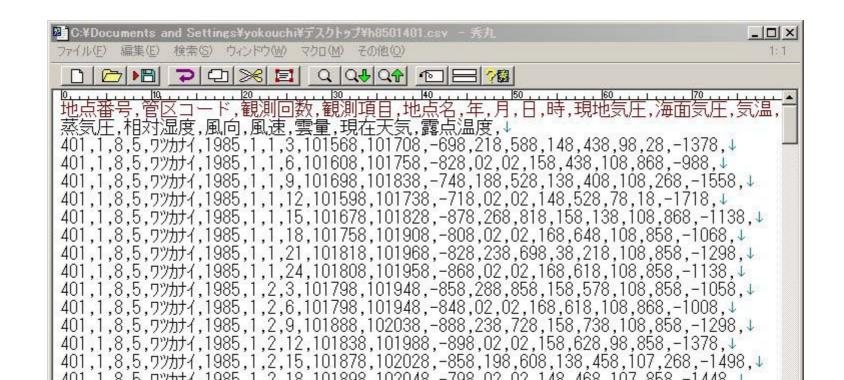
- •手作業中心
  - •表計算ソフトウェアや Rの十分な知識
- 編集作業者依存
  - •妥当性
  - インフォーマルな編 集記録

#### DandD環境

- •DandDインスタンス化中心
  - •DandDインスタンスにする手続きのみに特化
- オートマチックなログ機能
  - ●信頼性の確保
  - •編集手続きの自動化

## 実際例2

- 気象データの利用
  - 1961年から2003年まで
  - 144観測所の月別データ
  - 観測間隔3時間ごと
  - 約61000ファイル, 約 1.16 GB



### 表計算ソフトウェアとRを用いた 一般的な作業工程

- データの編集
- 適切なヘッダー
- 複数のシートを結合
- R への取り込み



- •大規模データに対する限界
  - •表計算ソフト
    - ●一部だけ
  - ●デフォルト設定のR
    - ●読み込めるがおよそ1年分 のデータで動作が不安定

#### DandD環境

- データの編集
- DandDインスタンスの生成
- R への取り込み(DandDR)



- 大規模データへの対応
  - •従来のDandD
    - RDB を背景に大規模データのストレージには対応
    - 無理やり1つのDandDインスタンスにまとめていた
  - 現在のDandD
    - •複数のDandDインスタンス
    - •必要な段階でまとめる

## リレーショナルデータベース+DandD

- データの高度利用
  - 表計算ソフトウェアよ, さようなら
    - 報告書に載せる表の作成が基本
      - モデルを創ることまでは考えていない
    - 手作業
    - 大規模データは扱えない
    - 編集手続き
  - DandDよ, こんにちは
    - 大規模データの扱い
      - 表計算ソフトやRだけでは限界
      - データのストレージはすべてRDB
    - Audit(監査)
      - フォーマル
      - 容易な編集の実現

DandDプロジェクトのホームページ

http://www.stat.math.keio.ac.jp/DandD/

## 必須属性(人間の手を煩わすか否か)

- DataVector(最小単位)の必須属性
  - データベクトル 1 本の外形を知る上で最低限必要なもの
  - 将来にわたって変更する必要性がないもの(言語の追加を除いて)
  - 個々の値に関わらない属性
- 構造(Relation) の必須属性
  - 1 つのRelationを外形を知る上で最低限必要なもの
  - 将来にわたって変更する必要性がないもの(言語の追加を除いて)
  - 個々のカラムの役割や結びつきに関わらない属性
- インスタンスの必須属性
  - インスタンス(データ全体)の外形を知る上で最低限必要なもの
  - 将来にわたって変更する必要性がないもの(言語の追加を除いて)
  - 個々のデータ構造に関わらない属性