

データの価値を究めるTRAD

データサイエンスコンソーシアム, 慶應義塾大学

柴田 里程

データの価値評価

- 手元に蓄積したデータ
 - どれだけの利用価値があるのか
 - どう活用できるか
 - だれが活用できるのか
- データの流通
 - 適切な価格
 - 需要先
 - マーケット

データ解析の一つの目的はこのような価値を見極めること

TRAD (TextilePlot, R and DandD)

- データサイエンスコンソーシアムの長年の蓄積と大勢の努力の賜物
- データサイエンスの基礎理論と実践の反映
- 業務使用に耐えるだけの完成度
- データサイエンスの健全な発展を願って無償で公開 (Windows, Mac)
 - <http://datascience.jp>
 - 日々アップデート, なにか気づいたことがあれば
query@datascience.jp までご一報のほど.


Visual Excel

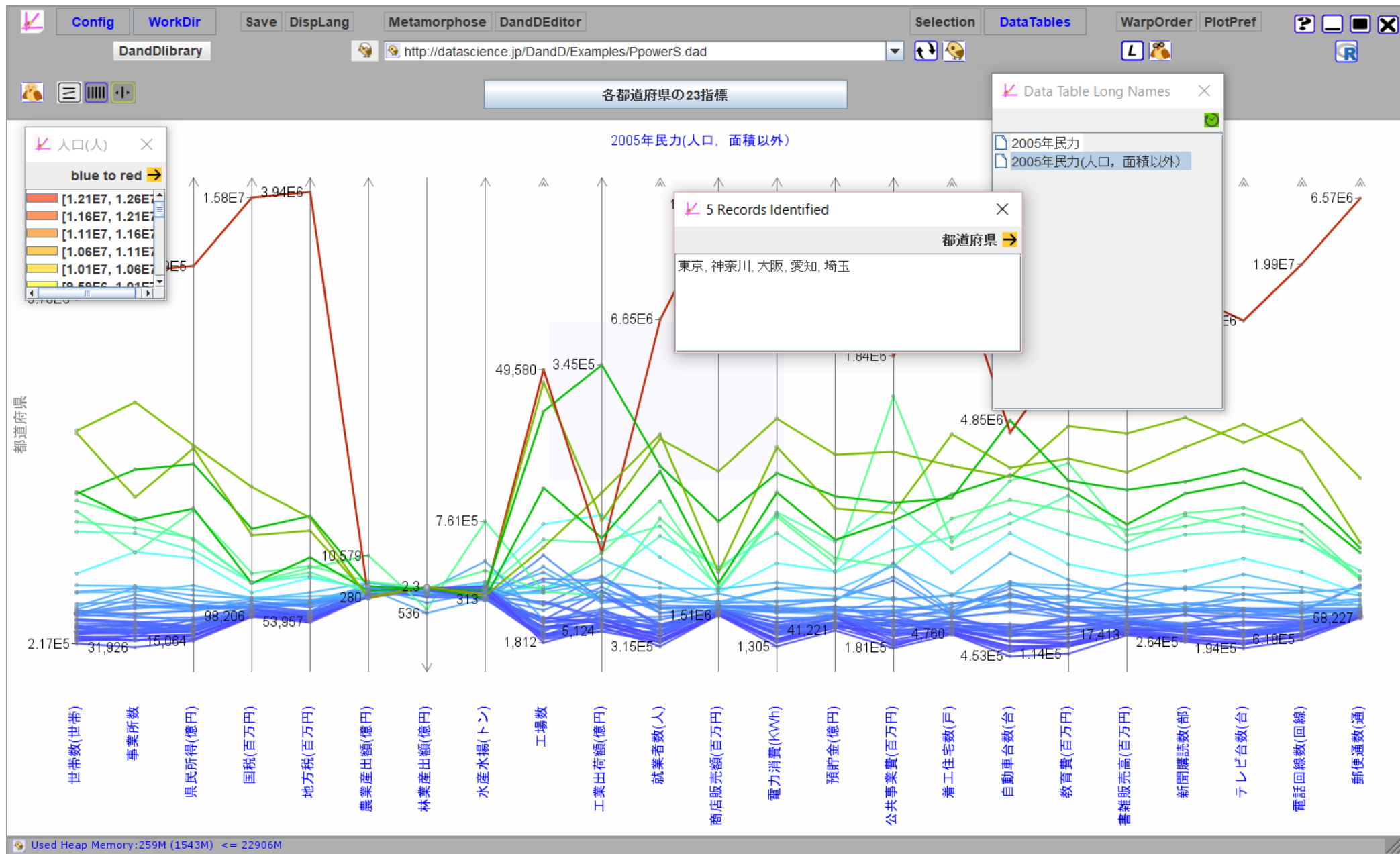
- 表計算ソフト
 - 簡単につかえ便利
 - 数字や文字の羅列を眺めていても, どう解析したらよいのか何のアイディアも浮かばない
 - とりあえずは, いろいろな図表を作ってみる(図表の海に溺れる)
 - 適当な方法を適用して結果らしきものを出す(自信, 確信はない)
- R
 - 駆使できるようになるにはかなりの修練が必要
 - 修練を積んでも手間がかかることには変わらない
 - データの全体像をつかむには手間と時間がかかる
 - データを眺めるというより, データを扱うための高度な機能を備えたソフトウェア

一年間の TRAD 研究開発の軌跡

- インストールの簡便化
- カラーリング
- Logicals to Mark
- Projections
- 高速化とメモリー使用量の削減 (Excel と同等かそれ以上の性能)
- Parallel Coordinate Plot の導入
- 線形モデル当てはめ結果の Parallel Coordinate Plot による表示

インストール

- あらかじめ Java Runtime Environment をインストールしておく(無償)
- TRADのインストールはワンクリック (Windows, MacOS)
 - 40MB 程度
 - 本体は約1MB, ソースコードで約 5 万行
 - 立ち上げ時メモリー消費量は約 100MB
- すでに R がインストールされていれば, 自動的に連携
- アイコン  をクリックするだけで起動
- オフラインマニュアルも電子ブック(ePub) 形式で提供



カラーリング (HSBモデル)

- 各記録あるいは水準に色相(Hue) を割り当てる
 - ハイライトと同じようにマウス操作で割り当てる
 - 割り当ての基準となるデータベクトルを選ぶ
- Lint (2つの軸に挟まれた部分での Weft の断片)
 - 画面の解像度に合わせて Thinning を行う段階で色相を合成する
 - ハイライトは彩度(Saturation)の違いで表現する
 - 基本の彩度や明度(Brightness) はユーザが自由に変更できる
- 描画は画面の下側から行う. デフォルトの不透明度は90 %
- 軸上のノードの色も同様な手順で行う. 必ずしも Lint 色と同等になるとは限らない

✖ Weft Color

Normal Weft

Saturation (%):

Brightness (%):

Highlight Weft

Saturation (%):

Brightness (%):

Other Color Settings

Background Whiteness (%):

Weft Opacity (%):

Logicals to Mark

- 複数選択可のアンケート回答
 - 選択したかどうかをあらわす TRUE (1) あるいは FALSE(0) の列
 - このままでは疎でありすぎるので, Mark 型に直したい
 - しかも回答項目が多いのでシステムティックに直したい

×

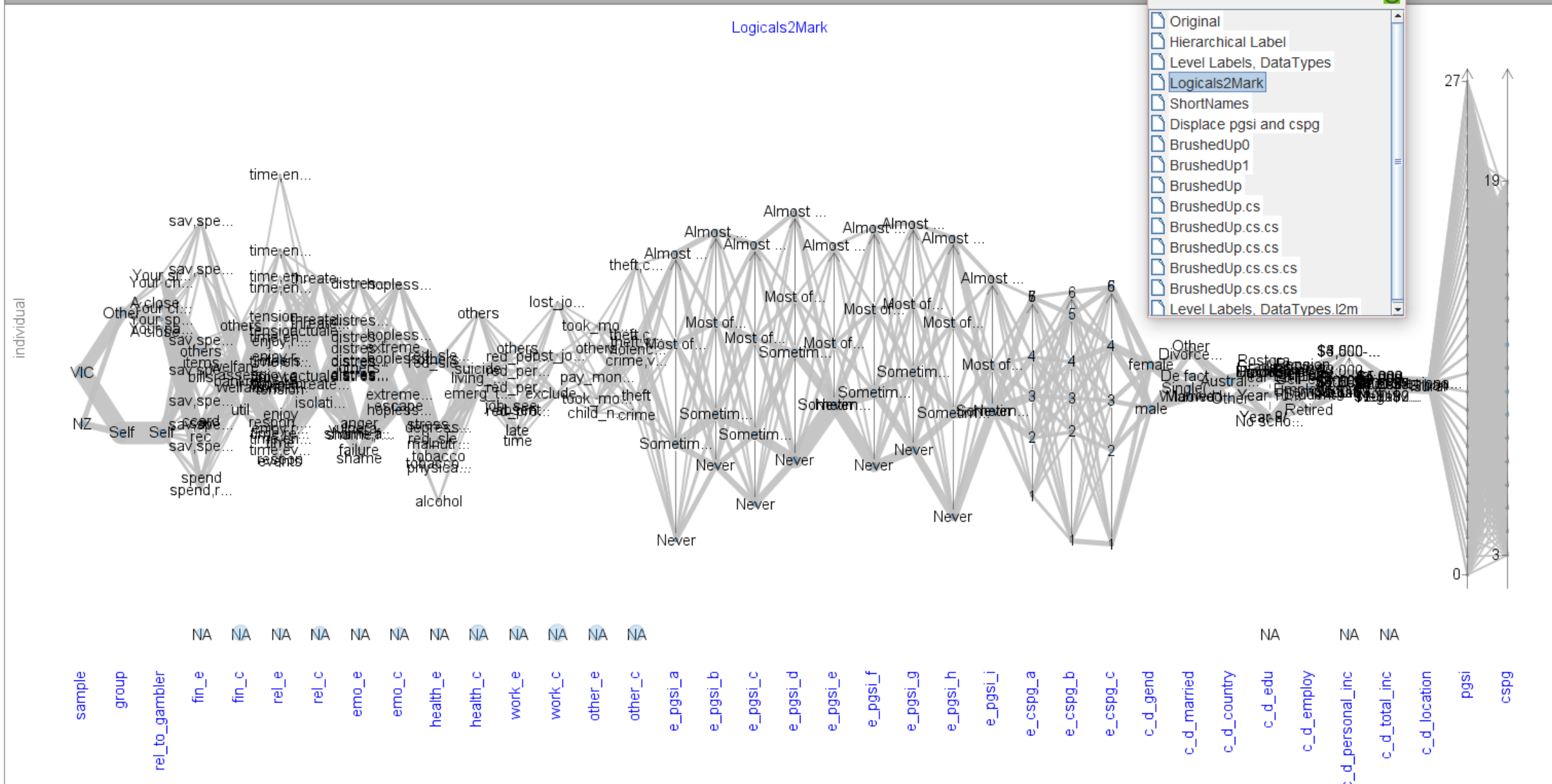
 γ

OK


×







Data Table Long Names X



Original

Hierarchical Label

Level Labels, DataTypes

Logicals2Mark

ShortNames

Displace pgsi and cspg

BrushedUp0

BrushedUp1

BrushedUp

BrushedUp.cs

BrushedUp.cs.cs

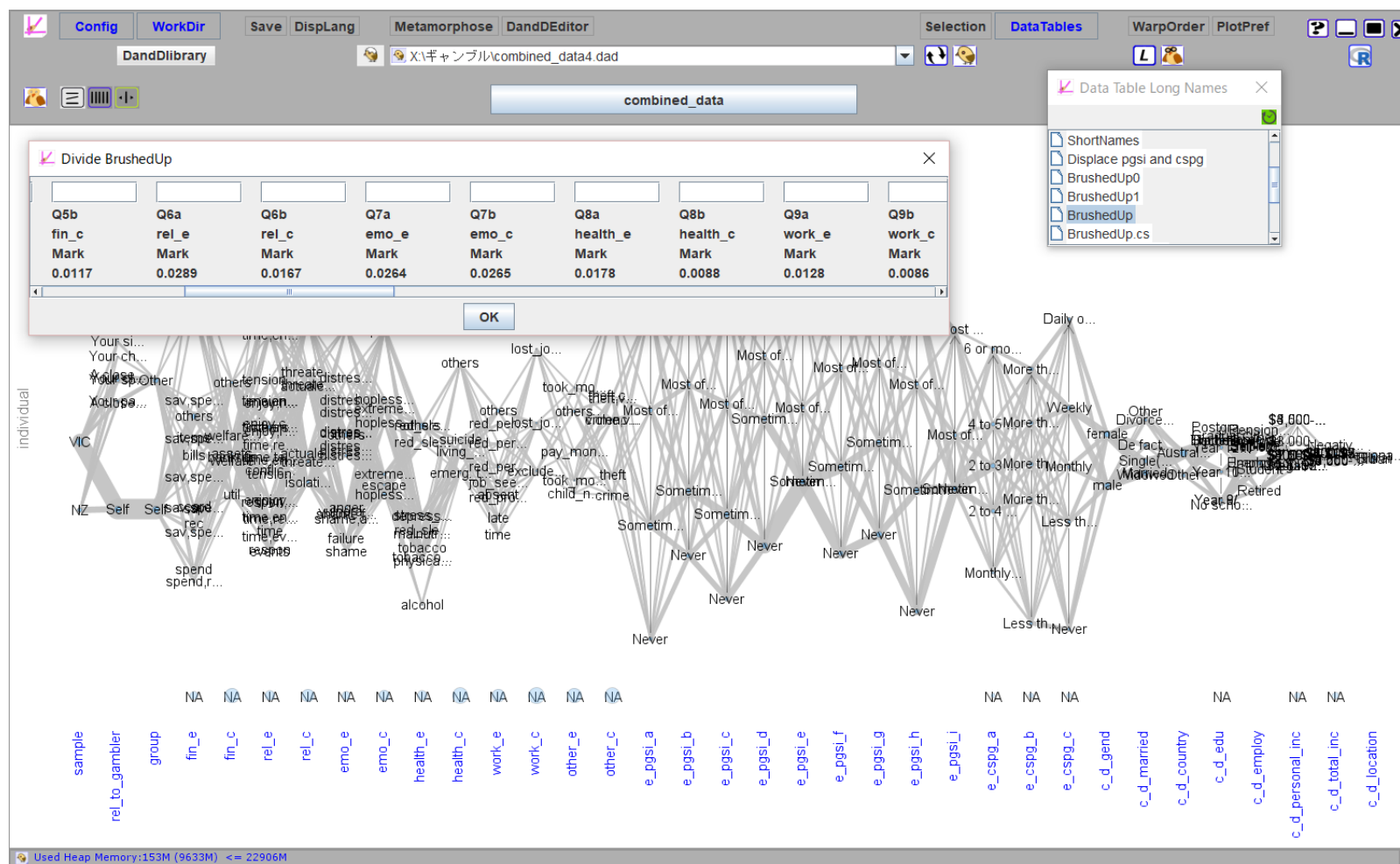
BrushedUp.cs.cs.cs

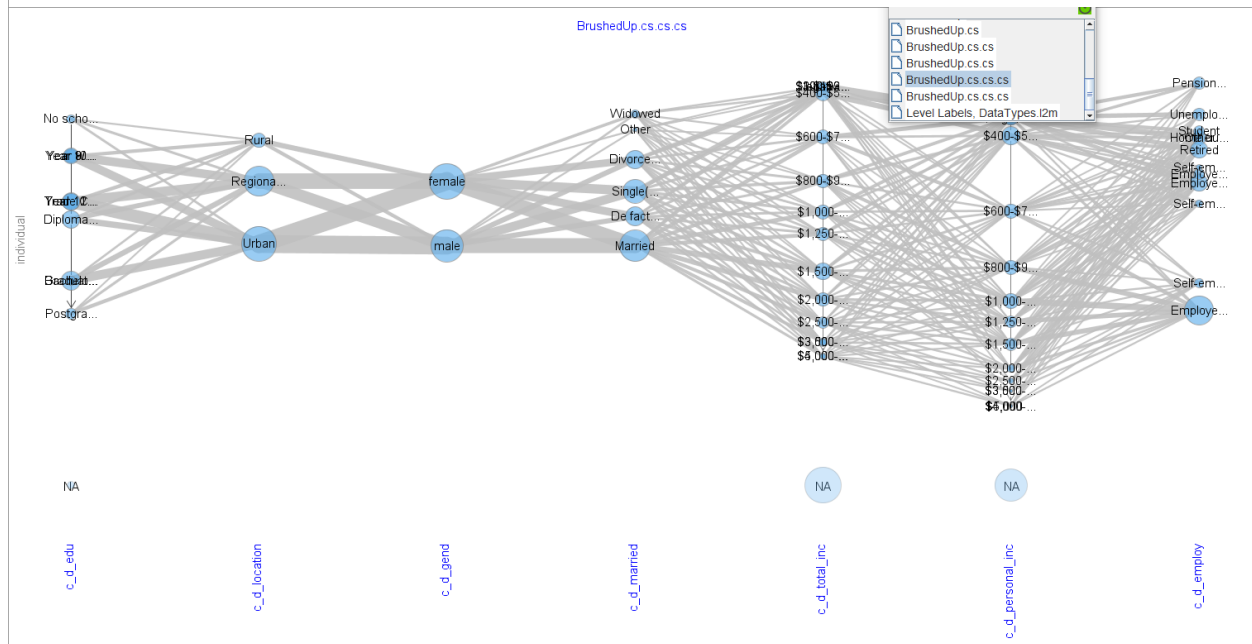
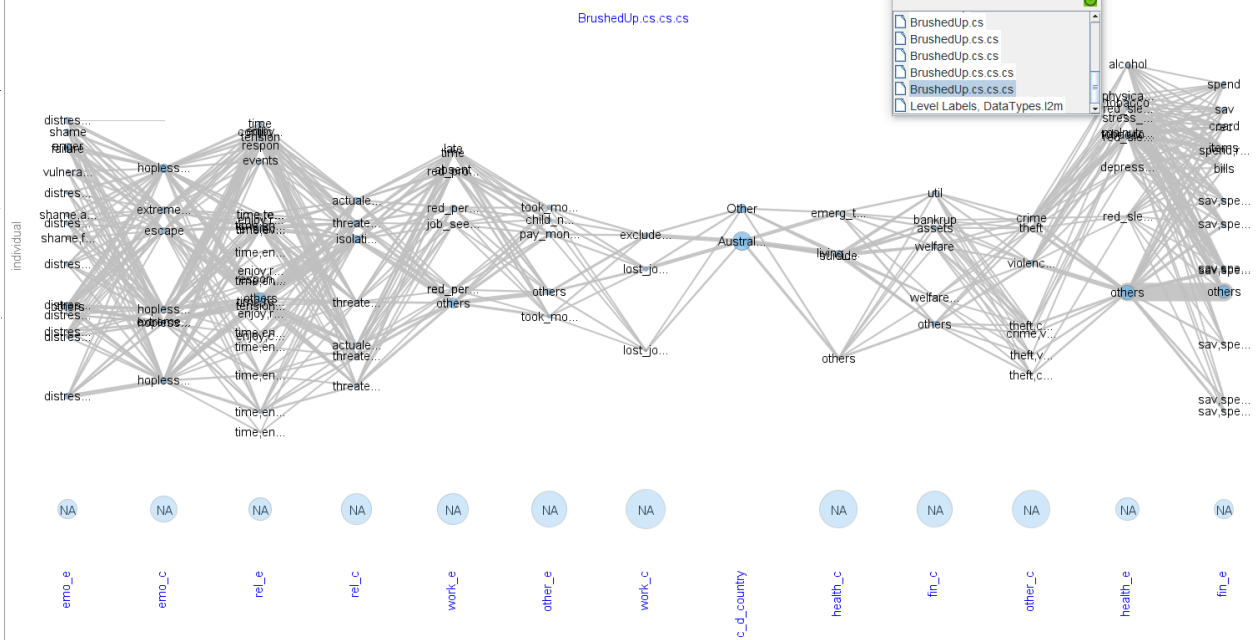
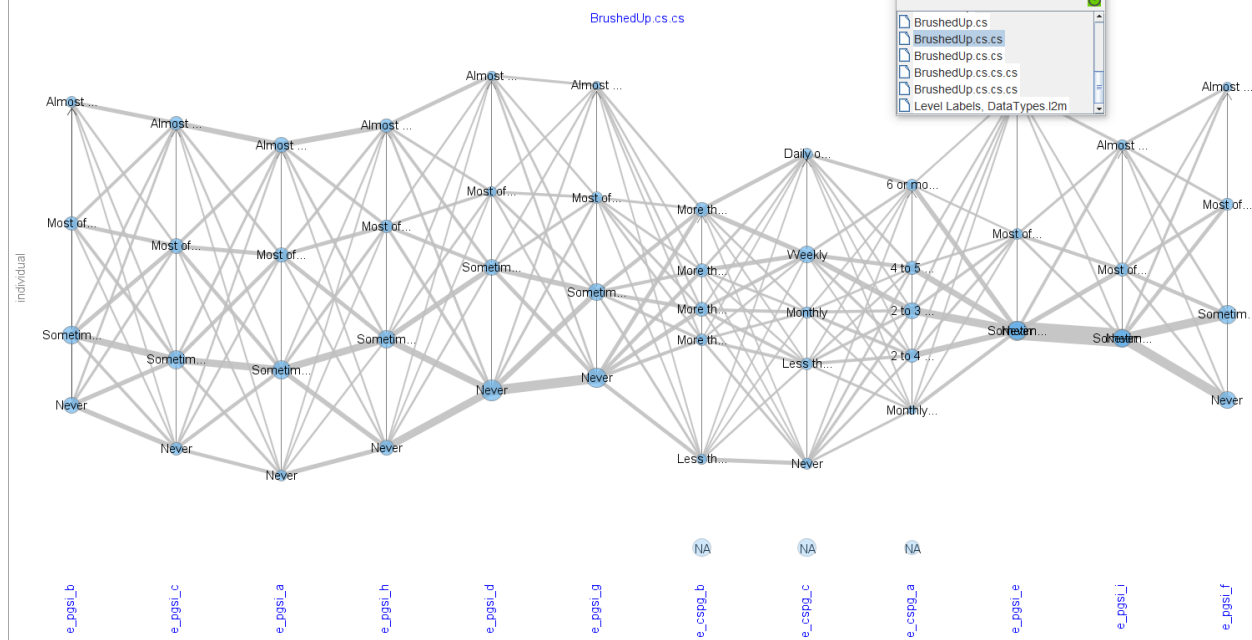
BrushedUp.cs.cs.cs.cs

Level Labels, DataTypes.l2m

Projections

- 変量(軸)をいくつか選んで、一つのデータテーブルからいくつかのサブデータテーブルを作り出す









高速化とメモリー使用量削減

- 高速化
 - 隠しプログラミングテクニック(Java)
 - 最終的にもっとも時間を費やすのは最後の描画
 - イメージとして保存して再利用
- メモリー使用量の削減
 - ガベージコレクションの効果的な適用
 - 負荷の少ない計算なら結果を保存せず必要になったときに求める(トレードオフ)
- 目安
 - 30列30万行 CSV ファイル(200M): 描画まで28秒, メモリー 1GB

PlotType の導入

- TextilePlot 
 - 各軸の位置尺度を水平性規準で定める
- Parallel Coordinate Plot
 - 各軸の位置尺度は任意
 - 画面の高さいっぱいまで伸ばす 
 - ママ (位置=0, 尺度=1) 
 - 位置・尺度を与える (例: 線形モデル当てはめ結果の表示) 

まとめ

- どのような機能が必要で, どんな形で実装すればよいのかわかっているわけではない
- 実験を重ね, 改良を重ねるしかない
- 芸術の側面, 美しく楽しい環境
- 様々なデータにする知見の蓄積
- 利用者からのフィードバック
- データの価値評価のための基本ソフトウェアになることを期待