

# TRADのいま

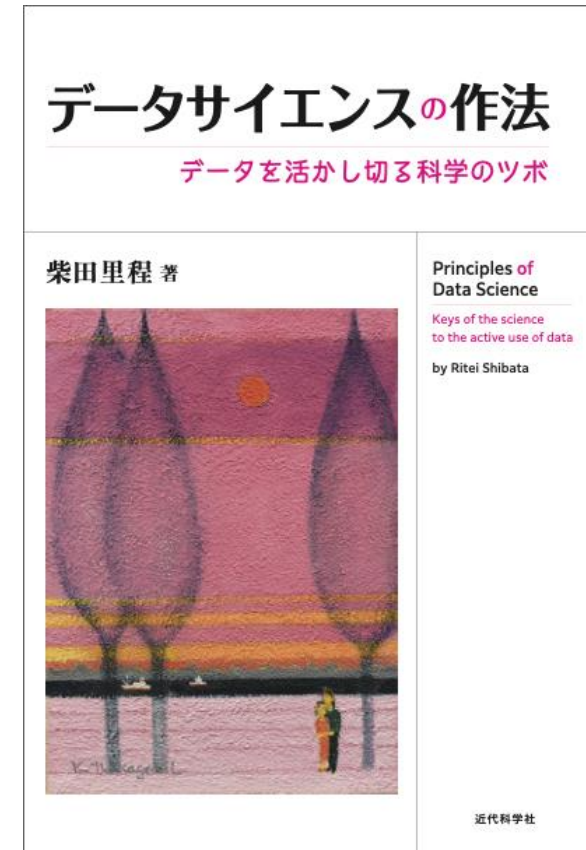
データサイエンスコンソーシアム, 慶應義塾大学 柴田里程

# データサイエンス研究

- 基本ソフトウェアの必要性
  - 基本的な研究環境として: 知見の蓄積
  - 研究スタイルの変化にあった環境
  - 研究者間のコミュニケーションの道具として
  - 研究者と実践者の間のコミュニケーションの道具として
- 公開ソフトウェア
  - R (CUI) <https://cran.r-project.org>
  - TRAD (TextilePlot, R and DandD) (GUI) <https://datascience.jp/TRAD.html>

# Data Portal : TRAD

- データの入口
  - データの概略を視覚的に把握
    - データと遊ぶ
    - 楽しむ, 夢想する
    - データの攻略策を練りあげる
- TRADはデータの入口でコンシェルジュの役割も果たす  
研究者の作りあげた本格的な公開ソフトウェア
  - R はCUIで, データ全体を眺めるのは苦手
  - TRADはデータ全体をそのまま眺められる
  - EVコード署名, Apple公証済みされた安全なソフトウェア
  - Easy to use, User Friendly
  - Java で書かれた Multi-platform software
  - RoboHelp による本格的なヘルプ
  - 多言語対応





TRADの設定

作業ディレクトリ

保存 表示言語

DandDエディタ 変容

フィルタ

経軸順序

表示



DandDライブラリ



なにも読み込まれていません



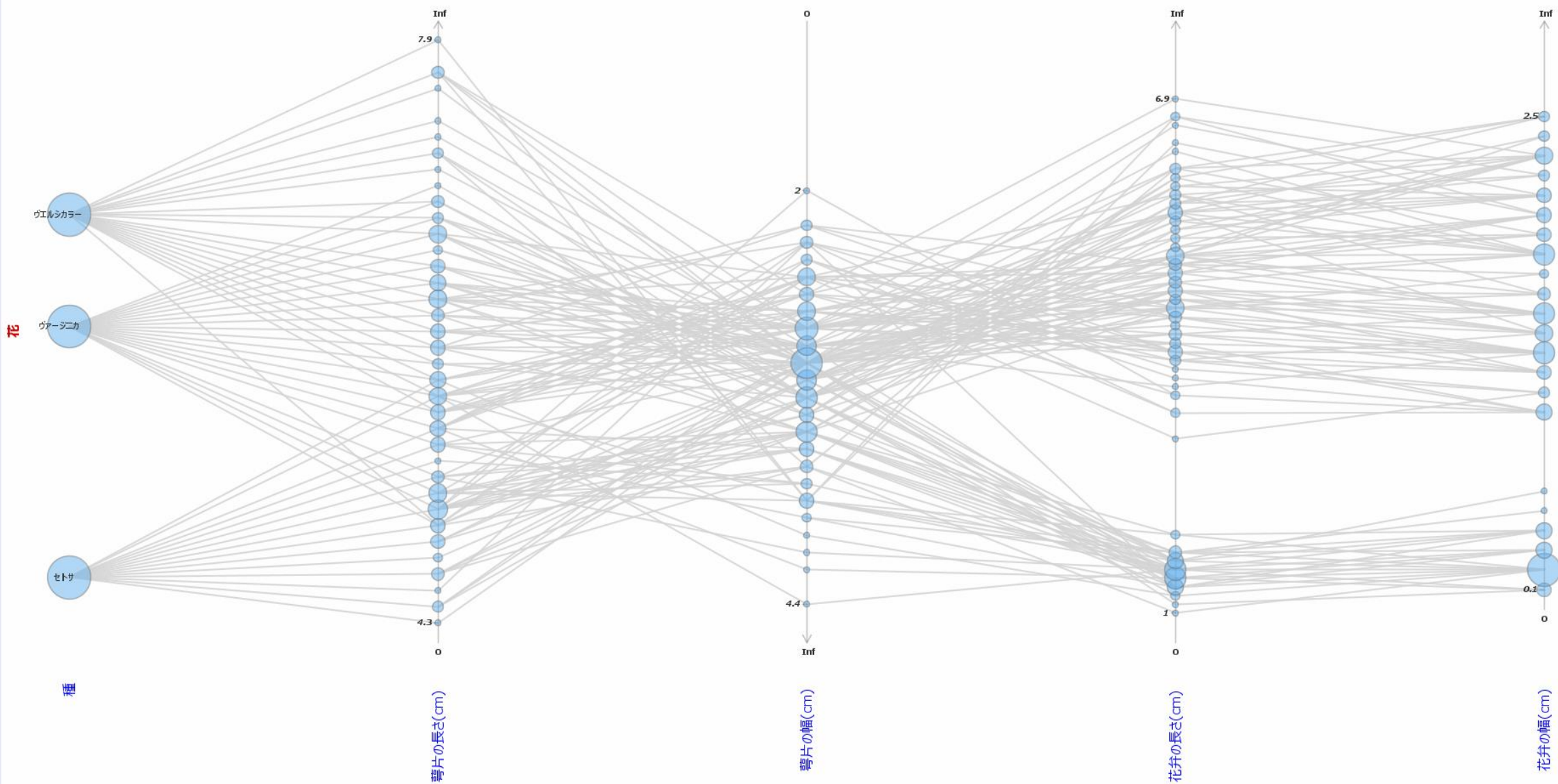


アイリスデータ



0.774

アイリスの花の特徴



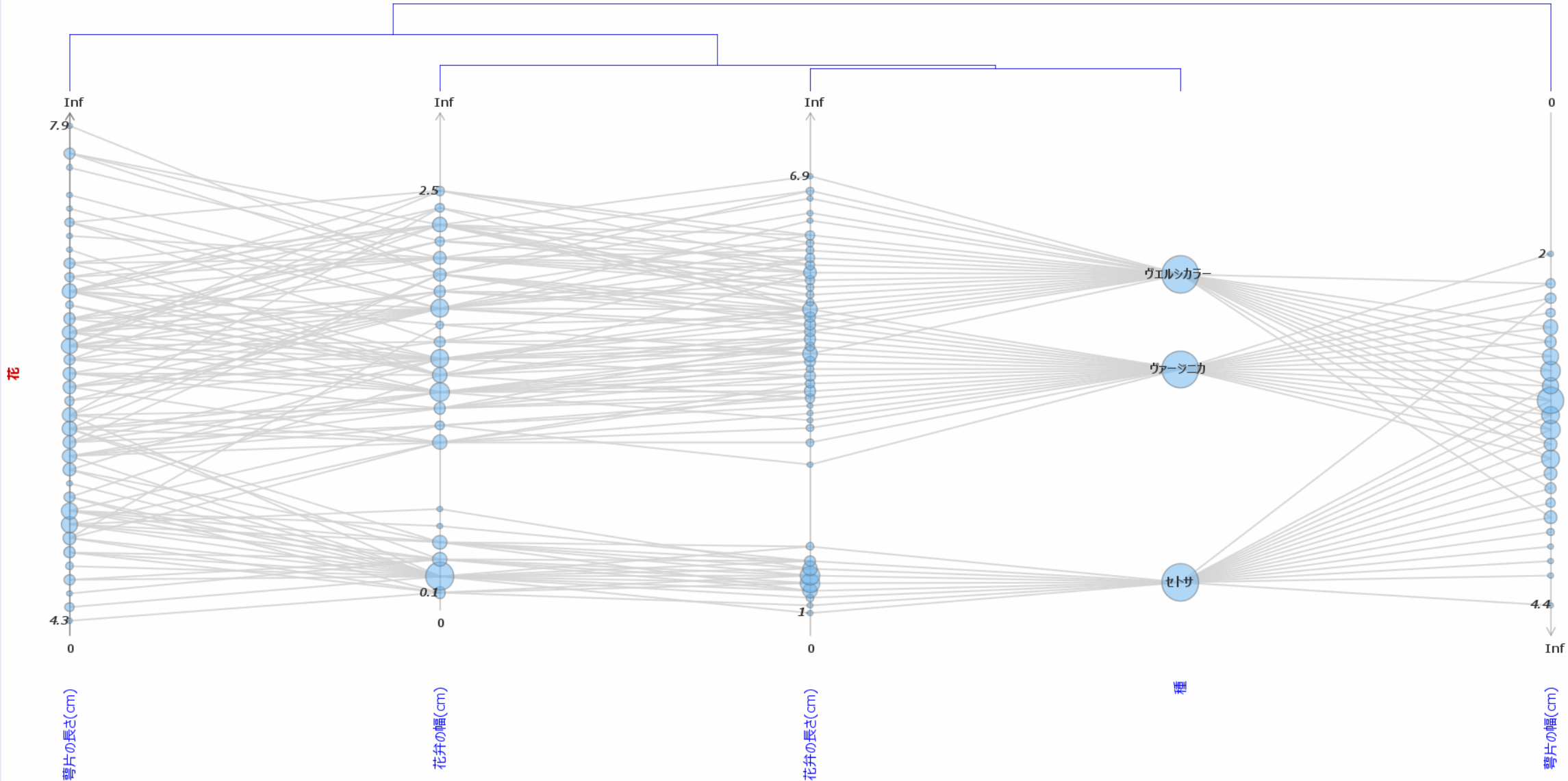


アイリスデータ



0.774

アイリスの花の特徴



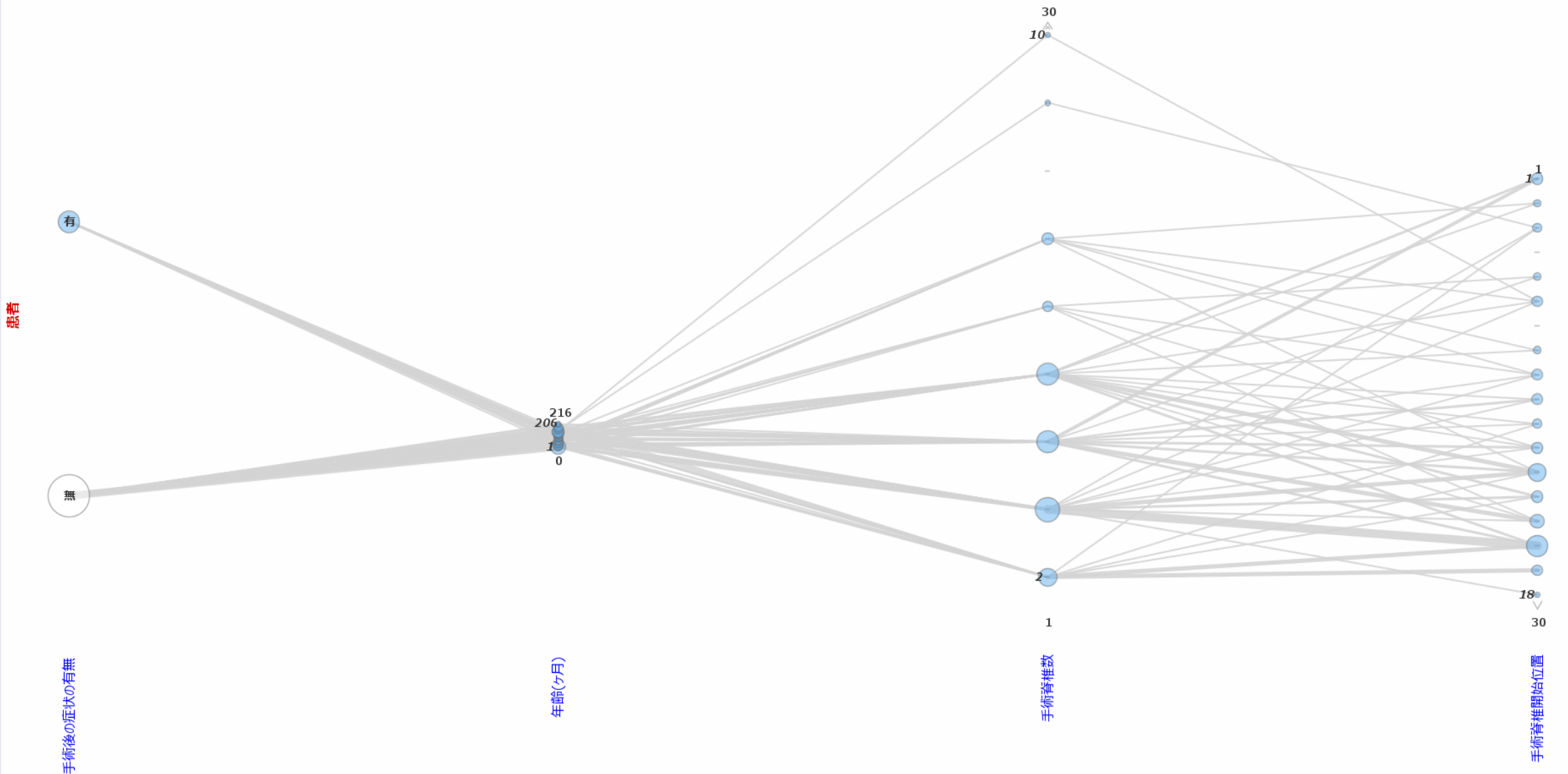


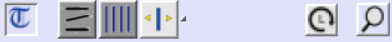
脊柱後弯症



0.461

脊柱後弯症手術の成否



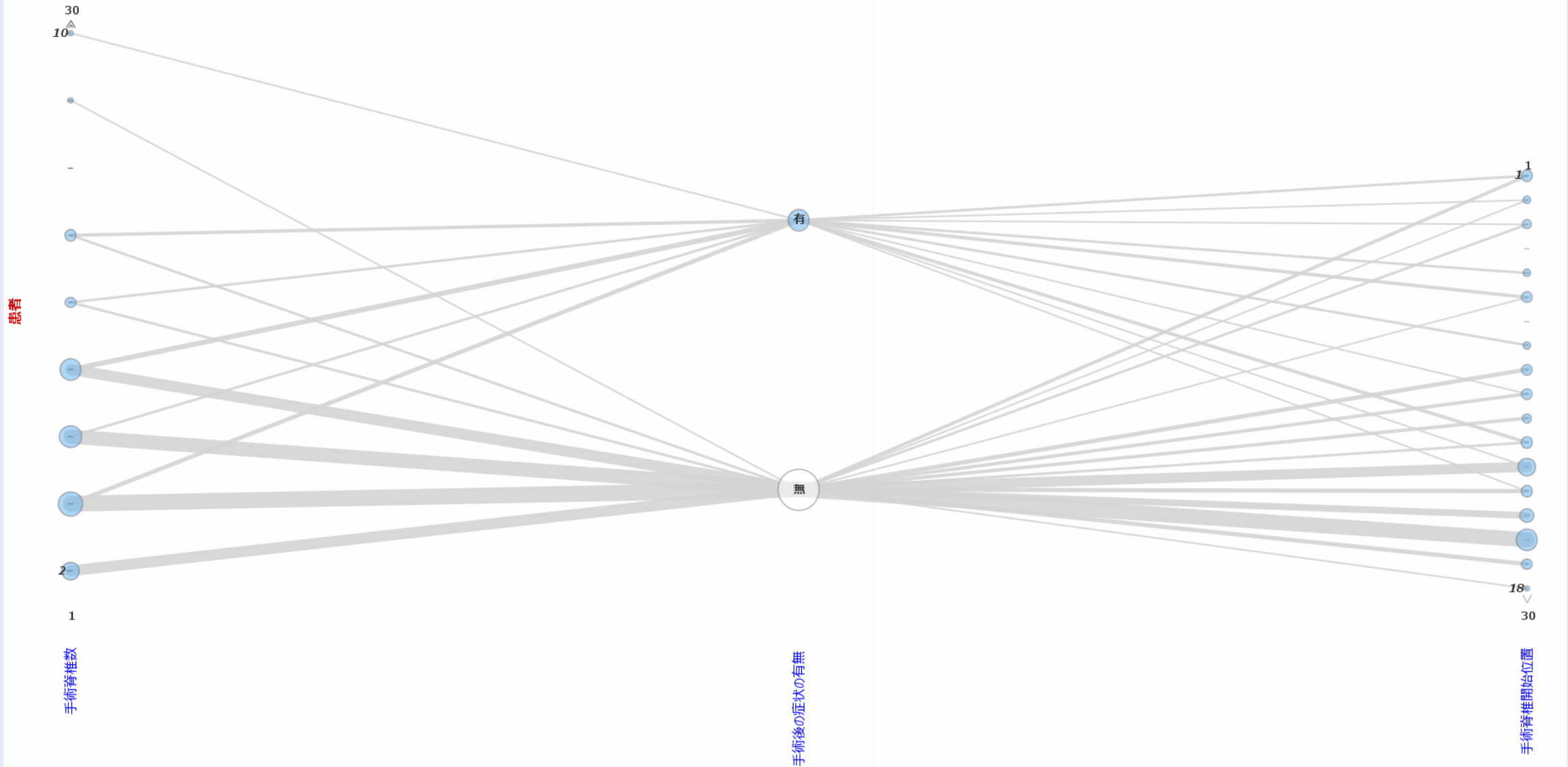


脊柱後弯症



0.614

### Kyphosis Data







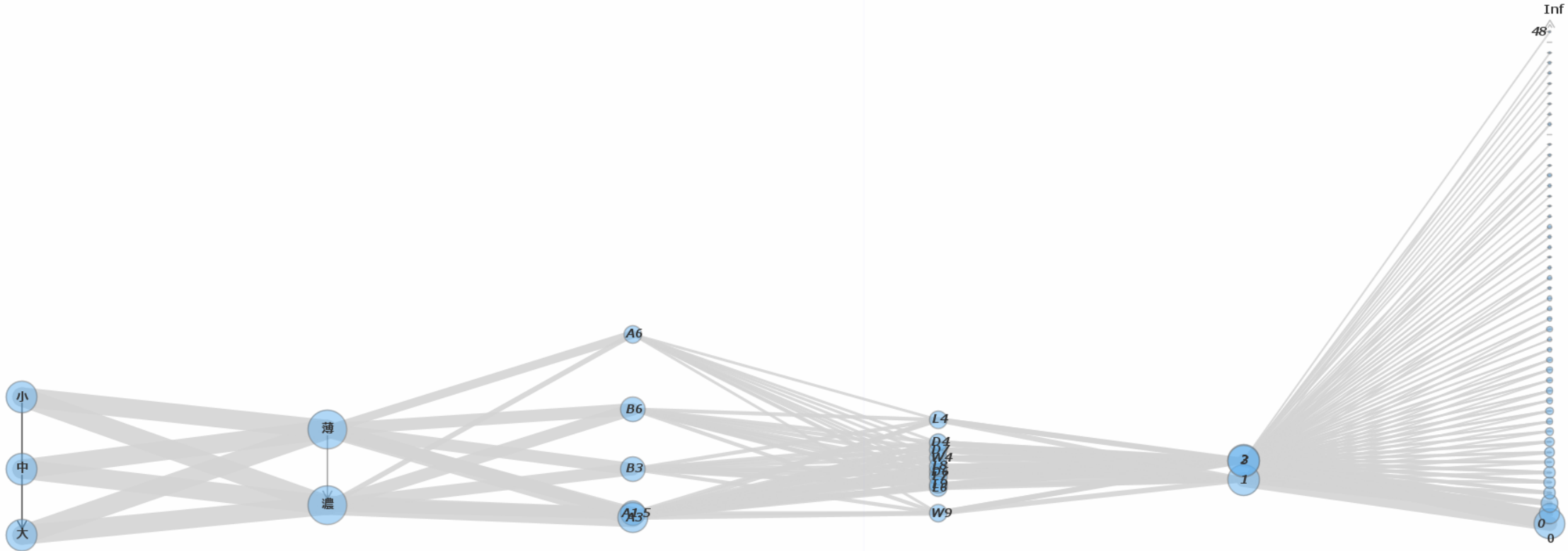
半田付け実験データ



0.312

不良要因

実験



マスクにあけた穴の大きさ

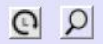
半田の量

マスクの種類と厚さ

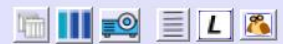
半田付けする部分の形状と大きさ

試験片の番号

不良箇所数

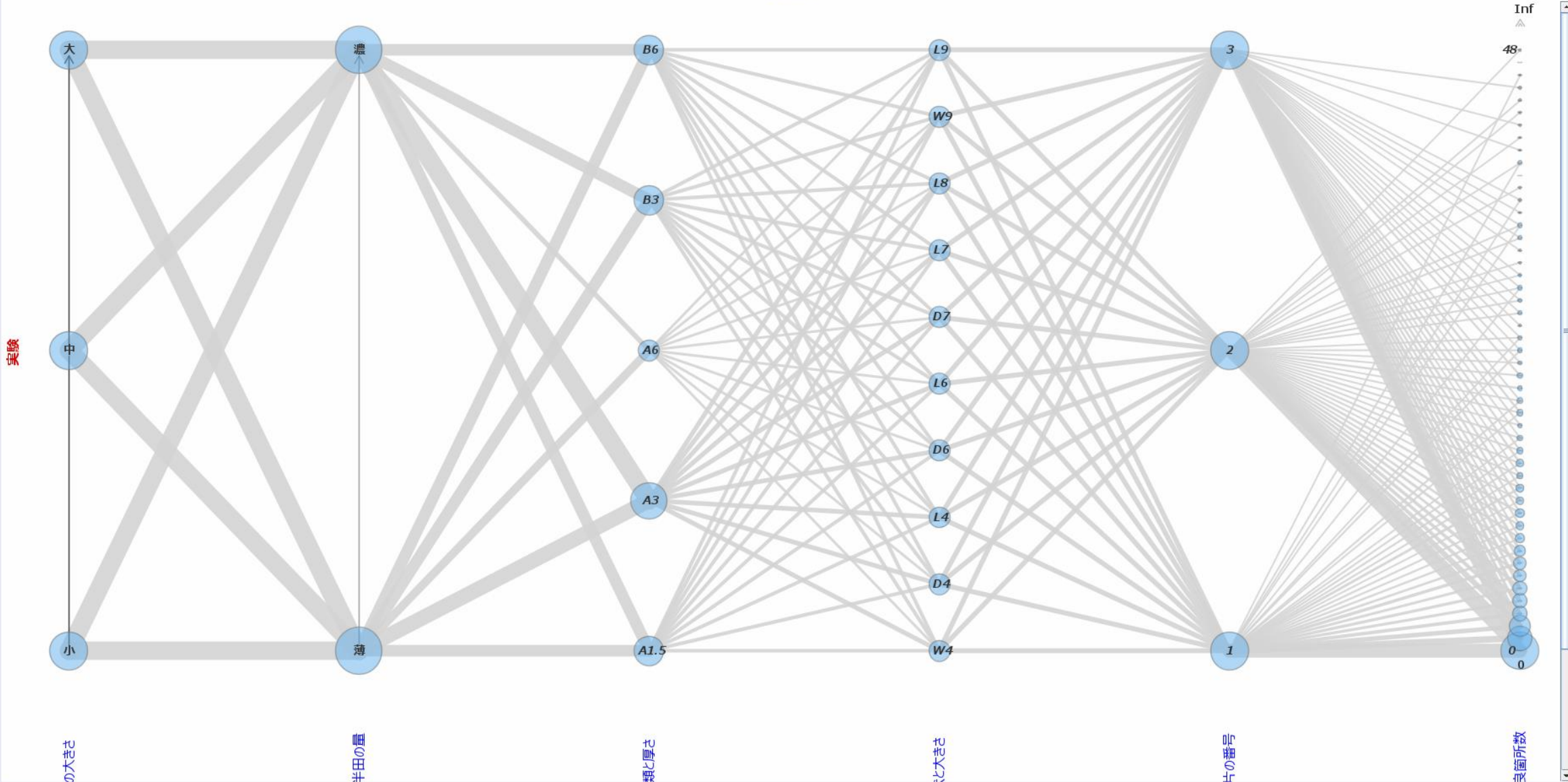


半田付け実験データ



和

不良要因



実験

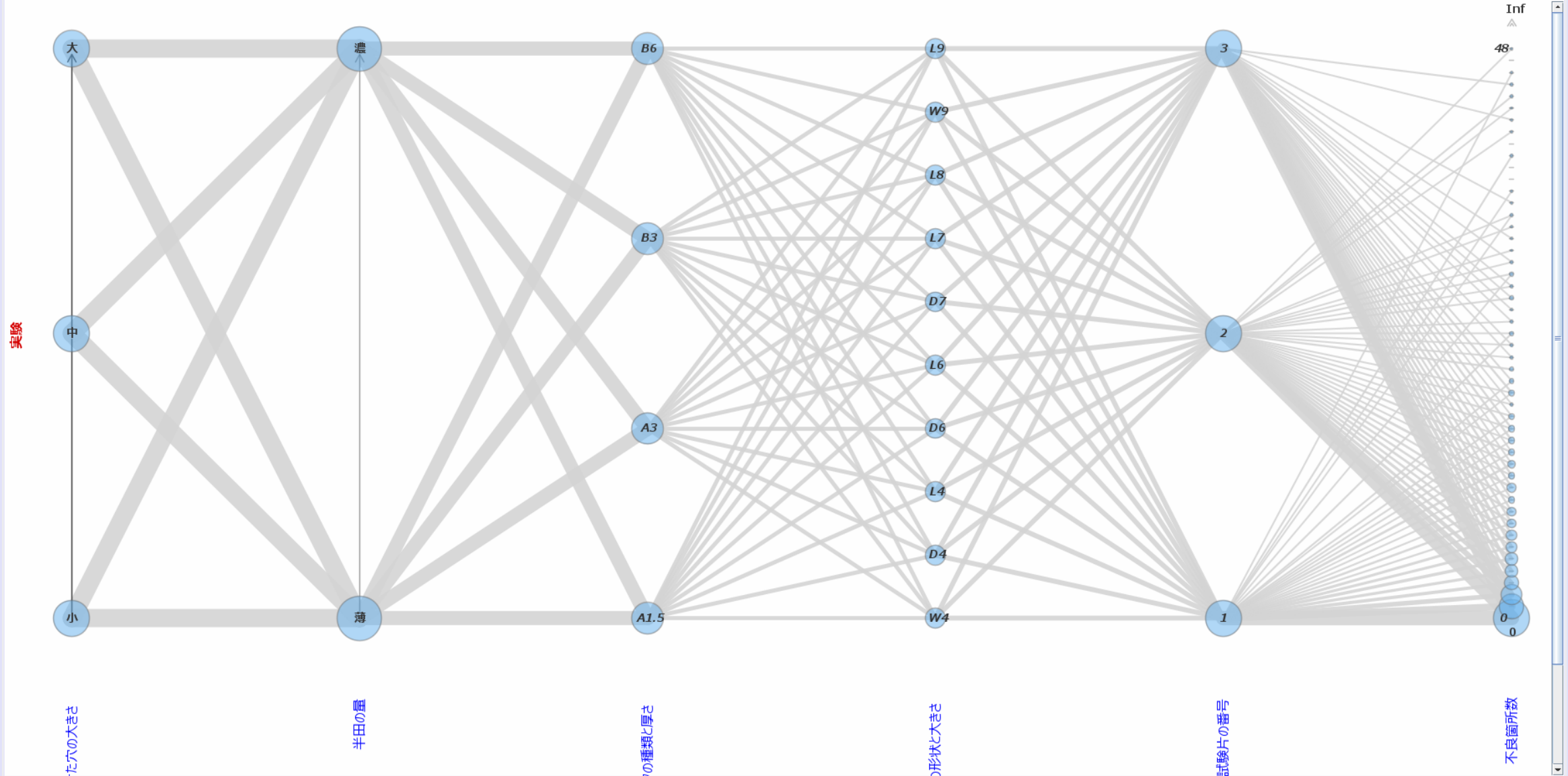
Inf

48

0

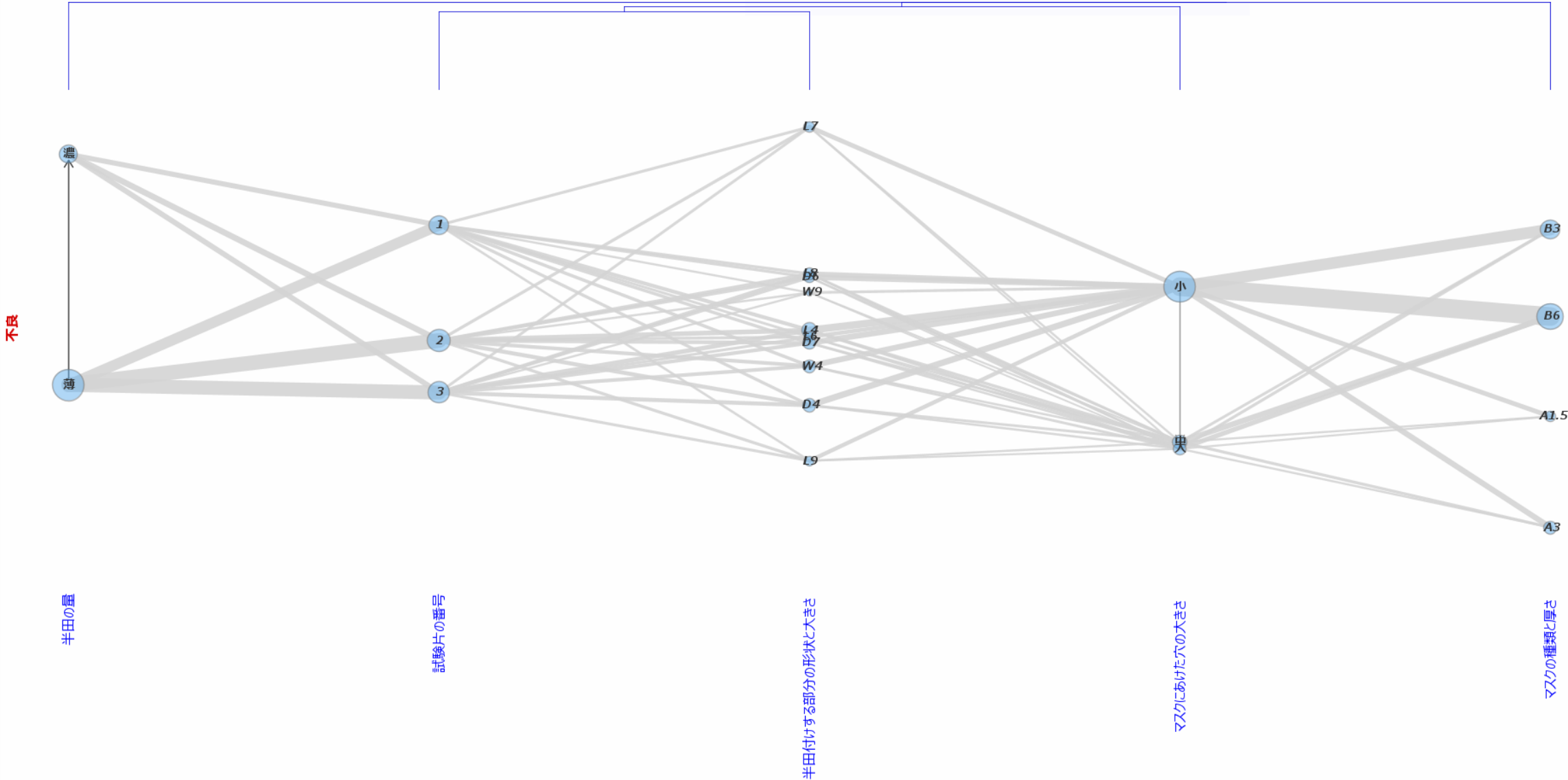
0

不良要因(Balanced)



0.242

型変更後



# TextilePlot表示

- 並行座標プロットの種類
- 座標を緯糸で結んで一記録を表している
- 各軸のローケーション, スケールを水平基準で緯糸がなるべく水平になるように定めているだけ.
- 軸(変量)のクラスタリングは, 軸を結ぶ緯糸の傾きの絶対和を距離とする最小距離法 (single linkage)
- 各変量の型を反映した表示

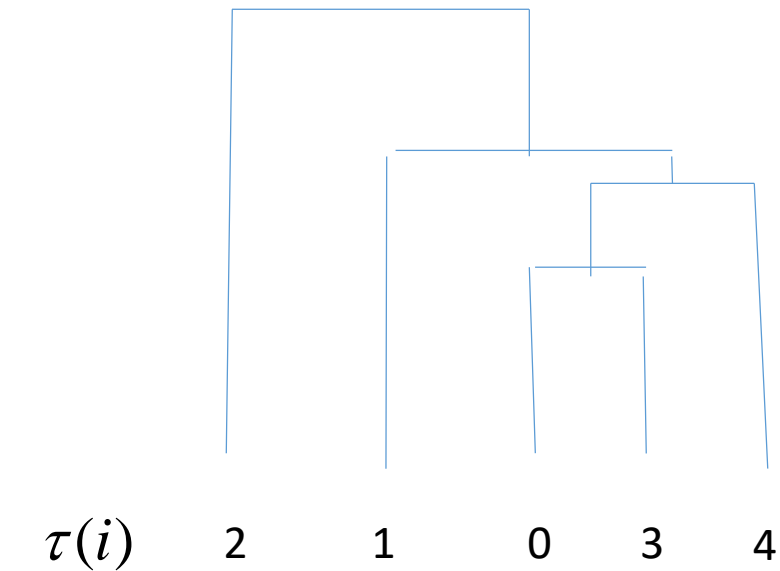
# 大規模データ

- 記録数大
  - Thinning, 縦スクロール
- 変量数大
  - 横スクロール
  - (変量の)クラスタリングアルゴリズム
    - 変量数を  $p$  としたとき
      - ナイーブなアルゴリズム: 計算時間  $O(p^3)$ , メモリ  $O(p)$
      - SLINK: 計算時間  $O(p^2)$ , メモリ  $O(p)$

# SLINK アルゴリズム (Sibson, 1973)

最小距離による階層的クラスタリング

例:  $p = 5, i = 0, 1, 2, 3, 4$

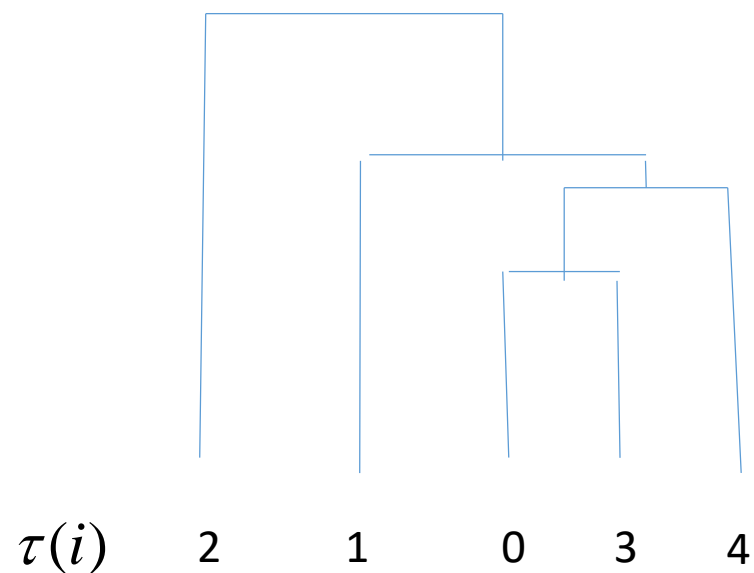
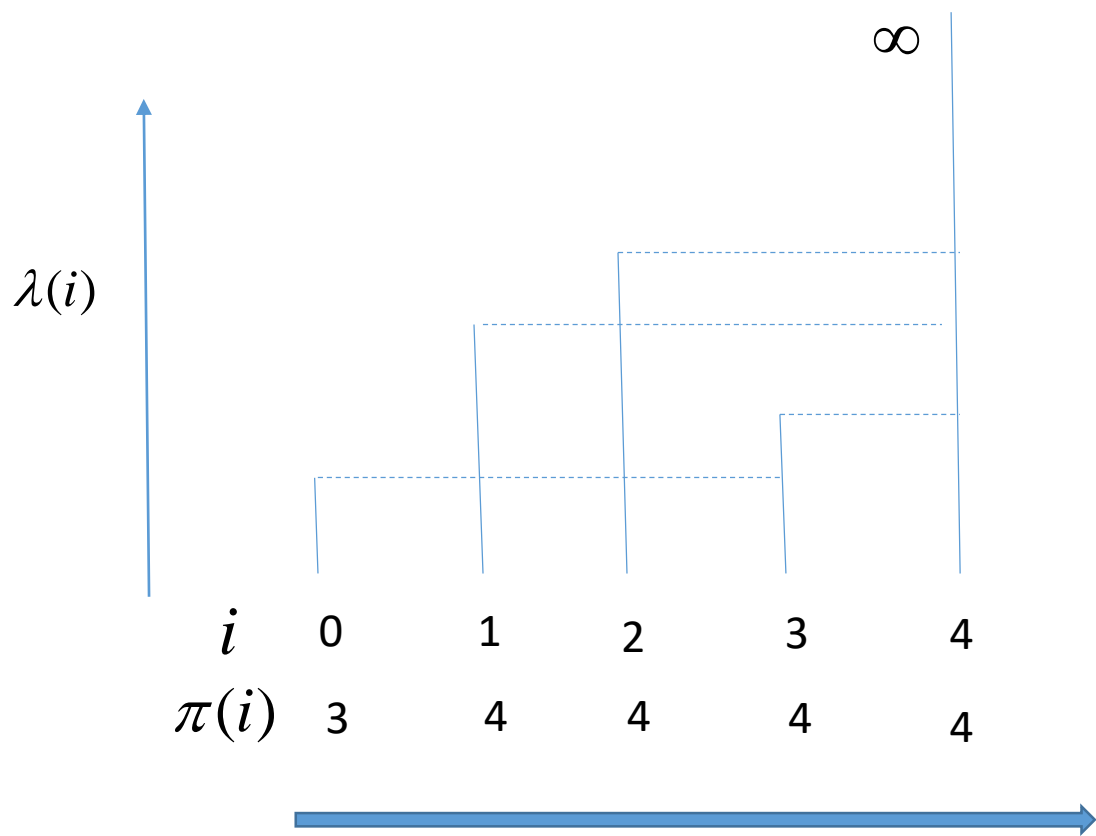


並べ替え



ナイーブなアルゴリズム

$\pi(i) = i$ があるクラスタに加わったときの  
 そのクラスタ内の最大要素番号(クラスタ番号)  
 $\lambda(i) =$ その時の距離





# 事後的に変量数が大きくなる例

$X_1, X_2, \dots, X_m$  : 説明変量 (条件),  $Z_1(t), Z_2(t), \dots, Z_p(t)$  : 時系列 (計測値)

このような実験が  $n$  回行われたとき, その関係を眺める一つの方法は  
各時系列をいくつかの時点  $t_1, t_2, \dots, t_k$  での値で代表させ,  $m + pk$  変量

$X_1, X_2, \dots, X_m, Z_1(t_1), Z_1(t_2), \dots, Z_1(t_k), Z_2(t_1), Z_2(t_2), \dots, Z_2(t_k), \dots, Z_p(t_1), Z_p(t_2), \dots, Z_p(t_k)$

の  $n$  記録として眺めること.

Dandドライブラリ

