# テキストマイニングツール TTM (TinyTextMiner) の理念と使い方



松村真宏・三浦麻子・金明哲

### テキストマイニングの登場

- テキストデータのような定性データは、大量のデータを分析することで安定した傾向が見いだせますが、人手で大量のテキストデータを分析することは現実的にはほとんど不可能でした
- テキストマイニングの登場によって、大量のデータを <u>統一的な視点・基準</u>から<u>少ない労力</u>で分析することが 可能になりました
- テキストマイニングは、世の中で流行っている話題 や、人々のニーズや不満を定量的に把握する手段とし て徐々に使われ始めています

### テキストマイニングの主な基盤技術

- 自然言語処理
  - 形態素解析,構文解析
- 統計解析
  - 多変量解析, 仮説検定
- データマイニング
  - 分類器, 予測器

### テキストマイニングの内側

- テキストマイニングの処理は、 「前処理」と「後処理」に大別できます。
- 前処理(テキストデータならではの処理)
  - 語の切り出しと集計
  - これが面倒…
- 後処理(多変量解析やデータマイニングと同じ処理)
  - 集計データの統計処理
  - 従来の手法が流用できます!

### テキストデータの特徴 (1/4)

- 語彙や表現の揺れ(漢字・仮名・カタカナ・大文字と 小文字・全角と半角・誤字・新語)が多い
  - 「内閣総理大臣」と「首相」
  - 「打ち合わせ」と「打合せ」
  - 「インタフェース」と「インタフェイス」
  - 「税金」と「血税」
  - 「スナナレ」「もしドラ」
  - 「ドコモ」と「DoCoMo」と「docomo」

### テキストデータの特徴 (2/4)

- 日本語は文法も曖昧
  - <u>クロール</u>で泳いでいる彼女を見た
  - 望遠鏡で泳いでいる彼女を見た
  - プールで泳いでいる彼女を見た
  - 先生とお酒を飲む
  - ビールとお酒を飲む

### テキストデータの特徴 (3/4)

- 語の境界に曖昧性がある
  - 「そこではきものをぬげ」
  - ▶ そこで/はきもの/を/ぬげ
  - ▶ そこでは/きもの/を/ぬげ
- うなぎ文
  - 「ぼくはウナギだ」
- こんにゃく文
  - 「こんにゃくは太らない」

### テキストデータの特徴 (4/4)

- 意味の文脈依存性
  - 「学校で遊ぶ」
  - ▶ このときの「学校」は場所としての学校
  - 「学校が談話を発表した」
    - ▶ このときの「学校」は法人的側面を表す

### 形態素解析

- 形態素解析は、自然言語で書かれた文章を語(形態素)に分割する処理のことです
- Chasen (奈良先端大), MeCab (工藤拓氏), JUMAN (京都大学) が公開しているオープンソースのソフトウェアが有名です
- 新聞記事を対象とした場合の<u>精度は99%以上</u>ですが、 話し言葉を対象とすると精度は下がります(それでも 十分実用的な精度です)
- 常に新しい言葉が生まれているので、<u>未知語</u>(辞書に 載っていない語)問題はなかなかやっかいです

### 形態素解析の実行例

• 「親譲りの無鉄砲で子供の時から損ばかりしている。」を MeCab にかけた結果です

```
親譲り
      名詞, 一般, *, *, *, *, 親譲り, オヤユズリ, オヤユズリ
      助詞,連体化,*,*,*,の,ノ,ノ
      名詞, 一般, *, *, *, *, 無鉄砲, ムテッポウ, ムテッポー助詞, 格助詞, 一般, *, *, *, で, デ, デ
      名詞, 一般, *, *, *, *, 子供, コドモ, コドモ
      助詞,連体化,*,*,*,の,ノ,ノ
      名詞, 非自立, 副詞可能, *, *, *, 時, トキ, トキ
      助詞, 格助詞, 一般, *, *, *, から, カラ, カラ
から
損
      名詞, 一般, *, *, *, *, 損, ソン, ソン
      助詞, 副助詞, *, *, *, ばかり, バカリ, バカリ
ばかり
      動詞, 自立, *, *, サ変・スル, 連用形, する, シ, シ
      助詞,接続助詞,*,*,*,*,て,テ,テ
      動詞, 非自立, *, *, 一段, 基本形, いる, イル, イル
いる
      記号,句点,*,*,*,*,。,。,。
```

### 機能語と内容語

- 語は、助詞や助動詞といった「<u>機能語</u>」と、名詞、形容詞、動詞、副詞といった「<u>内容語</u>」に大別できます
- 機能語は、それ単体では意味を持たない語なので、文章の内容を理解する際の助けにはなりません
- 内容語は、名称、性質、動作、状況など、文章の内容の一部を表しているので、内容を理解する際の助けになります。しかし、名詞と結びつかないと意味が特定できない場合が多いです
- したがって、名詞は必須で、分析の目的に応じて形容 詞、副詞、動詞を用いることが多いです

### 未知語について

- 形態素解析器の辞書に登録されていない語は「<u>未知</u> 語」として出力されます
- 未知語の品詞推定は研究レベルでは実装されていますが、まだ実用レベルには達していません
- 未知語は、単なるゴミであることも多いのですが、世の中の流行を反映した「新しい語」(例えば「婚活」や「H1N1」など)が含まれていることもあるので油断なりません
- なので、取り敢えず未知語は分析対象に加えて、不便 があれば臨機応変に対応することが多いです

### 構文解析

- 構文解析は、語と語の係り受け関係を分析する処理の ことです
- <u>CaboCha</u> (工藤拓氏), <u>KNP</u> (京都大学) が公開しているオープンソースのソフトウェアが有名です
- 新聞記事を対象とした場合でも<u>精度は80%くらい</u>ですが、確からしい結果だけを利用すれば十分使えます
- 特定の語と関係する語(例えば,「美味しい」の係り 受け先など)を見たいときなど,用途を限定した場合 にも十分使えます

### 構文解析の実行結果

• 「親譲りの無鉄砲で子供の時から損ばかりしている。」を CaboCha にかけた結果です

```
親譲りの-D
無鉄砲で-----D
子供の-D
時から----D
損ばかり-D
している。
```



- 「形態素解析は分かった、構文解析も分かった、それで、どうすればいいの?」という皆さんの心の声に答えるために を作りました
- はテキストデータを形態素解析器、構文解析器にかけて、その分析結果を読み込んで集計し、CSVファイルを出力するフリーウェアです
- はテキストマイニングの<u>前処理に特化</u>していますので、ここまでしか行いません、後処理は、みなさんの使い慣れたソフトウェアに読み込ませて、好きなように分析してもらいたいと思っています



### のスクリーンショット (1/2)



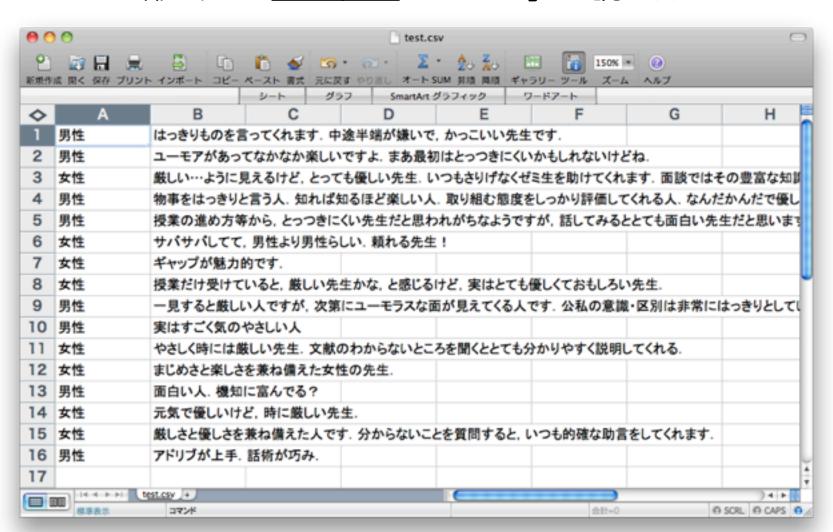


### のスクリーンショット (2/2)

TTM: TinyTextMiner	
解析 詳細設定 TTMについて	
キーワードファイル(任意)	
C:/Documents and Settings/mlab/My Documents/Downloads/keyword.txt	選択
同義語ファイル(任意)	
C:/Documents and Settings/mlab/My Documents/Downloads/synonym.txt	選択
不要語ファイル(任意)	
C:/Documents and Settings/mlab/My Documents/Downloads/noise.txt	選択
cabocha.exe(係り受け解析時は必須)	
C:/Program Files/CaboCha/bin/cabocha.exe	選択
品詞の選択(日本語のみ)	
▼ 名詞 ▼ 形容詞 ▼ キーワード ▼ 同義語 □ 動詞 □ 副詞	
その他	
English text	
□ 係り受け解析を行う(日本語のみ)	
0 語の出現頻度/出現件数の最小値	

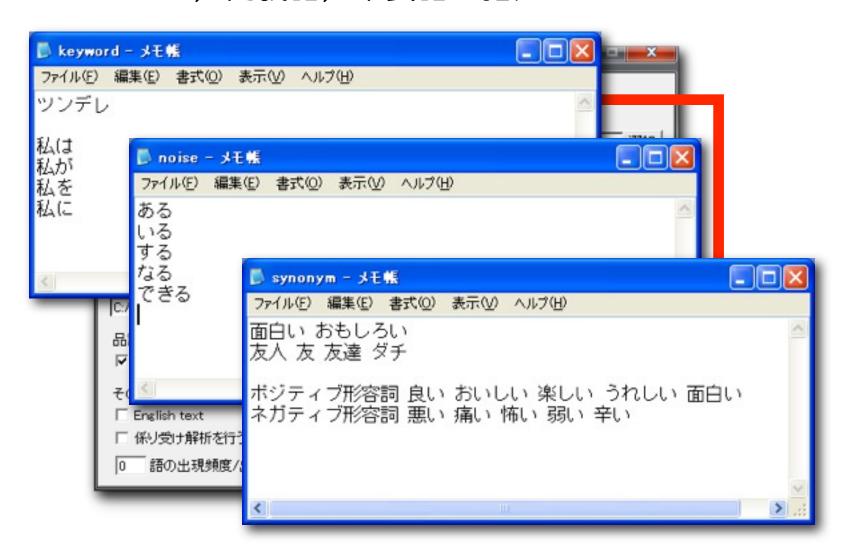
# の特徴 (1/4)

• CSV 形式の「<u>タグ付き</u>テキスト」を読み込みます



# の特徴 (2/4)

• キーワード、同義語、不要語を指定できます



# の特徴 (3/4)

• 品詞,係り受け解析,閾値,英文の設定もできます

TTM: TinyTextMiner	X
解析 詳細設定 TTMについて	
キーワードファイル(任意)	
C:/Users/matumura/Desktop/keyword.txt	選択
同義語ファイル(任意)	
C:/Users/matumura/Desktop/synonym.b/t	選択
不要語ファイル(任意)	
C:/Users/matumura/Desktop/noise.txt	選択
cabocha.exe(係り受け解析時は必須)	
C:/Program Files/CaboCha/bin/cabocha.exe	選択
品詞の選択(日本語のみ) マ 名詞 マ 形容詞 マ キーワード マ 同義語 厂 動詞 厂 副詞	
その他	
☐ English text	
□ 係り受け解析を行う(日本語のみ)	
0 語の出現頻度/出現件数の最小値	

# の特徴 (4/4)

• 6種類の出力ファイルを提供します

TTM: TinyTextMiner	( <b>-</b>   -   <b>×</b>
解析 詳細設定 TTMについて	
入力ファイル(必須)	
C:/Users/matumura/Desktop/test.csv	選択
出力フォルダ(必須)	
C:/Users/matumura/Desktop	選択
mecab.exe(必須)	
C:/Program Files/MeCab/bin/mecab.exe	選択
出力フォーマット	
▼ ttm1:語のタグ別集計(出現頻度)	
▼ ttm2:語のタグ別集計(出現件数)	
▼ ttm3:語×タグのクロス集計(出現頻度)	
▼ ttm4: 語×タグのクロス集計(出現件数)	
▼ ttm5:語×語のクロス集計(出現件数)	
▼ ttm6: テキスト×語のクロス集計(出現頻度)	
解析	中止  終了

# **かインストール**

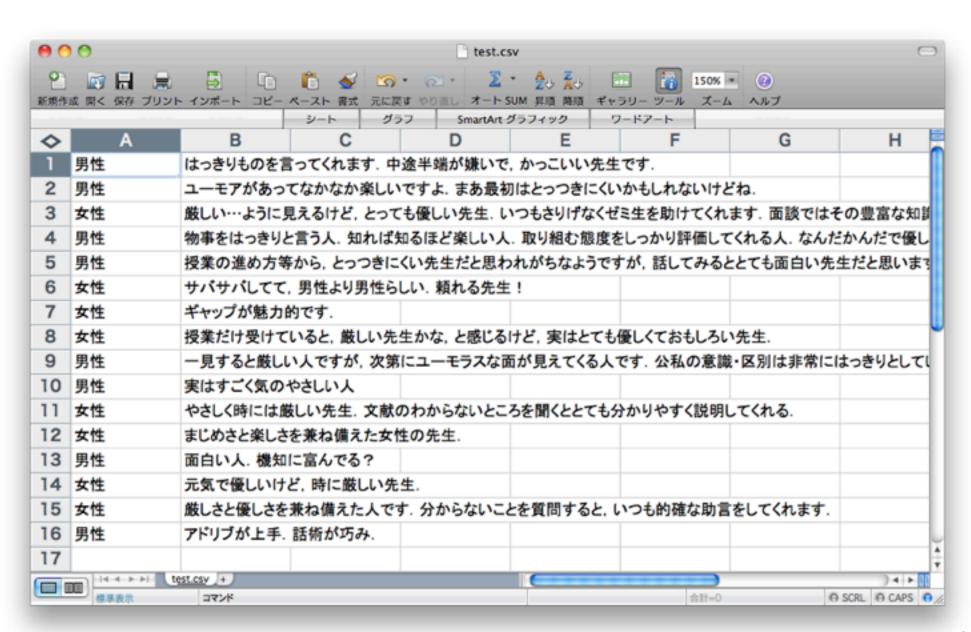
- Windows版とMac OSX版(10.5以降)があります
- http://mtmr.jp/ttm/ から ttm.exe をダウンロードするだけで 本体のインストールは終わりです
- 形態素解析を実行するためのソフトウェア MeCabを 別途インストールする必要があります(構文解析を行 うときは CaboCha もインストールします). 詳細は http://mtmr.jp/ttm/ をご覧ください
- Mac OSX版はOS内蔵のMeCabを使いますので
   MeCabを別途インストールする必要はありません

### サンプルデータ test.csv

• 三浦麻子先生のゼミに所属する16名(男女8名ずつ) の大学生が「三浦先生ってどんな人?」という質問に 対して自由に記述した文章

(http://mtmr.jp/ttm/test.csv からDLできます)

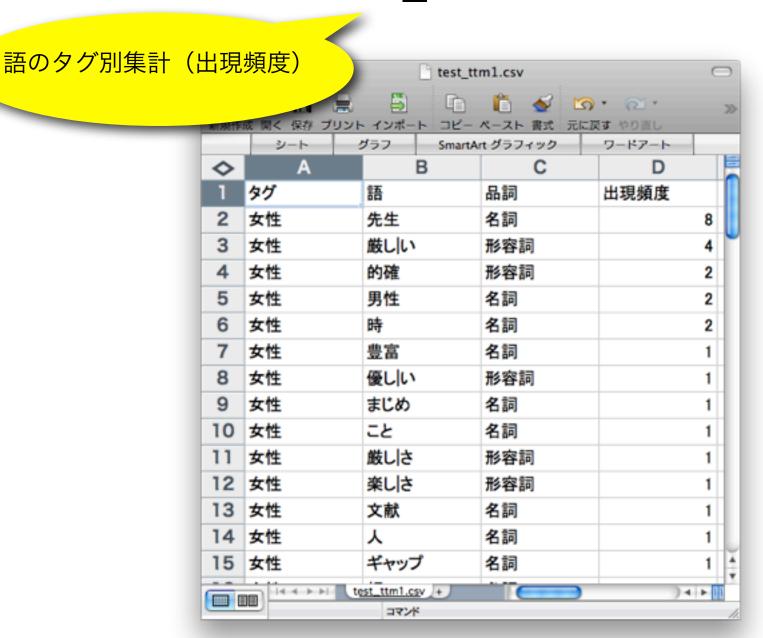
#### test.csv の内容



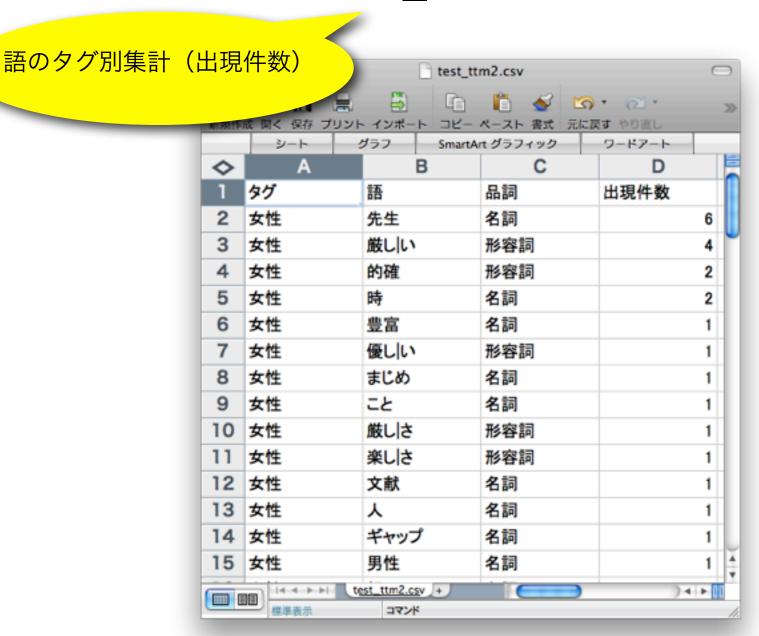
### 分析条件

- 次の条件で M で test.csv を分析してみましょう
  - 出力する品詞は「名詞,形容詞」
  - オプションファイルは「設定せず」
  - 出現頻度/出現件数の最小値は「O」

#### test\_ttm1.csv



### test\_ttm2.csv



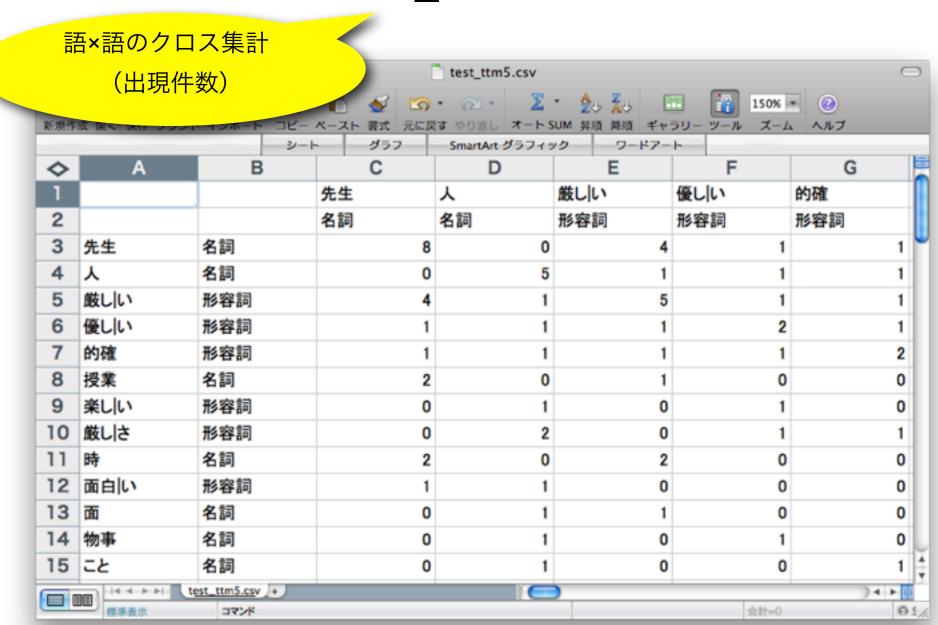
### test\_ttm3.csv



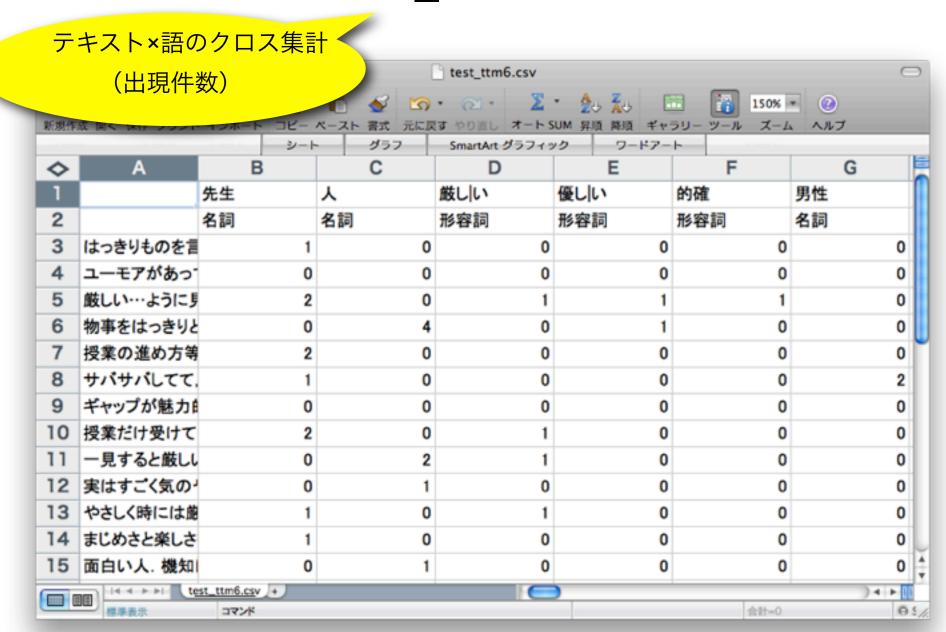
### test\_ttm4.csv



### test\_ttm5.csv



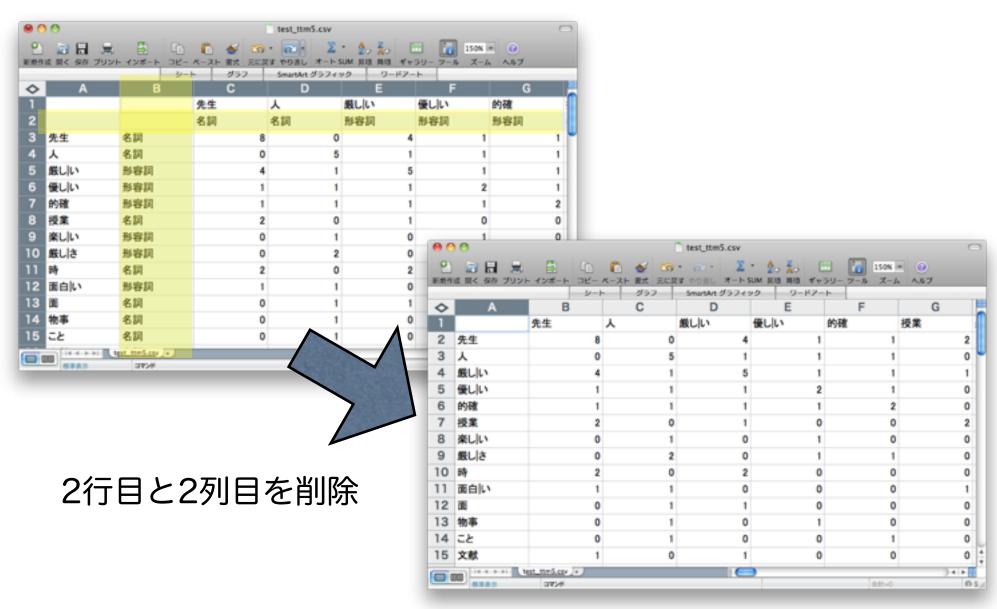
### test\_ttm6.csv



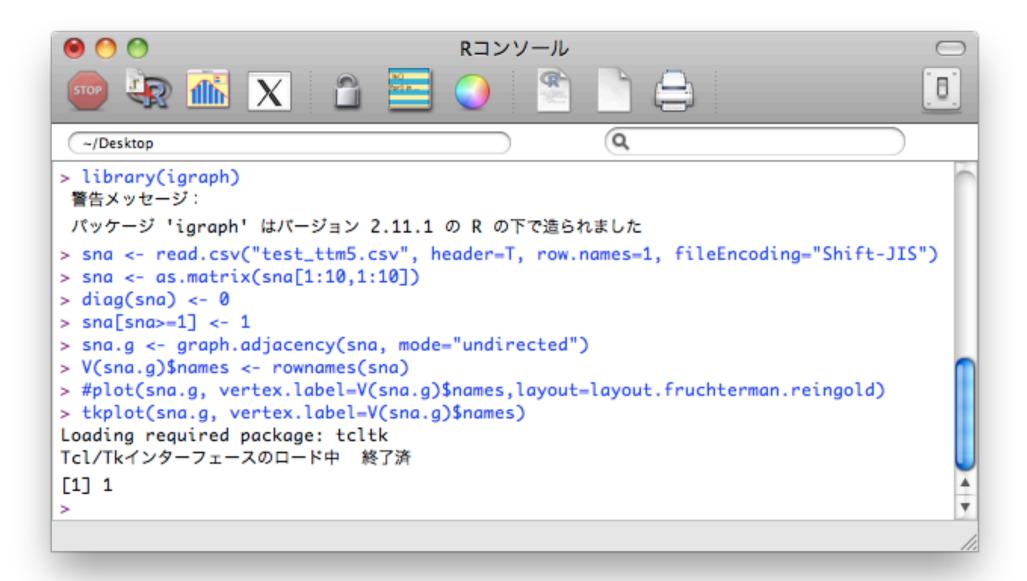


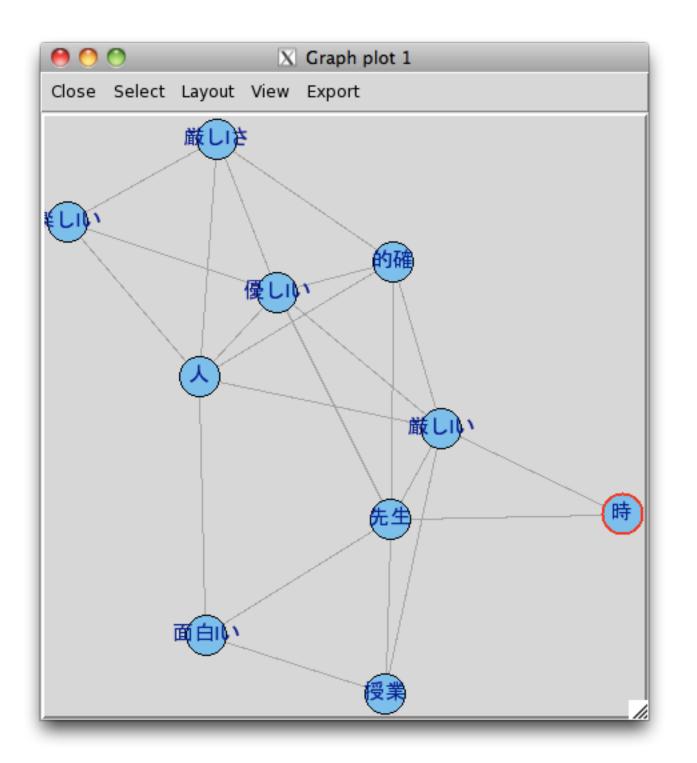
- 統計解析のフリーソフトウェアです
- いろんな人が便利な関数やパッケージを公開している (2009年6月17日現在, CRANには1849ものパッケージが登録されています) ので, コレーつで大抵のことはできます
- 形態素解析や構文解析を行うパッケージもあります
- データマイニングのパッケージもあります
- たくさんの書籍が出版されているので、マニュアルも 充実しています

### test\_ttm5.csvを編集

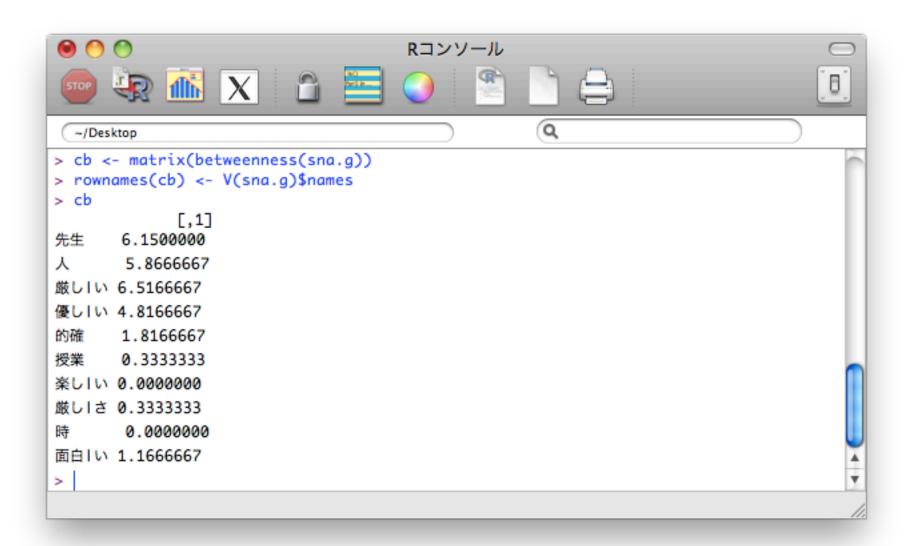


### 共起グラフの描画





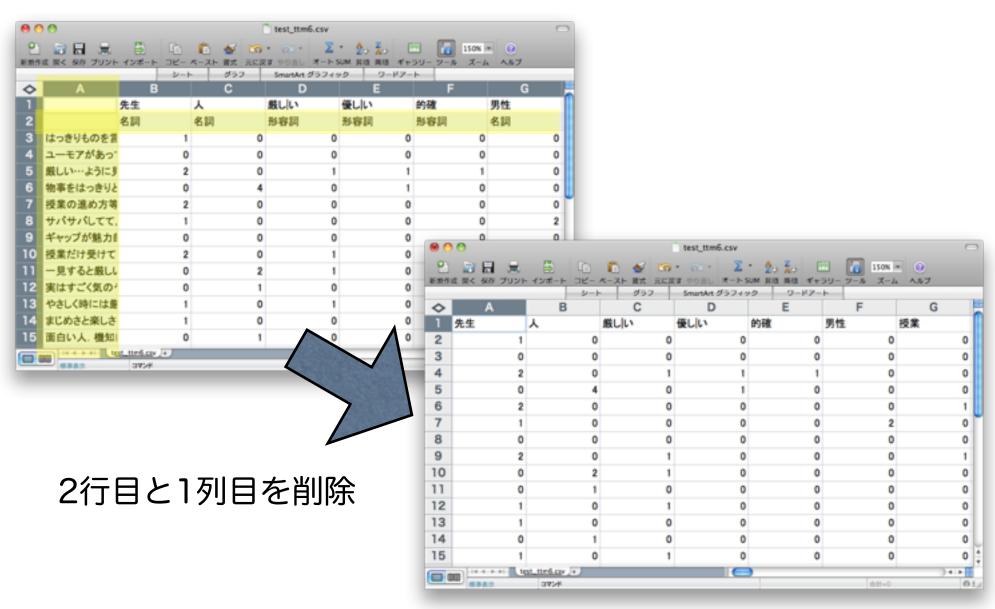
#### キーワード

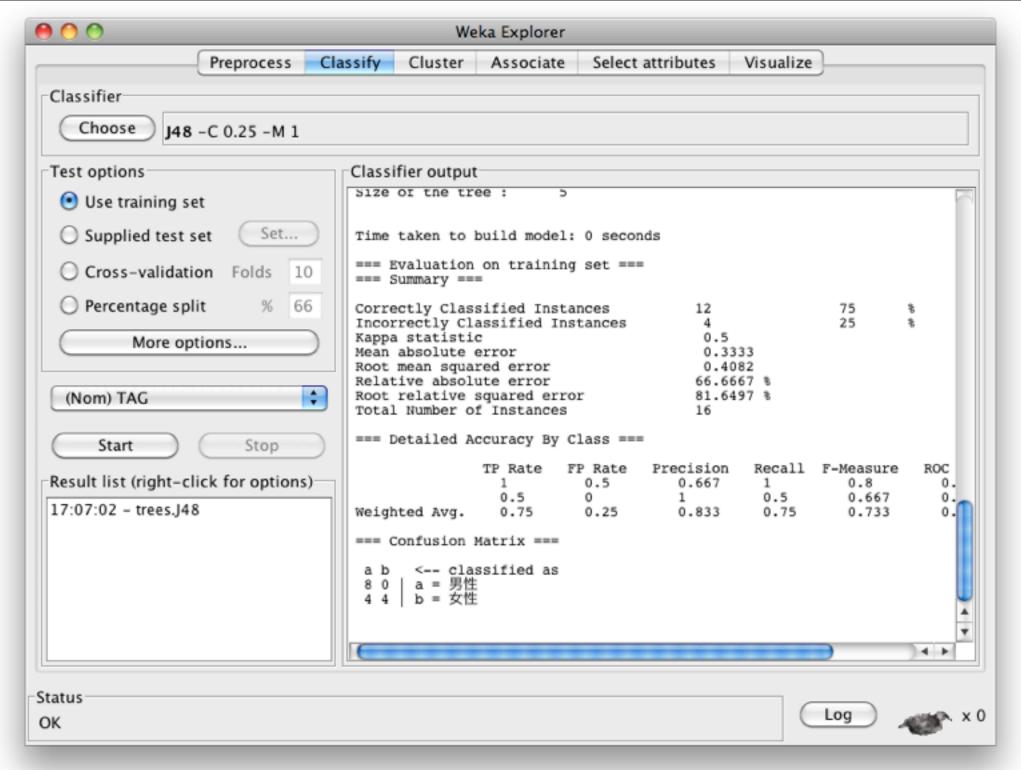


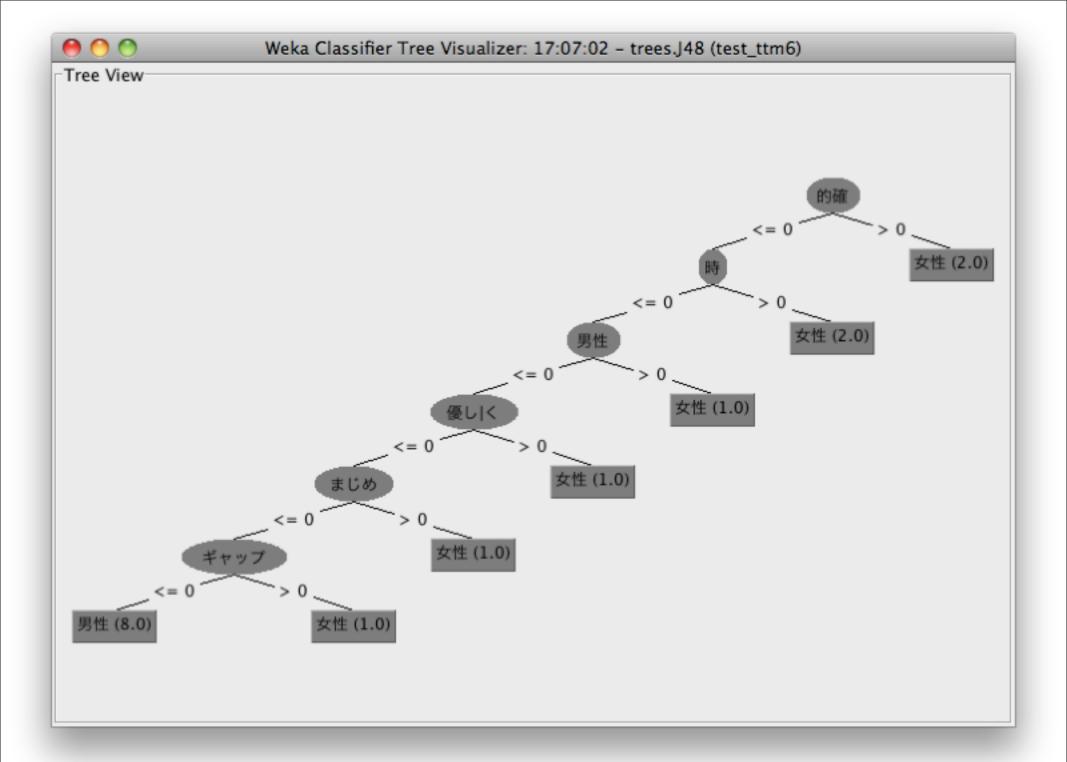


- データマイニングのフリーソフトウェアです
- 代表的なデータマイニングのアルゴリズムはほとんど 使えます
- GUIなので操作も簡単
- 最近はKNIMEやRapidMinerといったデータマイニングのフリーソフトウェアもあります

### test\_ttm6.csvを編集







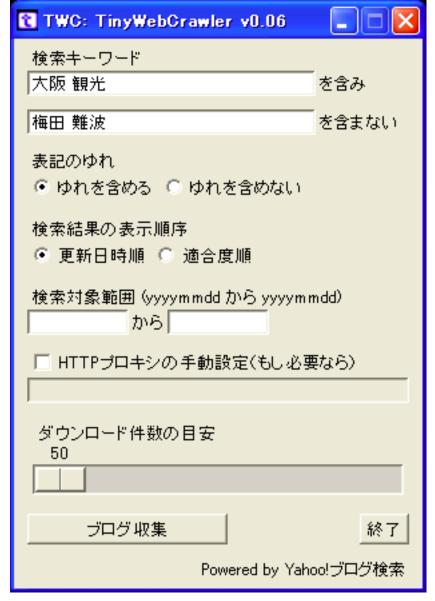
### おまけ:テキストデータの収集ソフト



TTC (TinyTweetCrawler)

http://mtmr.jp/ttc/

TWC (TinyWebCrawler) <a href="http://mtmr.jp/twc/">http://mtmr.jp/twc/</a>



### まとめ

- テキストマイニング
  - =〈前処理〉+〈多変量解析 or データマイニング〉
- テキストマイニングは語や表現の揺れが大きいので、 それを如何に吸収するかが重要になります
- 前処理さえ済めば、あとはRやWekaといった各自の 得意な土俵に持ち込んで勝負すればいいのです
- テキストマイニング恐るるに足らず!

### 宣伝

「人文・社会科学のためのテキストマイニング」 松村真宏・三浦麻子著,誠信書房 (2009) 2,520円



目次

第1章 序

第2章 TTMと関連ソフトウェアのインストール

第3章 TTMによるテキストデータの分析

第4章 Rを併用したテキストデータの統計解析

第5章 Wekaを併用したテキストデータのデータマイニング

第6章 テキストマイニングの応用事例

第7章 テキストマイニングの基盤技術