

## Weka 概要の補足

櫻井彰人  
慶應義塾大学工学部

## Weka



- 今回使用するソフトウェア
- ニュージーランドのワイカト大学が開発 (University of Waikato, New Zealand)
- Waikato Environment of Knowledge Analysis の略
- Weka: 探求心旺盛な飛べない鳥

## Weka の特徴

- Java言語で記述(使う人にとっては関係ないことですが)
  - しかし、そうはいつでも、すぐどこでも動くかつ安全なことは安心材料
- フリーソフト
  - 営利目的以外には自由に使用可能。改変可
- 機能の追加が可能

## Wekaの特徴(2)

- 日本語化が比較的容易 (Javaがそうだから)
- 欠点: 機能が少ない
  - 特に GUI (graphical user interface) が貧弱
  - 営利目的でない以上、ある程度は我慢すべし
  - 無保証(これは商用ソフトも似たようなもの)

## 対象とするデータ

@relation 天気とテニス  
 @attribute 天気予報 (晴, 曇, 雨)  
 @attribute 気温 real  
 @attribute 湿度 real  
 @attribute 風 (強, 弱)  
 @attribute テニス (行う, 止め)

@data  
 晴,29.85,弱,止め  
 晴,27.90,強,止め  
 曇,28.86,弱,行う  
 雨,21.96,弱,行う  
 雨,20.80,弱,行う  
 雨,18.70,強,止め  
 曇,18.65,強,行う  
 晴,22.95,弱,止め  
 晴,21.70,弱,行う  
 雨,24.80,弱,行う  
 晴,24.70,強,行う  
 曇,22.90,強,行う  
 曇,27.75,弱,行う  
 雨,22.91,強,止め

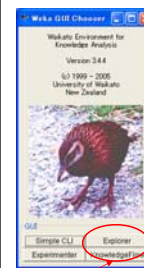
天気とテニス.arff の内容

Excel の表形式で書いたもの

天気予報	温度	湿度	風	テニス
晴	29	85	弱	止め
晴	27	90	強	止め
曇	28	86	弱	行う
雨	21	96	弱	行う
雨	20	80	弱	行う
雨	18	70	強	止め
曇	18	65	強	行う
晴	22	95	弱	止め
晴	21	70	弱	行う
雨	24	80	弱	行う
晴	24	70	強	行う
曇	22	90	強	行う
曇	27	75	弱	行う
雨	22	91	強	止め

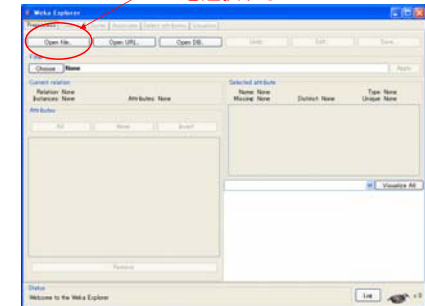
## 起動方法

■ Weka.pif をダブルクリック



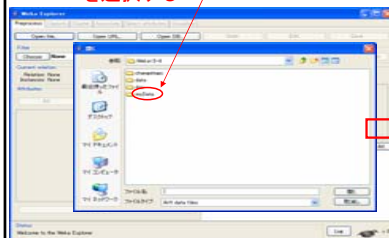
1. クリックしてExplorerを起動

2. クリックしてデータファイルを選択する

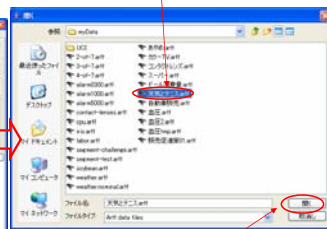


## 対象データファイルの指定

1. クリックしてmyDataフォルダを選択する



2. クリックして天気とテニス.arffファイル(どこかにある)を選択し、

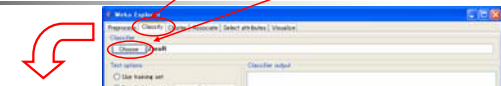


3. 「開く」をクリック、

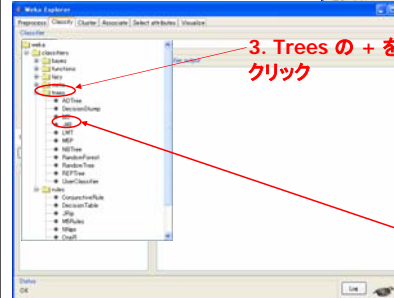
## 決定木の作成(計算)

1. Classify をクリック

2. Choose をクリック

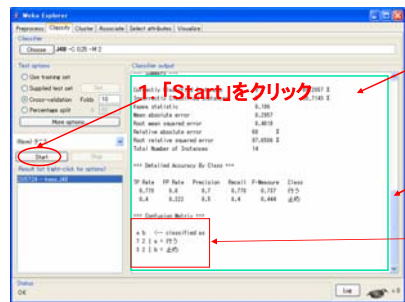


3. Trees の + をクリック



4. j48 をクリック

## 結果の確認



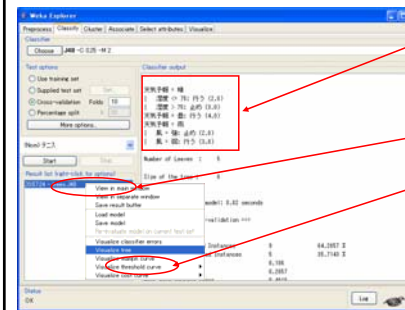
1. 「Start」をクリック

2. 結果はこのウィンドウに表示される

3. このバーを上ドラッグすると、最初の方が見れる

10重クロスバリデーションの結果の総和

## 結果の確認と図示



1. 決定木を文字列で表現したもの

2. この上で「右」クリック

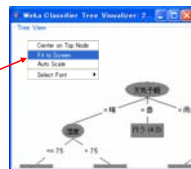
3. 「Visualize tree」の上でクリック

## 図示された木の変形



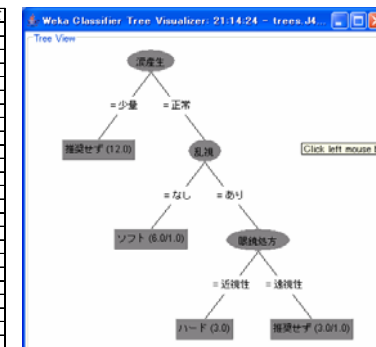
1. マウスマウスをこの角にもってくと、  
に変わる。その状態でドラッグすると、  
このウィンドウの形・大きさが変更できる

2. このスクリーン上で「右」クリック。Fit to Screen をクリックすると、スクリーンの大きさにあった大きさの木になり、Auto Scale でクリックすると木が適度にコンパクトになる。文字の大きさを変えるには Select Font でクリック、木をドラッグすることもできる



## コンタクトレンズの例

年齢	眼鏡処方	近視	遠産生	コンタクトレンズ
若年期	近視性	なし	少量	推奨せず
若年期	近視性	なし	正常	ソフト
若年期	近視性	あり	少量	推奨せず
若年期	近視性	あり	正常	ハード
若年期	遠視性	なし	少量	推奨せず
若年期	遠視性	あり	正常	ソフト
若年期	遠視性	あり	少量	推奨せず
若年期	遠視性	あり	正常	ハード
前老眼期	近視性	なし	少量	推奨せず
前老眼期	近視性	なし	正常	ソフト
前老眼期	近視性	あり	少量	推奨せず
前老眼期	近視性	あり	正常	ハード
前老眼期	遠視性	なし	少量	推奨せず
前老眼期	遠視性	あり	正常	ソフト
前老眼期	遠視性	あり	少量	推奨せず
老眼期	近視性	なし	少量	推奨せず
老眼期	近視性	なし	正常	推奨せず
老眼期	近視性	あり	少量	推奨せず
老眼期	近視性	あり	正常	ハード
老眼期	遠視性	なし	少量	推奨せず
老眼期	遠視性	なし	正常	ソフト
老眼期	遠視性	あり	少量	推奨せず
老眼期	遠視性	あり	正常	推奨せず



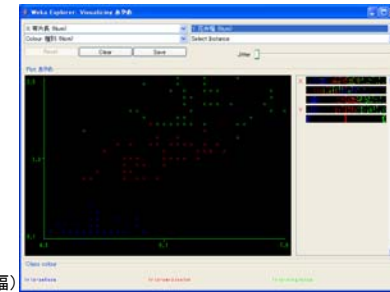
## 分類問題

- 分類問題は、統計的には「判別問題」として扱われるが結構難しい(Excel にはツールがない)
- 人工知能では古典的な課題である
- Fisher (統計学者)が扱った「あやめの分類問題」を考えてみる

## あやめの分類問題

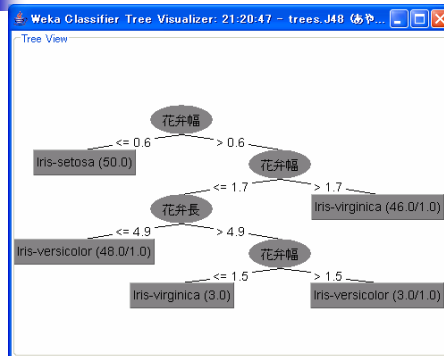
- 萼片長、萼片幅、花弁長、花弁幅とあやめ (setosa, versicolor, virginica の3種)の値が150組。

萼片長	萼片幅	花弁長	花弁幅	種別
5.1	3.5	1.4	0.2	Iris-setosa
4.9	3	1.4	0.2	Iris-setosa
4.7	3.2	1.3	0.2	Iris-setosa
4.6	3.1	1.5	0.2	Iris-setosa
5	3.6	1.4	0.2	Iris-setosa
5.4	3.9	1.7	0.4	Iris-setosa
4.6	3.4	1.4	0.3	Iris-setosa
5	3.4	1.5	0.2	Iris-setosa
4.4	2.9	1.4	0.2	Iris-setosa



(横軸:萼片長、縦軸:花弁幅)

## 分類結果



## 労使間交渉の決着状況

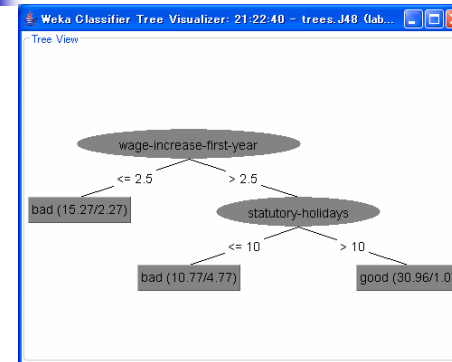
- カナダ労使間交渉の決着状況を、賃金・手当等との組みで表したもの
- 欠損値が多い(ごく普通の状況): 理論的・アルゴリズム的に困難な課題

## 労使間交渉データ

属性	型	1	2	3	40
継続期間	(年数)	1	2	3	2
賃上げ(第1年)	百分率	2	4	4.3	4.5
賃上げ(第2年)	百分率	?	5	4.4	4
賃上げ(第3年)	百分率	?	?	?	?
生活費保証	{none, tcf, tc}	none	tcf	?	none
労働時間/週	時間数	28	35	38	40
年金	{none, ret-allw, empl-cntr}	none	?	?	?
stand-by pay	百分率	?	13	?	?
変則勤務手当	百分率	?	5	4	4
教育手当	{あり, なし}	あり	?	?	?
士曜休業	休日数	11	15	12	12
休暇	{平均以下, 平均, 平均以上}	平均	平均以上	平均以上	平均
長期傷害助成	{あり, なし}	なし	?	?	あり
歯科診療保険助成	{なし, 半分, 完全}	なし	?	完全	完全
死別助成	{あり, なし}	なし	?	?	あり
健康保険助成	{なし, 半分, 完全}	なし	?	完全	半分
対応	{良い, 悪い}	悪い	良い	良い	良い

(縦横がこれまでと逆なので注意)

## 労使間交渉データの結果



## 判断値が数値のとき

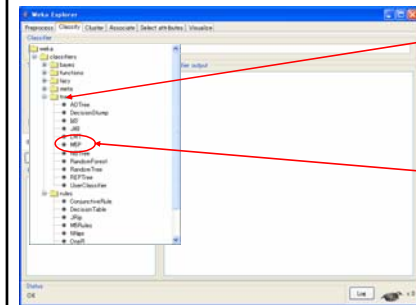
- これまでは、if ... then ... の then のあとがカテゴリ変数(クラス、分類)であった
- 数値のときを、次に扱う
- 回帰と類似であるが、説明変数にカテゴリ変数があること、一次式(直線)で説明できない場合を扱うことが特徴

## ファイルの選択

1. 販売促進01.arffファイル(どこかにある)をクリック、

月	日	曜日	天候	客数	備考
7	1	金	曇り	491	通常
7	2	土	雨	432	通常
7	3	日	晴	514	通常
7	4	月	晴	457	通常
7	5	火	曇り	451	通常
7	6	水	雨	441	通常
7	7	木	雨	604	通常
7	8	金	曇り	467	通常
7	9	土	晴	408	通常
7	10	日	雨	457	通常
7	11	月	雨	484	通常
7	12	火	雨	474	通常
7	13	水	晴	474	通常
7	14	木	晴	666	通常
7	15	金	雨	479	通常
7	16	土	曇り	478	通常
7	17	日	晴	540	通常
7	18	月	晴	497	通常
7	19	火	晴	473	通常
7	20	水	晴	468	通常
7	21	木	晴	875	オートコール
7	22	金	晴	829	オートコール
7	23	土	晴	597	通常
7	24	日	晴	634	通常
7	25	月	曇り	478	通常
7	26	火	曇り	480	通常
7	27	水	晴	408	通常
7	28	木	晴	544	通常
7	29	金	雨	365	通常
7	30	土	晴	380	通常
7	31	日	晴	448	通常

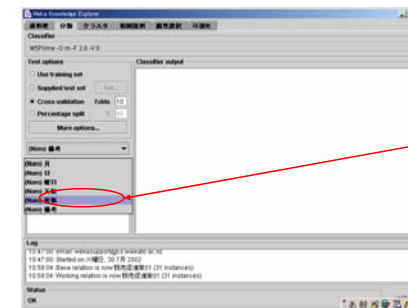
## 使うアルゴリズムの選択



1. Treeの右にある+をクリック

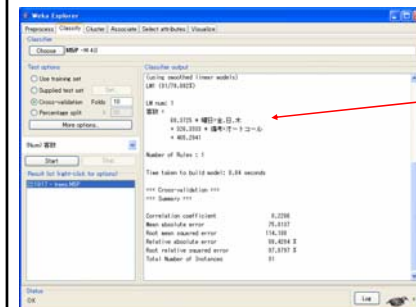
2. M5P というのをを選択する

## 被説明変数の指定



1. 「客数」の上でクリック  
 黙っているとデータ(表)のなかの最も右の属性が用いられる。  
 今回は、「最も右」ではないのでここで指定する

## 結果の解析



客数 =  
 $60.3725 * \text{曜日} = \text{金, 日, 木}$   
 $+ 326.3333 * \text{備考} = \text{オートコール}$   
 $+ 465.2941$

オートコールを行った方が客数が増加することがわかる

## 血圧の測定データ

血圧の測定データ				血圧の測定データ				血圧の測定データ				血圧の測定データ			
No.	性別	年齢	血圧 (mmHg)	No.	性別	年齢	血圧 (mmHg)	No.	性別	年齢	血圧 (mmHg)	No.	性別	年齢	血圧 (mmHg)
1	男	25	110	1	男	25	110	1	男	25	110	1	男	25	110
2	女	30	120	2	女	30	120	2	女	30	120	2	女	30	120
3	男	35	130	3	男	35	130	3	男	35	130	3	男	35	130
4	女	40	140	4	女	40	140	4	女	40	140	4	女	40	140
5	男	45	150	5	男	45	150	5	男	45	150	5	男	45	150
6	女	50	160	6	女	50	160	6	女	50	160	6	女	50	160
7	男	55	170	7	男	55	170	7	男	55	170	7	男	55	170
8	女	60	180	8	女	60	180	8	女	60	180	8	女	60	180
9	男	65	190	9	男	65	190	9	男	65	190	9	男	65	190
10	女	70	200	10	女	70	200	10	女	70	200	10	女	70	200
11	男	75	210	11	男	75	210	11	男	75	210	11	男	75	210
12	女	80	220	12	女	80	220	12	女	80	220	12	女	80	220
13	男	85	230	13	男	85	230	13	男	85	230	13	男	85	230
14	女	90	240	14	女	90	240	14	女	90	240	14	女	90	240
15	男	95	250	15	男	95	250	15	男	95	250	15	男	95	250
16	女	100	260	16	女	100	260	16	女	100	260	16	女	100	260
17	男	105	270	17	男	105	270	17	男	105	270	17	男	105	270
18	女	110	280	18	女	110	280	18	女	110	280	18	女	110	280
19	男	115	290	19	男	115	290	19	男	115	290	19	男	115	290
20	女	120	300	20	女	120	300	20	女	120	300	20	女	120	300

## Weka による分析結果

```

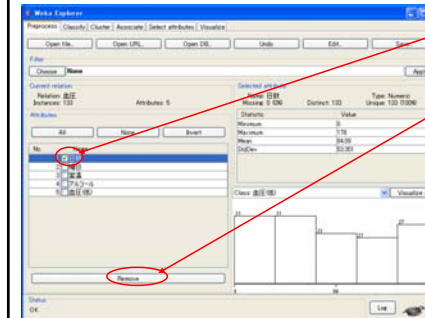
Classifier output
日数 <= 99 : LM1 (77/122,613)
日数 > 99 : LM2 (56/99,7423)

LM num: 1
血圧(低) =
+ 0.7918 * 曜日=金,日,木,水,月,火
+ 4.1259 * 曜日=日,木,水,月,火
+ 0.2459 * 曜日=水,月,火
+ 0.8361 * 室温
+ 0.336 * アルコール=少々,なし
+ 101.7857

LM num: 2
血圧(低) =
+ 0.1092 * 日数
+ 0.9493 * 曜日=金,日,木,水,月,火
+ 0.3167 * 曜日=水,月,火
+ 0.9244 * 室温
+ 2.19 * アルコール=少々,なし
+ 88.6543

Number of Rules : 2
    
```

## 日数をはずす



1. 日数のチェックボックスにチェック
2. 属性を remove するためクリック
3. 「分類」で M5Prime を Start

## 日数をはずした場合の結果

```

Classifier output
LM num: 1
血圧(低) =
+ 4.0606 * 曜日=金,日,木,水,月,火
+ 1.8615 * 曜日=木,水,月,火
+ 0.4319 * 室温
+ 2.2014 * アルコール=少々,なし
+ 93.9143

Number of Rules : 1

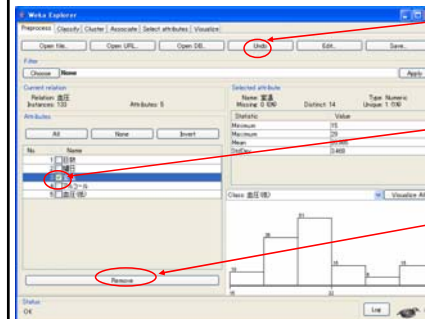
Time taken to build model: 0.1 seconds

*** Cross-validation ***
*** Summary ***

Correlation coefficient: 0.1648
Mean absolute error: 4.2615
Root mean squared error: 6.0325
Relative absolute error: 106.6214 %
Root relative squared error: 102.3489 %
Total Number of Instances: 123
    
```

Correlation coefficient 0.1648

## 室温をはずす



1. Undo をクリックすると日数が戻ってくる
2. 室温にチェックをつける
3. Removeする

## 室温をはずした場合の結果

```
Classifier output
日数 <= 93 : LM1 (77/124.576%)
日数 > 93 : LM2 (56/84.261%)

LM num: 1
血压(低) =
-0.0033 * 日数
+ 0.6118 * 曜日=金,日,木,水,月,火
+ 3.5396 * 曜日=日,木,水,月,火
+ 0.3149 * 曜日=木,水,月,火
+ 1.9447 * 曜日=月,火
+ 0.3771 * アルコール=少々,なし
+ 88.5818

LM num: 2
血压(低) =
0.0501 * 日数
+ 0.7928 * 曜日=金,日,木,水,月,火
+ 0.408 * 曜日=木,水,月,火
+ 3.2053 * アルコール=少々,なし
+ 79.8907

Number of Rules : 2
```

Correlation coefficient 0.2719

日数 <= 93 : LM1 (77/124.576%)  
日数 > 93 : LM2 (56/84.261%)

LM1: 血压(低) =  
-0.0033 \* 日数  
+ 0.6118 \* 曜日=金,日,木,水,月,火  
+ 3.5396 \* 曜日=日,木,水,月,火  
+ 0.3149 \* 曜日=木,水,月,火  
+ 1.9447 \* 曜日=月,火  
+ 0.3771 \* アルコール=少々,なし  
+ 88.5818

LM2: 血压(低) =  
0.0501 \* 日数  
+ 0.7928 \* 曜日=金,日,木,水,月,火  
+ 0.408 \* 曜日=木,水,月,火  
+ 3.2053 \* アルコール=少々,なし  
+ 79.8907

## 日数と室温との関係

```
Classifier output
日数 <= 111.5 : LM1 (88/67.068%)
日数 > 111.5 :
| 日数 <= 162.5 : LM2 (34/55.335%)
| 日数 > 162.5 : LM3 (11/16.813%)

LM num: 1
室温 =
0.007 * 日数
+ 18.7126

LM num: 2
室温 =
0.0513 * 日数
+ 16.6505

LM num: 3
室温 =
0.0785 * 日数
+ 13.5047
```

Correlation coefficient 0.8465

日数 <= 111.5 : LM1 (88/67.068%)  
日数 > 111.5 :  
| 日数 <= 162.5 : LM2 (34/55.335%)  
| 日数 > 162.5 : LM3 (11/16.813%)

LM1 室温 = 0.007 \* 日数 + 18.7126  
LM2 室温 = 0.0513 \* 日数 + 16.6505  
LM3 室温 = 0.0785 \* 日数 + 13.5047

## 日数と室温をはずすと

```
Classifier output
LM num: 1
血压(低) =
3.2053 * 曜日=金,日,木,水,月,火
+ 1.4925 * 曜日=木,水,月,火
+ 2.4464 * アルコール=少々,なし
+ 85.4822

Number of Rules : 1
Time taken to build model: 0.09 seconds

=== Cross-validation ===
=== Summary ===

Correlation coefficient -0.0089
Mean absolute error 4.921
Root mean squared error 6.1926
Relative absolute error 105.9657 %
Root relative squared error 104.2919 %
Total Number of Instances 129
```

残りの属性(曜日と前日のアルコール摂取量)ではうまく説明できないことがわかる

## 「血压」の総合的な結論

- 日数がたつにつれ、血压が上昇している
- しかし、それは日数がたったからか、気温が上昇したからかはわからない
- 土曜日に低い傾向はあるが、確信できず
- 前日のアルコール摂取量で低い傾向はあるが、確信度はもっと低い