

情報意味論(3) Wekaの紹介

Preslav Nakov (October 6, 2004)

<http://www.sims.berkeley.edu/courses/is290-2/f04/lectures/lecture11.ppt>

Eibe Frank

<http://prdownloads.sourceforge.net/weka/weka.ppt>

1

- **WEKA: Explorer**
- WEKA: Experimenter
- WEKA: 使ってみよう

2

WEKA: 好奇心旺盛な飛べない鳥



Copyright: Martin Kramer (mkramer@wxs.nl), University of Waikato, New Zealand

Slide adapted from Eibe Frank's

3

WEKA: 用語

Weka で用いる単語は、多少、普通の用法と異なる:

- **Attribute:** 属性, feature (普通)
- **Relation:** 事例の集合, collection of examples
- **Instance:** その時点で使用中の事例集合
- **Class:** クラス、範疇、category

4

WEKA: The Software Toolkit

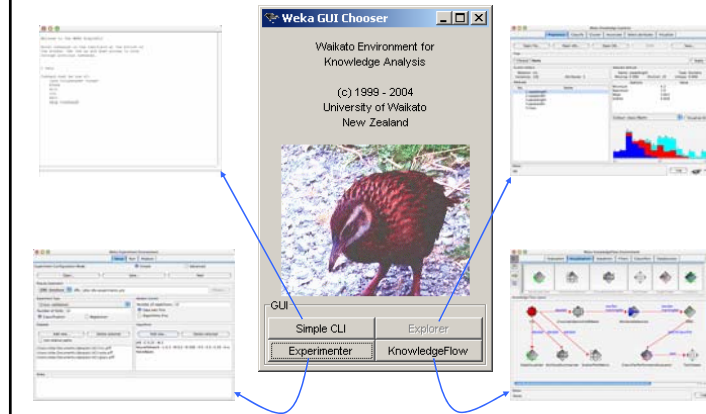
<http://www.cs.waikato.ac.nz/ml/weka>

- Java で書かれた、機械学習/データマイニングソフト
- GNU ライセンス
- 研究、教育、応用に用いられる
- Witten & Frank 著 "Data Mining" で使用
- 主な機能:
 - データ前処理ツール
 - 学習アルゴリズム
 - 評価手法
 - グラフィカルインタフェース (データ可視化含む)
 - 学習アルゴリズムの比較

Slide adapted from Eibe Frank's

5

WEKA GUI Chooser `java -Xmx1000M -jar weka.jar`



Slide adapted from Eibe Frank's

6

Explorer: データの前処理

- WEKA はデータがインポートできる:
 - ファイルから: ARFF, CSV, C4.5, binary
 - URL から
 - SQL データベースから (JDBCを用いて)
- 前処理ツール(フィルター)が使われるのは:
 - 離散化 discretization, 正規化 normalization, resampling, 属性選択 attribute selection, 属性の変換と結合, etc.

Slide adapted from Eibe Frank's

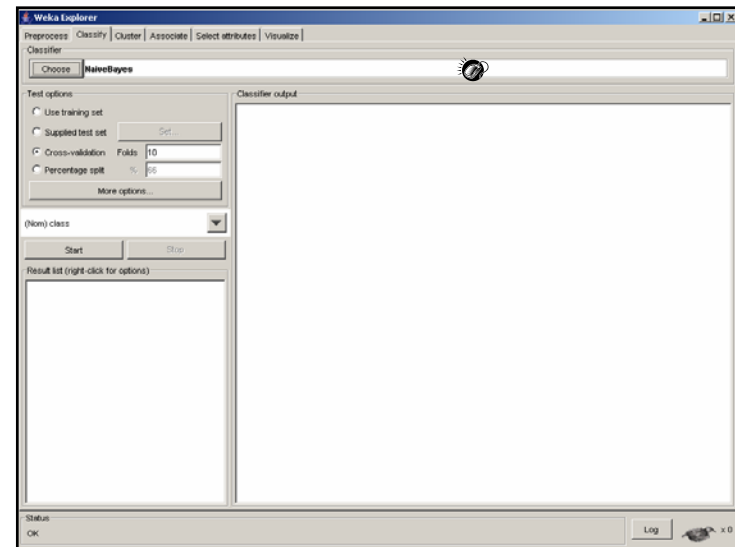
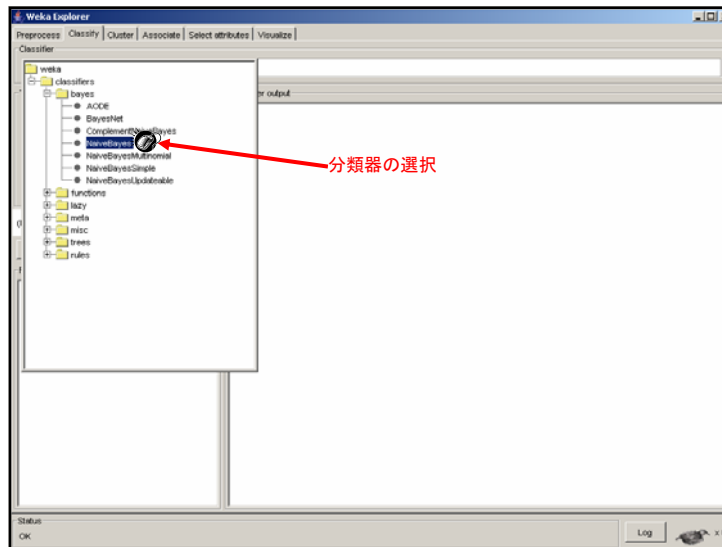
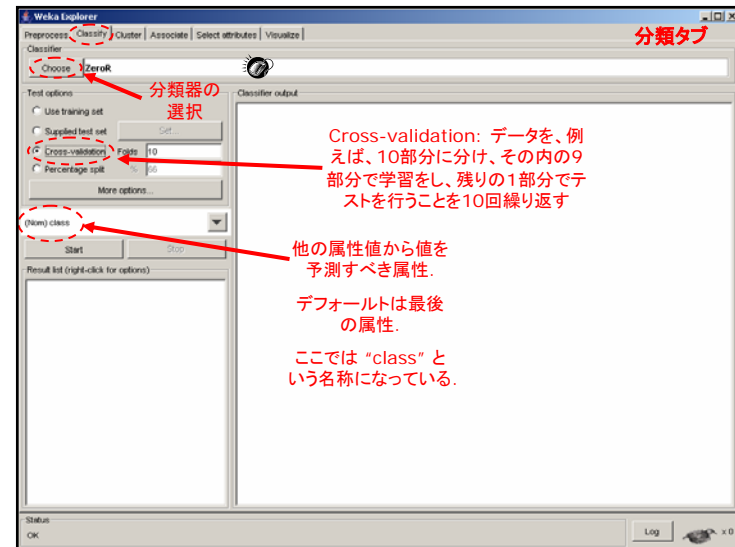
7

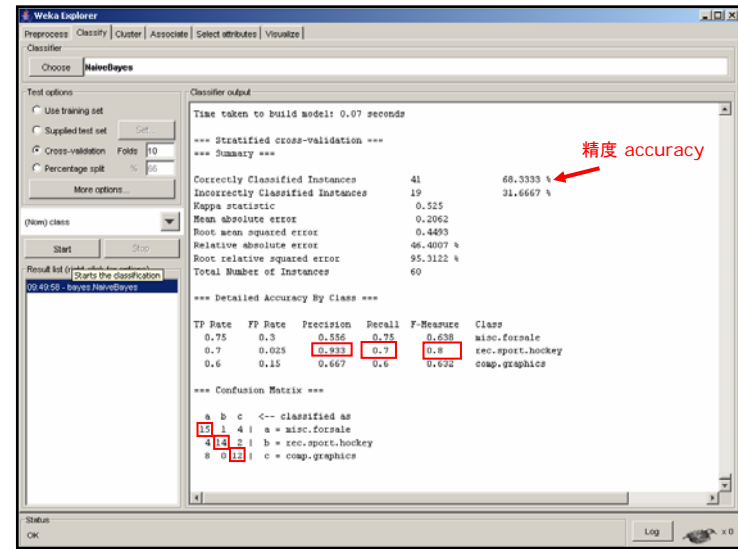
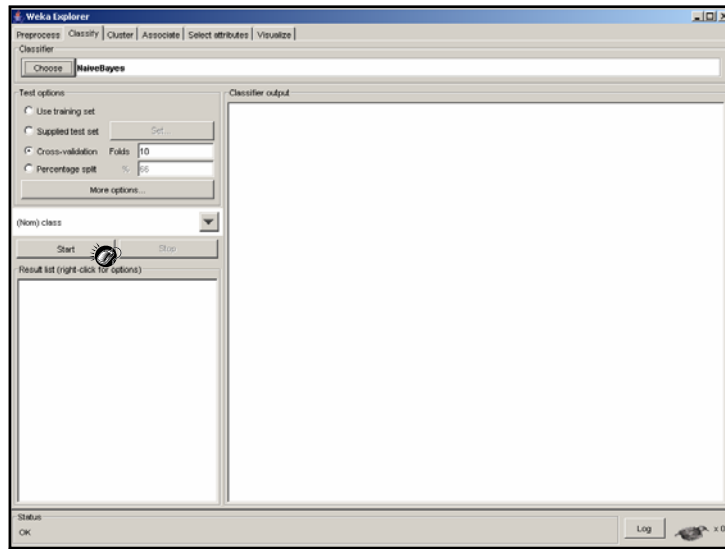
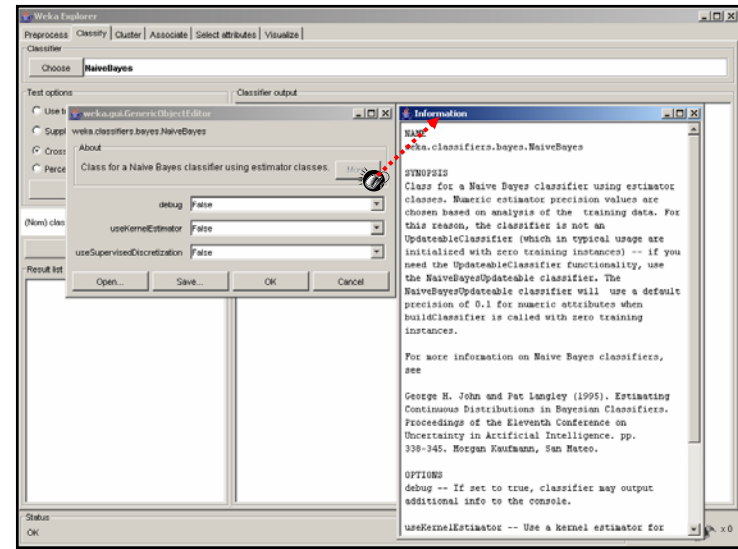
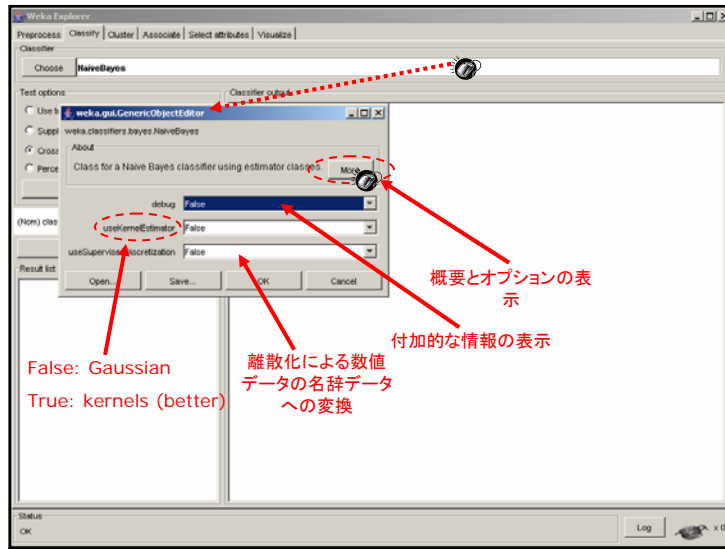
Explorer: 分類

- Weka の *Classifiers* は:
 - 分類 classification (名辞で表すクラスの予測)
 - 回帰 regression (数値的な量の予測)
- 使用する学習アルゴリズム:
 - Naive Bayes, 決定木 decision trees, kNN, support vector machines, ニューラルネットワーク, logistic regression, 等.
- メタ分類:
 - 単独で使うものではない
 - 学習アルゴリズムと組み合わせて使用
 - 例: boosting, bagging 等.

Slide adapted from Eibe Frank's

9





Confusion matrix

- 予測した分類と本当の分類との違いを数値化したもの
- これから様々な値が得られる:

- accuracy: $(a+d)/(a+b+c+d)$
- recall: $d/(c+d) \Rightarrow R$
- precision: $d/(b+d) \Rightarrow P$
- F-measure: $2PR/(P+R)$
- false positive (FP) rate: $b/(a+b)$
- true negative (TN) rate: $a/(a+b)$
- false negative (FN) rate: $c/(c+d)$

		予測値	
		-	+
真	-	a	b
	+	c	d

- これらは2クラス以上にも適用できる

17

予測出力

Classifier output

Time taken to build model: 0.19 seconds

=== Stratified cross-validation ===

=== Summary ===

Classifier evaluation options

Incorrect: 0.525

Wrong att: 0.2062

Root mean: 0.4492

Relative: 6.4007 %

Root relc: 5.3122 %

Total Mis: 10

=== Detail ===

TP Rate: 0.75

FP Rate: 0.7

FN Rate: 0.6

Random seed for XVal / % Split: 0.620

F-Measure: 0.632

Class: misc.forzaile, rec.sport.hockey, comp.graphics

=== Confusion Matrix ===

a b c <-- classified as

15 1 4 | a = misc.forzaile

4 14 2 | b = rec.sport.hockey

0 0 12 | c = comp.graphics

各事例に対しそう分類する確からしさ(確率)を出力する

Predictions Output

=== Predictions on test data ===

inst#	actual	predicted	error	probability distribution
1	3:comp.gra	1:misc.for	+ *1	0 0 0
2	3:comp.gra	3:comp.gra	0	0 0 *1
3	2:rec.spor	3:comp.gra	+ 0	0.212 *0.788
4	2:rec.spor	2:rec.spor	0	*1 0 0
5	1:misc.for	1:misc.for	*1	0 0 0
6	1:misc.for	1:misc.for	*1	0 0 0
1	3:comp.gra	1:misc.for	+ *0.992	0 0.008
2	3:comp.gra	3:comp.gra	0	0 0 *1
3	2:rec.spor	1:misc.for	+ *1	0 0 0
4	2:rec.spor	1:misc.for	+ *1	0 0 0
5	1:misc.for	1:misc.for	*1	0 0 0
6	1:misc.for	1:misc.for	*1	0 0 0
1	3:comp.gra	3:comp.gra	0	0 *1
2	3:comp.gra	3:comp.gra	0	*1 0 0
3	2:rec.spor	2:rec.spor	0	*1 0 0
4	2:rec.spor	2:rec.spor	0	*1 0 0
5	1:misc.for	1:misc.for	*1	0 0 0
6	1:misc.for	1:misc.for	*1	0 0 0
1	3:comp.gra	1:misc.for	+ *0.96	0 0.04

誤分類事例への確率付与の例:
予測1
正解3

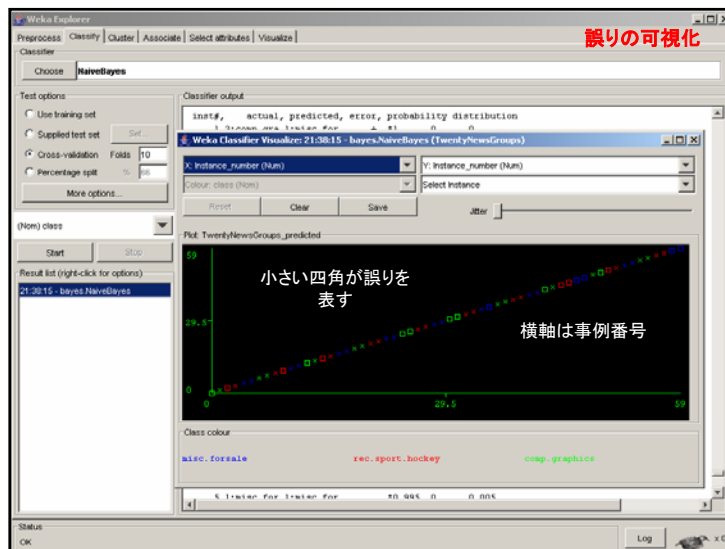
Naive Bayesは条件付独立性を無条件に仮定しているため、通常は過信ぎみの結果を出す(結果の正誤とは関係なく)

誤りの可視化

Classifier output

inst#, actual, predicted, error, probability distribution

1	3:comp.gra	1:misc.for	+ *1	0 0 0
2	3:comp.gra	3:comp.gra	0	0 0 *1
3	2:rec.spor	3:comp.gra	+ 0	0.212 *0.788
4	2:rec.spor	2:rec.spor	0	*1 0 0
5	1:misc.for	1:misc.for	*1	0 0 0
6	1:misc.for	1:misc.for	*1	0 0 0
1	3:comp.gra	3:comp.gra	0	0 *1
2	3:comp.gra	3:comp.gra	0	*1 0 0
3	2:rec.spor	2:rec.spor	0	*1 0 0
4	2:rec.spor	1:misc.for	+ *1	0 0 0
5	1:misc.for	1:misc.for	*1	0 0 0
6	1:misc.for	1:misc.for	*1	0 0 0
1	3:comp.gra	1:misc.for	+ *0.992	0 0.008
2	3:comp.gra	3:comp.gra	0	0 *1
3	2:rec.spor	2:rec.spor	0	*1 0 0
4	2:rec.spor	1:misc.for	+ *1	0 0 0
5	1:misc.for	1:misc.for	*1	0 0 0
6	1:misc.for	1:misc.for	*1	0 0 0
1	3:comp.gra	1:misc.for	+ *0.96	0 0.04
2	3:comp.gra	3:comp.gra	0	*1 0 0
3	2:rec.spor	2:rec.spor	0	*1 0 0
4	2:rec.spor	2:rec.spor	0	*1 0 0
5	1:misc.for	1:misc.for	*1	0 0.004

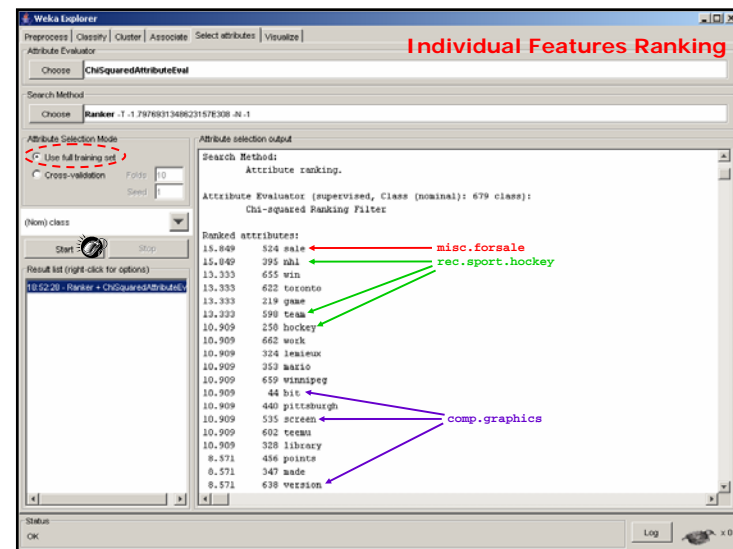
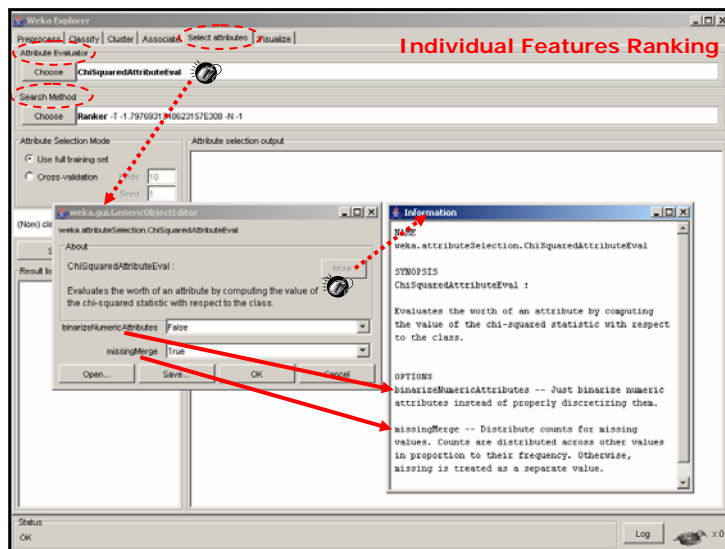


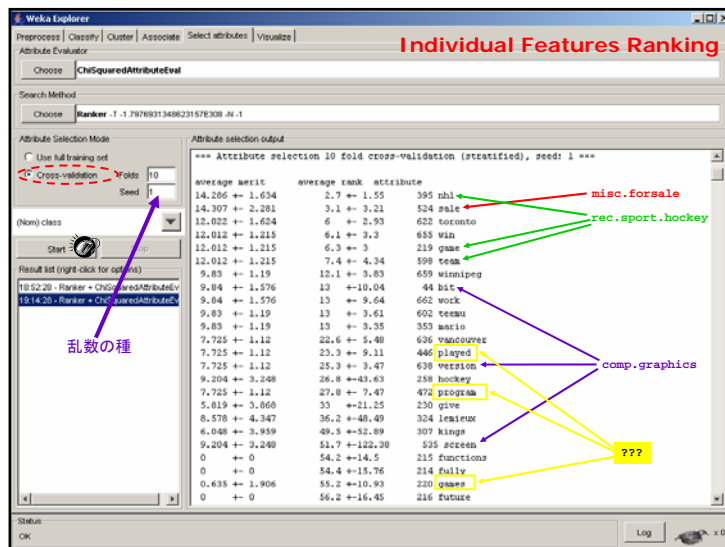
Explorer: 属性選択

- 予測力の高い属性を探す
 - 予測力の低い属性を、予め、排除したい
- 2種類の方法:
 - 探索法:
 - best-first, forward selection, random, 網羅探索
 - exhaustive, 遺伝的アルゴリズム genetic algorithm, ranking
 - 評価法:
 - information gain, χ 自乗, etc.
- WEKA は (殆どの) 任意の組合せができる

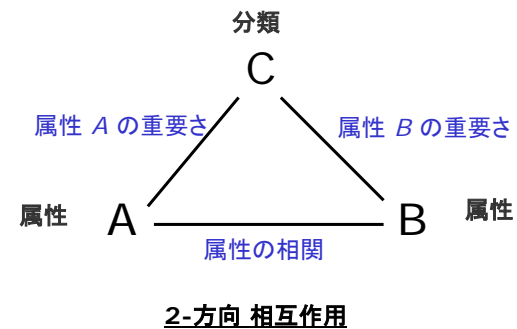
Slide adapted from Eibe Frank's

22





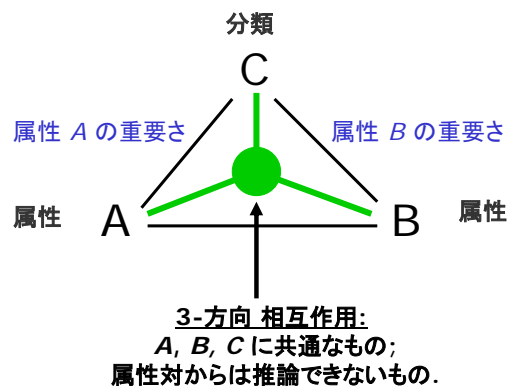
属性間の相互作用



Slide adapted from Jakulin, Bratko, Smrke, Demšar and Zupan's

26

属性間の相互作用



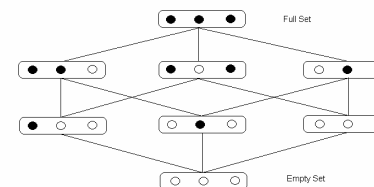
Slide adapted from Jakulin, Bratko, Smrke, Demšar and Zupan's

27

属性部分集合の選択

■ 問題の例図

- 全体集合 Full set
- 空集合 Empty set
- 数え上げ Enumeratic



■ 探索

- 網羅的/完全 (数え上げ/branch&bounding)
- ヒューリスティック (逐次的に 前進方向/後退方向)
- 確率的 stochastic (生成/評価)
- 個々の属性や部分集合の生成/評価

Slide adapted from Guozhu Dong's

28

属性部分集合選択

Choose: **CfsSubsetEval**

Search Method: Choose: **BestFirst -D1 -N5**

Attribute: **weka.attributeSelection.CfsSubsetEval**

(Nom) class: **locallyPredictive**

Result list: **MissingSeparate**

Information:

NAME
weka.attributeSelection.CfsSubsetEval

SYNOPSIS
CfsSubsetEval :

Evaluates the worth of a subset of attributes by considering the individual predictive ability of each feature along with the degree of redundancy between them.

Subsets of features that are highly correlated with the class while having low intercorrelation are preferred.

OPTIONS

locallyPredictive -- Identify locally predictive attributes. Iteratively adds attributes with the highest correlation with the class as long as there is not already an attribute in the subset that has a higher correlation with the attribute in question

MissingSeparate -- Treat missing as a separate value. Otherwise, counts for missing values are distributed across other values in proportion to their frequency.

0 + 0 666.2 + 1.69
0 + 0 669.4 + 0.88
0 + 0 669.8 + 4.15
0 + 0 670.2 + 6.34
0 + 0 672.7 + 3.69
0 + 0 673.3 + 2.41
0 + 0 673.4 + 3.07
0 + 0 672.8 + 6.01
0 + 0 676.2 + 1.47

属性部分集合選択

Choose: **CfsSubsetEval**

Search Method: Choose: **BestFirst -D1 -N5**

Attribute Selection Mode: **Use full training set**

Cross-validation: **Folds: 10**

(Nom) class: **misc.forsale**

Attribute selection output:

Selected attributes: 44, 57, 124, 219, 230, 258, 283, 328, 353, 395, 472, 524, 535, 564, 593, 622, 636, 638, 646, ...

bit
brand
cup
game
give
hockey
integrated
library
mario
nhl
program
sale
screen
software
system
toronto
vancouver
verizon
wanted
work
worst

misc.forsale
rec.sport.hockey
comp.graphics

17,309 個の部分集合を検討し
21 個の属性が選択された

選択した属性の保存

Choose: **ClassifierSubsetEval -B weka.classifiers.bayes.NaiveBayes -T -H "Click to set hold out or test instances" --**

Search Method: Choose: **BestFirst -D1 -N5**

Attribute Selection Mode: **Use full training set**

Cross-validation: **Folds: 10**

(Nom) class: **misc.forsale**

Attribute selection output:

Selected attributes: 44, 57, 124, 219, 230, 258, 283, 328, 353, 395, 472, 524, 535, 564, 593, 622, 636, 638, 646, ...

bit
brand
cup
game
give
hockey
integrated
library
mario
nhl
program
sale
screen
software
system
toronto
vancouver
verizon
wanted
work
worst

このタブから出来ることは
バッファをテキストファイルで
保存すること. 使いやすいと
はいえない...

属性選択は、前処理のス
テップでもできる...
(次のスライド参照)

View in main window
View in separate window
Save result buffer
Visualize reduced data

前処理における属性選択

Files: **CfsSubsetEval -S "weka.attributeSelection.BestFirst -D1 -N5"**

Attribute selection output:

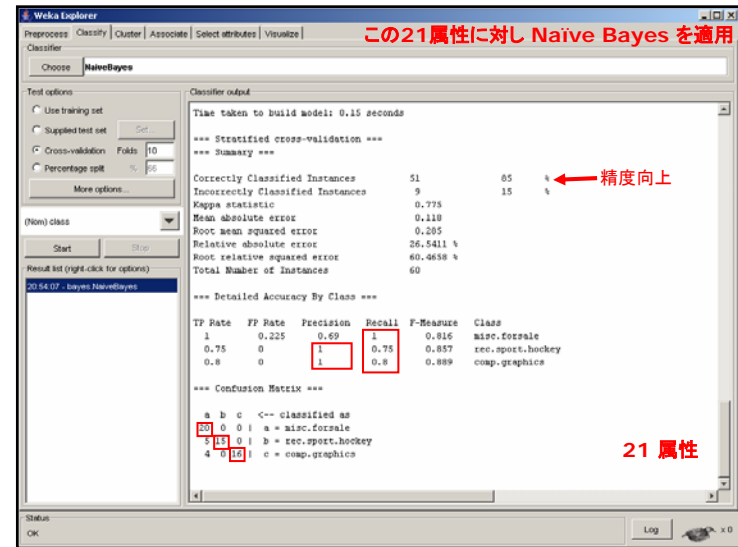
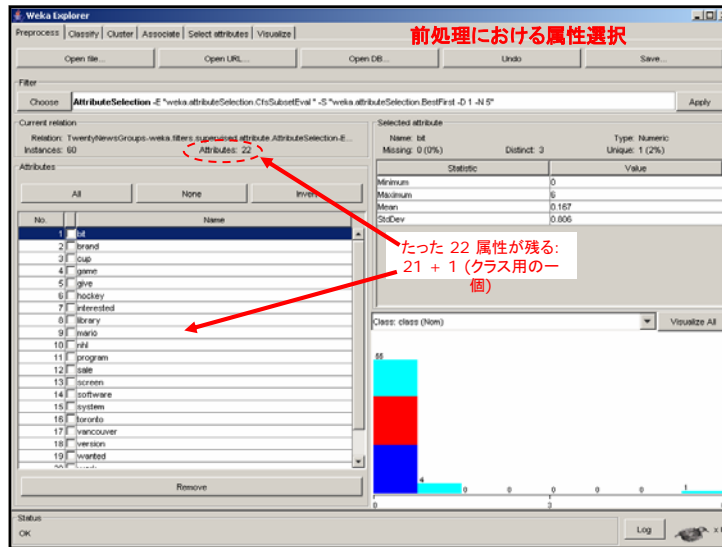
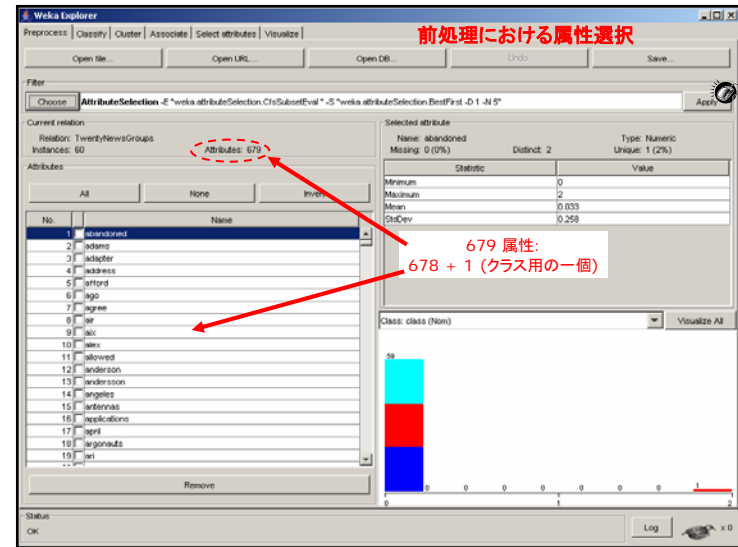
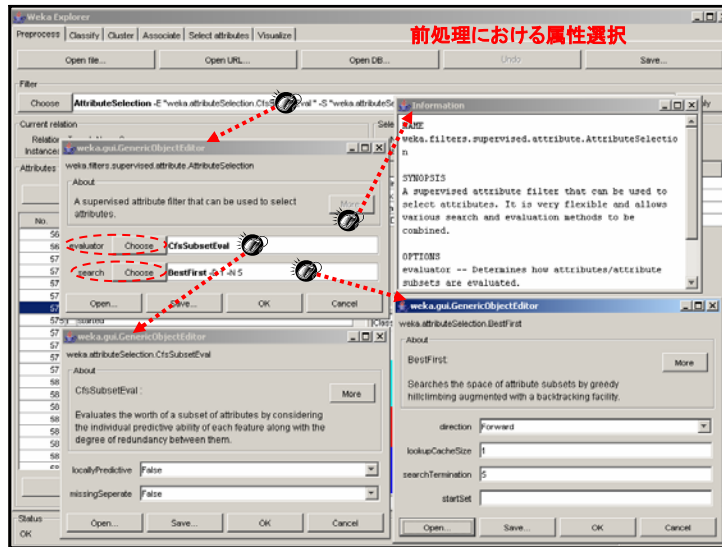
Selected attribute:

Name	star	Distinct	2	Type	Numeric
Missing	0 (0%)				
Statistic					
Minimum	0				
Maximum	2				
Mean	0.887				
StdDev	0.312				

Class: class (Nom)

Visualize All

57



Weka Explorer (再掲) 全属性を用いた Naive Bayes

Classifier: NaiveBayes

Test options: Use training set, Supplied test set, Cross-validation Folds 10, Percentage split % 60

Classifier output

Time taken to build model: 0.07 seconds

=== Stratified cross-validation ===
 === Summary ===

Correctly Classified Instances	41	68.3333 %
Incorrectly Classified Instances	19	31.6667 %
Kappa statistic	0.525	
Mean absolute error	0.2062	
Root mean squared error	0.4493	
Relative absolute error	46.4007 %	
Root relative squared error	95.3122 %	
Total Number of Instances	60	

← 精度

=== Detailed Accuracy By Class ===

TP Rate	FP Rate	Precision	Recall	F-Measure	Class
0.75	0.3	0.556	0.75	0.638	misc.foxmale
0.7	0.025	0.833	0.7	0.8	rec.sport.hockey
0.6	0.15	0.667	0.6	0.652	comp.graphics

=== Confusion Matrix ===

a	b	c	←- classified as
17	1	4	a = misc.foxmale
4	13	2	b = rec.sport.hockey
8	0	12	c = comp.graphics

全679 属性使用 (スライド再掲)

重要なアルゴリズム

- WEKA はアルゴリズムに時々変わった名称をつけている
- 主なアルゴリズム:
 - Naive Bayes: `weka.classifiers.bayes.NaiveBayes`
 - Perceptron: `weka.classifiers.functions.VotedPerceptron`
 - Winnow: `weka.classifiers.functions.winnow`
 - Decision tree: `weka.classifiers.trees.J48`
 - Support vector machines: `weka.classifiers.functions.SMO`
 - k nearest neighbor: `weka.classifiers.lazy.IBK`
- これらのものには、古典的なアルゴリズムではなくより洗練されたものがある
 - e.g. WEKA には古典的な Naive Bayes は見つけられなかった (5 種類もの実装されているのに).

38

- WEKA: Explorer
- WEKA: Experimenter**
- WEKA: 使ってみよう

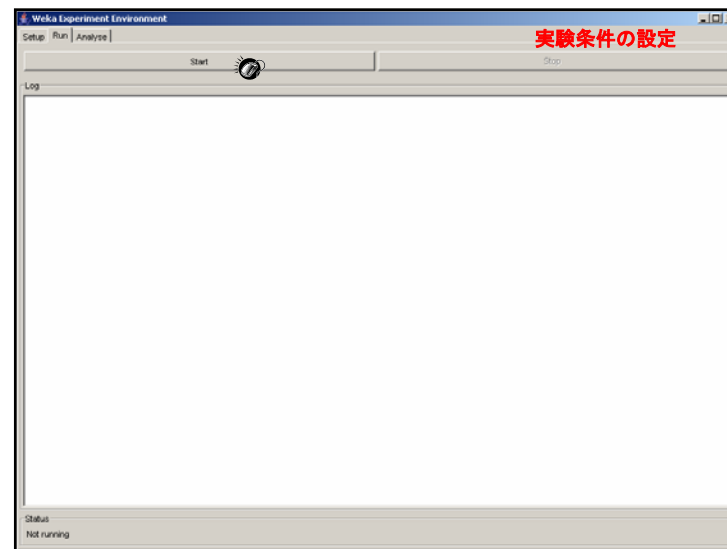
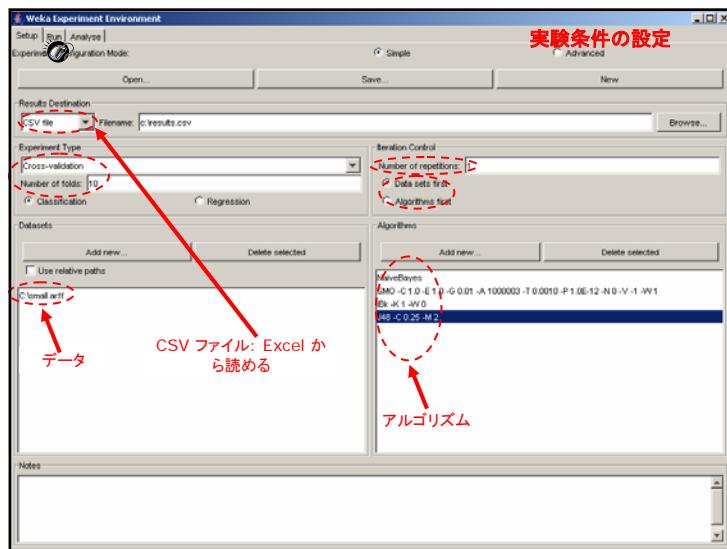
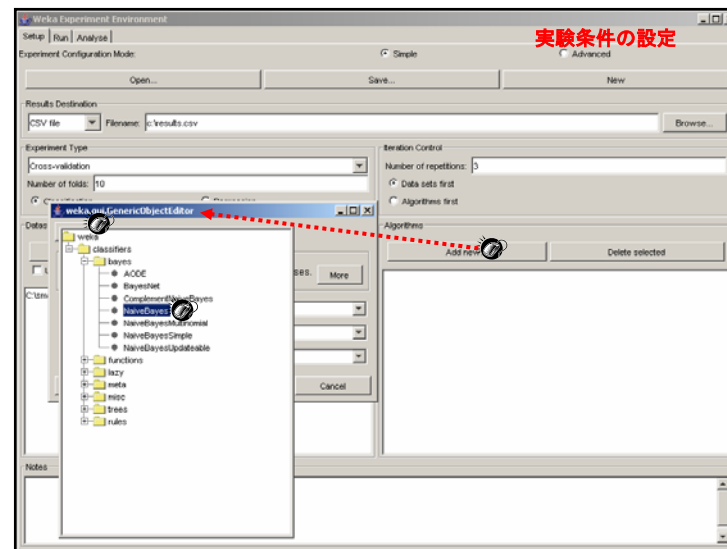
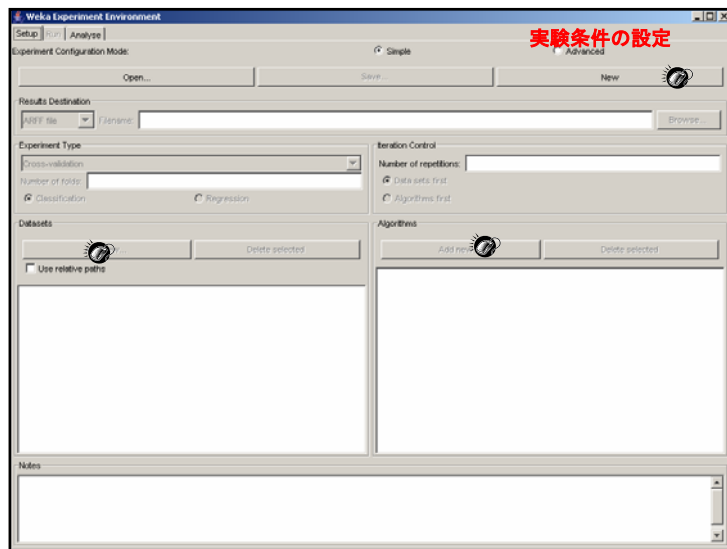
39

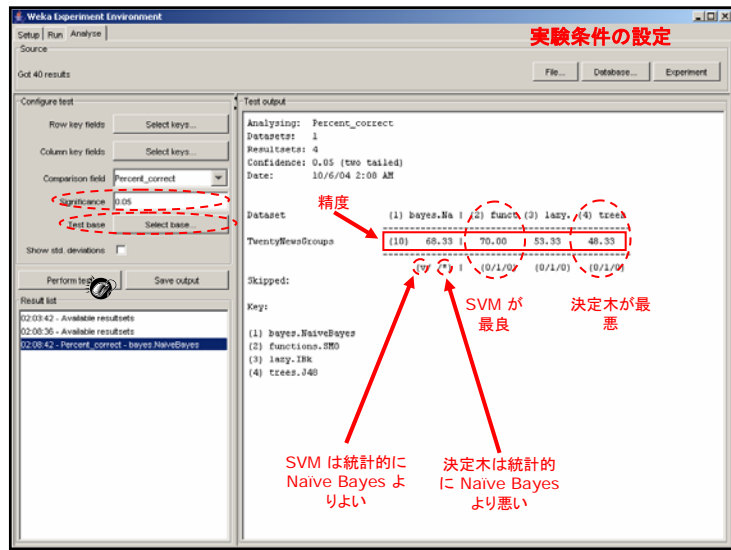
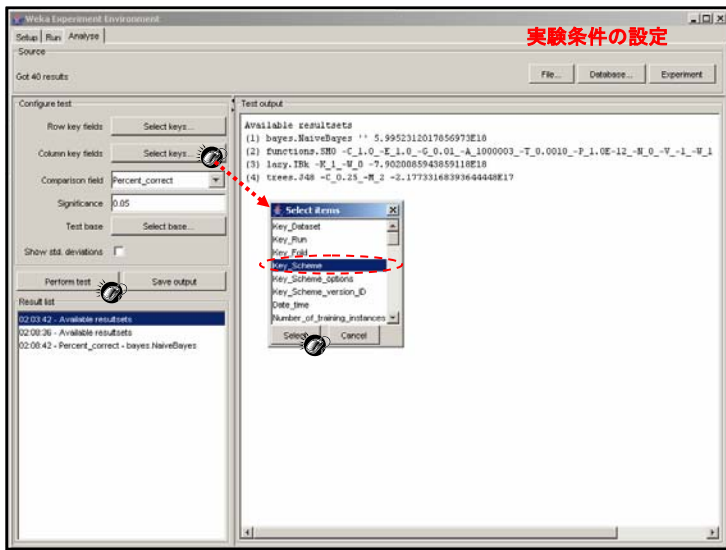
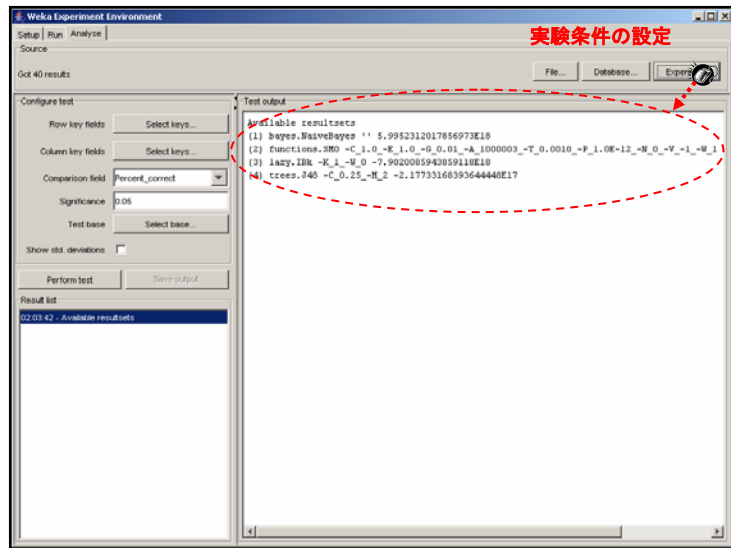
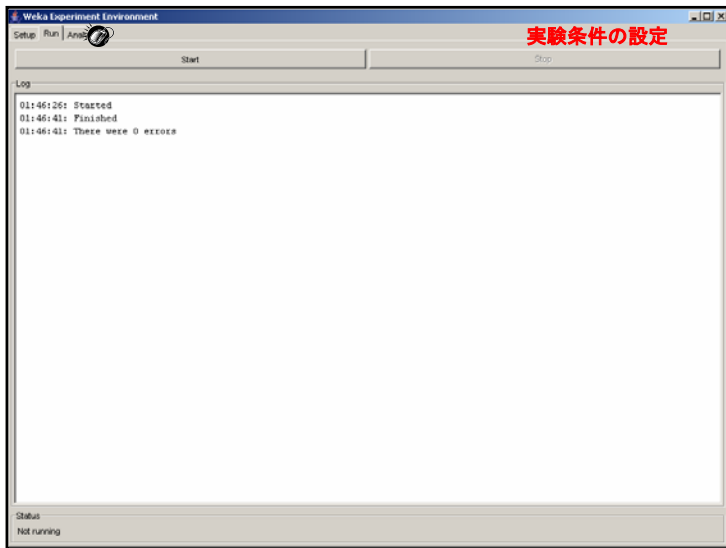
実験を行う

- Experimenter* を用いると、異なる学習アルゴリズムを比較することが容易にできる
- 扱う問題:
 - 分類
 - 回帰
- 結果: ファイルに書き込まれる
- 評価方法の種類:
 - cross-validation
 - 学習曲線 learning curve
 - hold-out
- パラメータを変えて繰り返すことができる
- 有意性検定が組み込み

Slide adapted from Eibe Frank's

40





Microsoft Excel - results.csv

実験: Excel

File Edit View Insert Format Tools Data Window Help

Key Dataset

Key_Dataset	Key_Run	Key_Fold	Key_Sch	Key_Sch	Key_Sch	Date_time	Number_of	Number_of	Number_of	Number_of	Number_of	Percent_of	Percent_of
TwentyNev	1	1	weka.class*	BE+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	2	weka.class*	BE+18	2.00E+07	54	6	5	1	0	83.33333	16.66667	33.33333
TwentyNev	1	3	weka.class*	BE+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	4	weka.class*	BE+18	2.00E+07	54	6	6	0	0	100	0	0
TwentyNev	1	5	weka.class*	BE+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	6	weka.class*	BE+18	2.00E+07	54	6	3	3	0	50	50	50
TwentyNev	1	7	weka.class*	BE+18	2.00E+07	54	6	6	0	0	100	0	0
TwentyNev	1	8	weka.class*	BE+18	2.00E+07	54	6	1	5	0	16.66667	83.33333	33.33333
TwentyNev	1	9	weka.class*	BE+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	10	weka.class*	BE+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	1	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	2	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	5	1	0	83.33333	16.66667	33.33333
TwentyNev	1	3	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	4	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	5	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	6	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	7	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	8	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	9	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	10	weka.class-C 1.0-E	9.22E+18	2.00E+07	54	6	4	2	0	83.33333	16.66667	33.33333
TwentyNev	1	1	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	2	4	0	33.33333	66.66667	33.33333
TwentyNev	1	2	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	3	3	0	50	50	50
TwentyNev	1	3	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	3	3	0	50	50	50
TwentyNev	1	4	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	3	3	0	50	50	50
TwentyNev	1	5	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	3	3	0	50	50	50
TwentyNev	1	6	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	5	1	0	83.33333	16.66667	33.33333
TwentyNev	1	7	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	3	3	0	50	50	50
TwentyNev	1	8	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	2	4	0	33.33333	66.66667	33.33333
TwentyNev	1	9	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	3	3	0	50	50	50
TwentyNev	1	10	weka.class-K 1-W 0	-7.9E+18	2.00E+07	54	6	5	1	0	83.33333	16.66667	33.33333
TwentyNev	1	1	weka.class-C 0.25-M	-2.2E+17	2.00E+07	54	6	1	5	0	16.66667	83.33333	33.33333
TwentyNev	1	2	weka.class-C 0.25-M	-2.2E+17	2.00E+07	54	6	1	5	0	16.66667	83.33333	33.33333

結果は CSV ファイルに書き出すことができ、それは Excel に読み込むことができる