

Data Mining Tutorial

Mark A. Austin

University of Maryland

austin@umd.edu

ENCE 688P, Fall Semester 2021

October 16, 2021

Working with Weka

Data Mining

Examples

Example 1. Will Customer Buy Computer?

Input datafile (arff format)

```

1  % =====
2  % ENCE 688P: Classification for Buy Computer?
3  % =====
4
5  @relation 'computer'
6  @attribute id real
7  @attribute age { young, middle, senior}
8  @attribute income { low, medium, high}
9  @attribute student {yes, no}
10 @attribute credit { fair, excellent}
11 @attribute purchase { no, yes}
12
13 @data
14 1,young,high,no,fair,no
15 2,young,high,no,excellent,no
16 3,middle,high,no,fair,yes
17 4,senior,medium,no,fair,yes
18 5,senior,low,yes,fair,yes
19 6,senior,low,yes,excellent,no
20 7,middle,low,yes,excellent,yes
21 8,young,medium,no,fair,no
22 9,young,low,yes,fair,yes
23 10,senior,medium,yes,fair,yes
24 11,young,medium,yes,excellent,yes
25 12,middle,medium,no,excellent,yes
26 13,middle,high,yes,fair,yes
27 14,senior,medium,no,excellent,no

```

Example 1. Will Customer Buy Computer?

Java Program Source Code

See: java-code-ml-weka2018/src/ence688p/ClassificationTask.java

Abbreviated Program Output (J48 unpruned tree)

```
age = young
| student = yes: yes (2.0)
| student = no: no (3.0)
age = middle: yes (4.0)
age = senior
| credit = fair: yes (3.0)
| credit = excellent: no (2.0)
```

Number of Leaves : 5

Size of the tree : 8

Example 1. Will Customer Buy Computer?

Classification Accuracy wrt Training Dataset

Correctly Classified Instances	14	100 %
Incorrectly Classified Instances	0	0 %
Kappa statistic	1	
Mean absolute error	0	
Root mean squared error	0	
Relative absolute error	0 %	
Root relative squared error	0 %	
Total Number of Instances	14	

=== Confusion Matrix ===

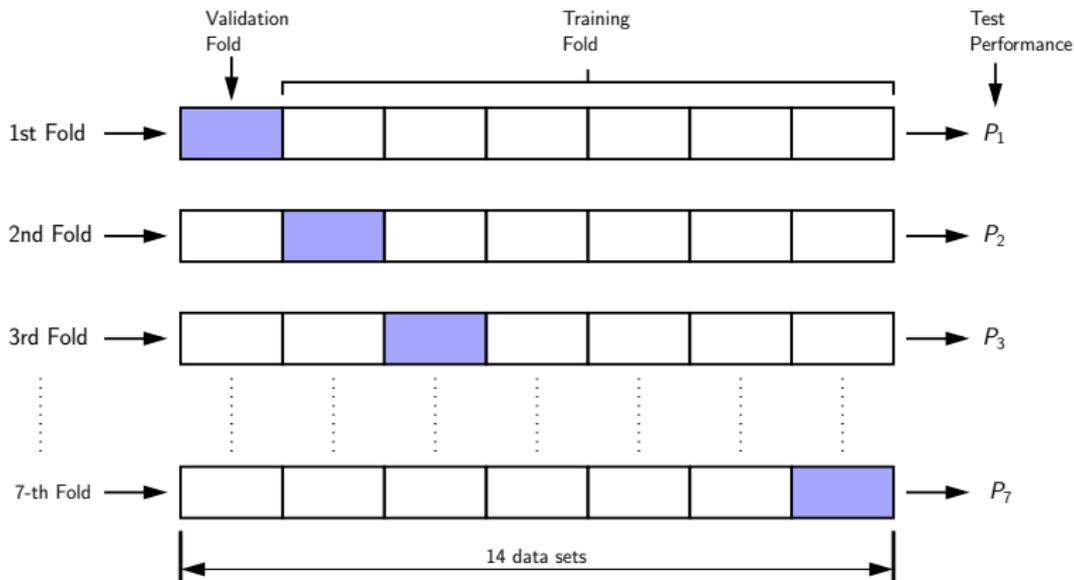
```

a b  <-- classified as
5 0 | a = no
0 9 | b = yes

```


Example 1. Will Customer Buy Computer?

Cross Validation Model (nofolds = 7)



Example 1. Will Customer Buy Computer?

Cross Validation Model (after classification) (nofolds = 7)

Correctly Classified Instances	10	71.4286 %
Incorrectly Classified Instances	4	28.5714 %
Kappa statistic	0.3778	
Mean absolute error	0.2798	
Root mean squared error	0.4393	
Relative absolute error	58.3333 %	
Root relative squared error	88.6322 %	
Total Number of Instances	14	

=== Confusion Matrix ===

```

a b  <-- classified as
3 2 | a = no
2 7 | b = yes

```

Example 1. Will Customer Buy Computer?

