Lect2: Classification #1 Basic methods Tokyo Tech. Intro. to Comp. & Data Lecture week2

Basic and well-known "methods" for classification

- 1. Introduction to classification.
- 2. Naive Bayes.
- 3. Decision Tree.
- 4. On Exercise #2. (Intro. to classification, again)
- * Some of the slide materials (in particular, green ones) are from the slides of the authors of the textbook and their group at the University of Waikato.
- * In this course, terms that I use for the class (which may not be standard sometimes) are quoted, e.g., "method".

1. Introduction to classification

Classification is a task to obtain a model for predicting the <u>class value</u> of given instances.

glossaries

Usually, 1 (yes, positive) or 0 (no, negative), but sometimes {0, 1, 2, 3}, etc.

yes

a[19] != 'k'

no

"method" = model + ML algo.

class or class value

= the target attribute of classification.

model

= a rule for predicting the class of each instance.

E.g., decision list

machine learning algorithm

= a way to compute a model from a given example set. The mushroom exercise is a typical example of classifications

11111111122 123456789012345678901 pfsnfsfcnbtsspwpwoewvd pksefffcnbtkkwppwoewvp e
ffwfnfwbntffwwpwoeksq <mark>e</mark>xswtlfcbnesswwpwopnng <mark>e</mark>xyufnfcnpesfwwpwofhvd exygtnfcbntsspgpwopkyd <mark>e</mark>ffetnfcbutsspwpwopnvd pfyyfnfwnyeyyyypyoewcl ebsytlfcbnesswwpwopksm p x s w f c f w n g e s s w w p w o p k v d esfqfnfcnkesswwpwopkyu pxsnfyfcnbtkkpppwoewvd exfgtnfcbptsswgpwopkvd UC Irvine ML Repository, 1987

/2

1. Introduction to classification

glossaries (cont.)

```
instance = an example data.
```

data set (or, sample) = a set of instances.

- attribute = a feature of each instance.
- nominal (or, categorical) attribute
 - = an attribute whose value is from a finite set.

numeric atteribute

= an attribute

whose value	Outlook	Temperature	Humidity	Windy	Play
	Sunny	85	85	False	No
is a (real) num	Der _{Sunny}	80	90	True	No
	Overcast	83	86	False	Yes
	Rainy	75	80	False	Yes

A data set of examples of the relation between whether and "play"

2. Naive Base

2. Naive Bayes

Naive Bayes is a ML algorithm to obtain a probabilistic classification model.

A probabilistic classification model

= probabilities of class values for a given instance.

Pr[P = Yes | (Weather, T, H, W) = (sunny, 80, 90, windy)]Pr[P = No | (Weather, T, H, W) = (sunny, 80, 90, windy)]

7	Outlook	Temperature	Humidity	Windy	Play ←	class value
	Sunny	85	85	False	No	
	Sunny	80	90	True	No	
	Overcast	83	86	False	Yes	
	Rainy	75	80	False	Yes	

A data set of examples of the relation between whether and "play"

2. Naive Base

2.1 Basics on probability

in general,

2.1. Basics on probability (from Probability Theory)

For each *probabilistic event* E, we write Pr[E] to denote the *probability* that E occurs. More concretely, we assume random variables and probability distributions on these random variables and specify E in terms of the event that such random variables take some values, e.g., $E \Leftrightarrow X = 5$, $E \Leftrightarrow Y \in [0.5, 1.5].$

The *joint probabilistic* $Pr[A \land B]$ is defined by

 $\Pr[A \land B] = \Pr[A|B] \times \Pr[B]$

where $\Pr[A | B]$ is conditional probability.

could be < or >We say that A and B are independent if $\Pr[A \land B] = \Pr[A] \cdot \Pr[B]$ or equivalently, $\Pr[A \mid B] \stackrel{*}{=} \Pr[A]$ because

2. Naive Base

2.1 Basics on probability

Bayes's rule is a way to compute conditional probability Pr[H/E] as follows.

$$\Pr[H|E] = \frac{\Pr[E|H] \times \Pr[H]}{\Pr[E]}$$

because...

$$\Pr[H|E] = \frac{\Pr[H \land E]}{\Pr[E]} = \frac{\Pr[E|H] \times \Pr[H]}{\Pr[E]}$$



Thomas Bayes 1702—1761 in England

Pr[H]: a priori probability i.e., the probability of *H* before evidence *E* is seen.

 $\Pr[H | E]$: a posteriori probability i.e., the probability of *H* after evidence *E* is seen.

2. Naive Base 2.2 Prob. assumptions

2.2. Probabilistic assumptions for Data Mining

We assume a certain *probability distribution* \mathcal{D} on instances (i.e., tuples of attribute values) of an assumed "domain", and our data set (or, sample) is a good approximation of \mathcal{D} .



2. Naive Base

2.2 Prob. assumptions



What a naive assumption !

 $\Pr_{\mathcal{D}}[Weather = \text{sunny \& } W = \text{windy }] \approx 0.2 \times 0.1 = 0.02$

Naive Bayes is a method to obtain a probabilistic classification model.

For example, suppose we want to compute

 $\Pr[P = no | (Weather, T, H, W) = (sunny, 80, 90, windy)]$

= $\Pr[P = no | Weather = sunny \& T = 80 \& H = 90 \& W = windy]$

Note that we can compute, e.g.,

 $Pr_{\mathcal{D}}[W = windy | P = no] \Rightarrow Pr_{\widehat{\mathcal{D}}}[W = windy | P = no]$ $= \frac{15}{120} = 0.125.$ Thus, our task is to compute
of no = 120
of windy & no = 15

 $\Pr[P = no | A \& B \& C \& D]$ based on $\Pr[P = no | A]$, etc.



Thus, we can compute

Pr[P = no | (Weather, T, H, W) = (sunny, 80, 90, windy)]

= $\Pr[P = no | Weather = sunny \& T = 80 \& H = 90 \& W = windy]$

 $\frac{\Pr[Wth=\operatorname{sny}|P=n] \cdot \Pr[T=80|P=n] \cdot \Pr[H=90|P=n] \cdot \Pr[W=wind|P=n] \cdot \Pr[P=n]}{\Pr[Wth=\operatorname{sny} \& T=80 \& H=90 \& W=wind]}$

Also, we have Pr[P = yes | (Weather, T, H, W) = (sunny, 80, 90, windy)]

 $\frac{\Pr[Wth=\operatorname{sny}|P=y] \cdot \Pr[T=80|P=y] \cdot \Pr[H=90|P=y] \cdot \Pr[W=\operatorname{wind}|P=y] \cdot \Pr[P=y]}{\Pr[Wth=\operatorname{sny} \& T=80 \& H=90 \& W=\operatorname{wind}]}$

Note that we do not have to compute

Pr[*Wth*=sny & *T*= 80 & *H*=90 & *W*=wind]

for predicting the class value.

3. Decision tree

Decision tree is a classification model by using a "decision tree" like this..



- "Divide-and-conquer" approach produces tree
- Nodes involve testing a particular attribute
- Usually, attribute value is compared to constant
- Other possibilities:
 - Comparing values of two attributes
 - Using a function of one or more attributes
- Leaves assign classification, set of classifications, or probability distribution to instances
- Unknown instance
 - is routed down the tree



Decision tree
 Some subclasses of decision trees

Decision stump = one level decision tree

(also called 1R in Weka)

- ✤ 1R was described in a paper by Holte (1993)
- Simplicity first pays off!

Very Simple Classification Rules Perform Well on Most Commonly Used Datasets

Robert C. Holte, Computer Science Department, University of Ottawa









3. Decision tree

Decision list

3.1 Variations

3. Decision tree

3. Decision tree

3.2 Construction

3.2. How to construct decision trees (only intro.)

Here we discuss the very simple one for constructing a decision stump or 1R. (We will discuss an important ML algorithm (which is a basis of J4.8 in Weka) later in this course.)

For each attribute,
For each value of the attribute, make a rule as follows:
 count how often each class appears
 find the most frequent class
 make the rule assign that class to this attribute-value
 Calculate the error rate of the rules
Choose the rules with the smallest error rate

Example:

3. Decision tree

3.2 Construction

Evaluating the weather attributes

Outlook	Temp	Humidity	Windy	Play
Sunny	Hot	High	False	No
Sunny	Hot	High	True	No
Overcast	Hot	High	False	Yes
Rainy	Mild	High	False	Yes
Rainy	Cool	Normal	False	Yes
Rainy	Cool	Normal	True	No
Overcast	Cool	Normal	True	Yes
Sunny	Mild	High	False	No
Sunny	Cool	Normal	False	Yes
Rainy	Mild	Normal	False	Yes
Sunny	Mild	Normal	True	Yes
Overcast	Mild	High	True	Yes
Overcast	Hot	Normal	False	Yes
Rainy	Mild	High	True	No

Attribute	Rules	Errors	Total errors
Outlook	Sunny \rightarrow No	2/5	4/14
	$Overcast \to Yes$	0/4	
	Rainy \rightarrow Yes	2/5	
Temp	$Hot\toNo^*$	2/4	5/14
	$Mild \to Yes$	2/6	
	$Cool \to Yes$	1/4	
Humidity	$High \to No$	3/7	4/14
	Normal \rightarrow Yes	1/7	
Windy	$False \to Yes$	2/8	5/14
	True \rightarrow No*	3/6	

* indicates a tie

4. On Exercise #2

Classification rule discovery project:

- Introduction to Weka^{*1}.
- Use 1R, Naive Bayes, and Decision Tree methods.

Task #1: Try to use Weka

(a) Prepare m8126.arff.

(b) Compute classification models on m8126.arff.

Task #2: Understand basics on classification tasks and decision trees

(c) Use breast-cancer.arff (given as a sample data in Weka) to learn how to create an appropriate decision tree.

Task #2: Understand basics on classification tasks and decision trees (c) Use breast-cancer.arff (given as a sample data in Weka) to learn how to create an appropriate decision tree. A standard flow of classification tasks Learning phase training set data set 70% ML algo. classification 30% Most difficult part test set of data mining !! our starting point Test phase 18

Intro. to classification task (cont.)

4. On Exercise #2

4. On Exercise #2

Intro. to classification



Important Points:

1. The performance of the obtained model should be tested on a separate data, i.e., test set.

e.g., a decision tree with 300 leaves for 300 instances

4. On Exercise #2

Intro. to classification

 Avoid overfitting. It is rather easy to obtain a model that is good (or almost perfect) on a training set, but it is not useful in general.

The simpler is the better.

Occam's Razor — Plurality must never be posited without necessity

