

Algorithms for Classification:

The Basic Methods

Outline

- Not really a classifier: **0R** or **ZeroR**
- Simplicity first: **1R** or **OneR**
- **Naïve Bayes**

Classification

- Task: Given a set of pre-classified examples, build a model or *classifier* to classify new cases.
- *Supervised* learning: classes are known for the examples used to build the classifier.
- A classifier can be a set of rules, a decision tree, a neural network, etc.
- Typical applications: credit approval, direct marketing, fraud detection, medical diagnosis, ...



Simplicity first

- Simple algorithms often work very well!
- There are many kinds of simple structure, eg:
 - Majority class classifier
 - One attribute does all the work
 - All attributes contribute equally & independently
 - A weighted linear combination might do
 - Instance-based: use a few prototypes
 - Use simple logical rules
- Success of a method depends on the domain

Inferring rudimentary rules

- 0R: predicts majority class
-
- 1R: learns a 1-level decision tree
 - I.e., rules that all test one particular attribute
 - Basic version
 - One branch for each value
 - Each branch assigns most frequent class
 - Error rate: proportion of instances that don't belong to the majority class of their corresponding branch
 - Choose attribute with lowest error rate

(assumes nominal attributes)

Pseudo-code for 1R

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

- Note: “missing” is treated as a separate attribute value

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 → No	2/5	4/14*
	Overcast → Yes	0/4	
	Rainy → Yes	2/5	
Temp	Hot → No*	2/4	5/14
	Mild → Yes	2/6	
	Cool → Yes	1/4	
Humidity	High → No	3/7	4/14*
	Normal → Yes	1/7	
Windy	False → Yes	2/8	5/14
	True → No*	3/6	


* indicates a tie

Dealing with numeric attributes

- Discretize numeric attributes
- Divide each attribute's range into intervals
 - Sort instances according to attribute's values
 - Place breakpoints where the class changes (the majority class)
 - This minimizes the total error
- Example: *temperature* from weather data

Outlook	Temperature	Humidity	Windy	Play
Sunny	85	85	False	No
Sunny	80	90	True	No
Overcast	83	86	False	Yes
Rainy	75	80	False	Yes
...

64 65 68 69 70 71 72 72 75 75 80 81 83 85





Yes | No | Yes Yes Yes | No No  Yes Yes Yes | No | Yes Yes | No

The problem of overfitting

- This procedure is very sensitive to noise
 - One instance with an incorrect class label will probably produce a separate interval
- Also: *time stamp* attribute will have zero errors
- Simple solution:
enforce minimum number of instances in majority class per interval

Discretization example

- Example (with min = 3):

64	65	68	69	70	71	72	72	75	75	80	81	83	85						
Yes		No		Yes	Yes	Yes		No	No	Yes	Yes	Yes		No		Yes	Yes		No

- Final result for temperature attribute

64	65	68	69	70	71	72	72	75	75	80	81	83	85		
Yes	No	Yes	Yes	Yes		No	No	Yes	Yes	Yes		No	Yes	Yes	No

With overfitting avoidance

- Resulting rule set:

Attribute	Rules	Errors	Total errors
Outlook	Sunny \rightarrow No	2/5	4/14
	Overcast \rightarrow Yes	0/4	
	Rainy \rightarrow Yes	2/5	
Temperature	$\leq 70.5 \rightarrow$ Yes	1/5	5/14
	> 70.5 and $\leq 77.5 \rightarrow$ Yes	2/5	
	$> 77.5 \rightarrow$ No*	2/4	
Humidity	$\leq 82.5 \rightarrow$ Yes	1/7	3/14
	> 82.5 and $\leq 95.5 \rightarrow$ No	2/6	
	$> 95.5 \rightarrow$ Yes	0/1	
Windy	False \rightarrow Yes	2/8	5/14
	True \rightarrow No*	3/6	

Discussion of 1R

- 1R was described in a paper by Holte (1993)
 - Contains an experimental evaluation on 16 datasets (using *cross-validation* so that results were representative of performance on future data)
 - Minimum number of instances was set to 6 after some experimentation
 - 1R's simple rules performed not much worse than much more complex decision trees
- Simplicity first pays off!

Very Simple Classification Rules Perform Well on Most Commonly Used Datasets

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





Bayesian (Statistical) modeling

- “Opposite” of 1R: use all the attributes
- Two assumptions: Attributes are
 - *equally important*
 - *statistically independent* (given the class value)
 - I.e., knowing the value of one attribute says nothing about the value of another (if the class is known)
- Independence assumption is almost never correct!
- But ... this scheme works well in practice

Probabilities for weather data

Outlook			Temperature			Humidity			Windy			Play	
<i>Yes No</i>			<i>Yes No</i>			<i>Yes No</i>			<i>Yes No</i>			<i>Yes</i>	<i>No</i>
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

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

Probabilities for weather data

Outlook			Temperature			Humidity			Windy			Play	
Yes No			Yes No			Yes No			Yes No			Yes	No
Sunny	2	3	Hot	2	2	High	3	4	False	6	2	9	5
Overcast	4	0	Mild	4	2	Normal	6	1	True	3	3		
Rainy	3	2	Cool	3	1								
Sunny	2/9	3/5	Hot	2/9	2/5	High	3/9	4/5	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	Mild	4/9	2/5	Normal	6/9	1/5	True	3/9	3/5		
Rainy	3/9	2/5	Cool	3/9	1/5								

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

■ A new day:

Likelihood of the two classes

$$\text{For "yes"} = 2/9 \times 3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0053$$

$$\text{For "no"} = 3/5 \times 1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0206$$

Conversion into a probability by normalization:

$$P(\text{"yes"}) = 0.0053 / (0.0053 + 0.0206) = 0.205$$

$$P(\text{"no"}) = 0.0206 / (0.0053 + 0.0206) = 0.795$$

Bayes's rule

- Probability of event H given evidence E :

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

- *A priori* probability of H : $\Pr[H]$
 - Probability of event *before* evidence is seen
- *A posteriori* probability of H : $\Pr[H | E]$
 - Probability of event *after* evidence is seen

from Bayes “Essay towards solving a problem in the doctrine of chances” (1763)

Thomas Bayes

Born: 1702 in London, England

Died: 1761 in Tunbridge Wells, Kent, England



Naïve Bayes for classification

- Classification learning: what's the probability of the class given an instance?
 - Evidence E = instance
 - Event H = class value for instance
- Naïve assumption: evidence splits into parts (i.e. attributes) that are *independent*

$$\Pr[H \mid E] = \frac{\Pr[E_1 \mid H] \Pr[E_2 \mid H] \dots \Pr[E_n \mid H] \Pr[H]}{\Pr[E]}$$

Weather data example

Outlook	Temp.	Humidity	Windy	Play
Sunny	Cool	High	True	?

← ***Evidence E***

***Probability of
class “yes”***

$$\begin{aligned}\Pr[\text{yes} \mid E] &= \Pr[\text{Outlook} = \text{Sunny} \mid \text{yes}] \\ &\quad \times \Pr[\text{Temperature} = \text{Cool} \mid \text{yes}] \\ &\quad \times \Pr[\text{Humidity} = \text{High} \mid \text{yes}] \\ &\quad \times \Pr[\text{Windy} = \text{True} \mid \text{yes}] \\ &\quad \times \frac{\Pr[\text{yes}]}{\Pr[E]} \\ &= \frac{\frac{2}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{3}{9} \times \frac{9}{14}}{\Pr[E]}\end{aligned}$$

The “zero-frequency problem”

- What if an attribute value doesn't occur with every class value?
(e.g. “Humidity = high” for class “yes”)
 - Probability will be zero! $\Pr[Humidity = High \mid yes] = 0$
 - *A posteriori* probability will also be zero! $\Pr[yes \mid E] = 0$
(No matter how likely the other values are!)
- Remedy: add 1 to the count for every attribute value-class combination (*Laplace estimator*)
- Result: probabilities will never be zero!
(also: stabilizes probability estimates)

*Modified probability estimates

- In some cases adding a constant different from 1 might be more appropriate
- Example: attribute *outlook* for class *yes*

$$\frac{2 + \mu/3}{9 + \mu}$$

Sunny

$$\frac{4 + \mu/3}{9 + \mu}$$

Overcast

$$\frac{3 + \mu/3}{9 + \mu}$$

Rainy

- Weights don't need to be equal (but they must sum to 1)

$$\frac{2 + \mu p_1}{9 + \mu}$$

$$\frac{4 + \mu p_2}{9 + \mu}$$

$$\frac{3 + \mu p_3}{9 + \mu}$$

Missing values

- Training: instance is not included in frequency count for attribute value-class combination
- Classification: attribute will be omitted from calculation
- Example:

Outlook	Temp.	Humidity	Windy	Play
?	Cool	High	True	?

Likelihood of "yes" = $3/9 \times 3/9 \times 3/9 \times 9/14 = 0.0238$

Likelihood of "no" = $1/5 \times 4/5 \times 3/5 \times 5/14 = 0.0343$

$P(\text{"yes"}) = 0.0238 / (0.0238 + 0.0343) = 41\%$

$P(\text{"no"}) = 0.0343 / (0.0238 + 0.0343) = 59\%$

Numeric attributes

- Usual assumption: attributes have a *normal* or *Gaussian* probability distribution (given the class)
- The *probability density function* for the normal distribution is defined by two parameters:

- *Sample mean* μ

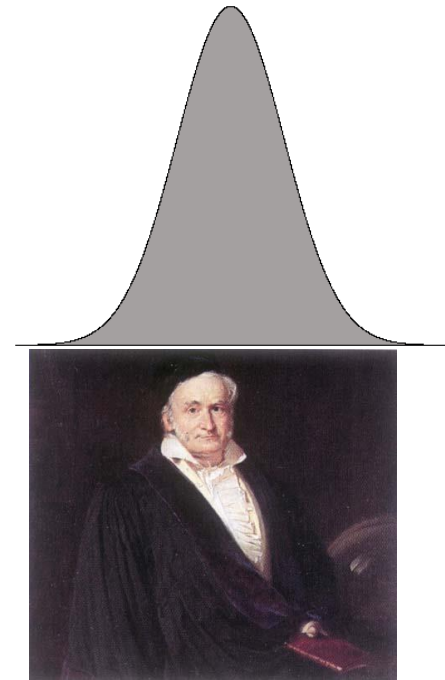
$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

- *Standard deviation* σ

$$\sigma = \frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2$$

- Then the density function $f(x)$ is

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$



Karl Gauss, 1777-1855
great German mathematician

Statistics for weather data

Outlook			Temperature		Humidity		Windy			Play	
			Yes	No	Yes	No	Yes	No		Yes	No
Sunny	2	3	64, 68,	65, 71,	65, 70,	70, 85,	False	6	2	9	5
Overcast	4	0	69, 70,	72, 80,	70, 75,	90, 91,	True	3	3		
Rainy	3	2	72, ...	85, ...	80, ...	95, ...					
Sunny	2/9	3/5	$\mu = 73$	$\mu = 75$	$\mu = 79$	$\mu = 86$	False	6/9	2/5	9/14	5/14
Overcast	4/9	0/5	$\sigma = 6.2$	$\sigma = 7.9$	$\sigma = 10.2$	$\sigma = 9.7$	True	3/9	3/5		
Rainy	3/9	2/5									

- Example density value:

$$f(\text{temperature} = 66 \mid \text{yes}) = \frac{1}{\sqrt{2\pi} 6.2} e^{-\frac{(66-73)^2}{2*6.2^2}} = 0.0340$$

Classifying a new day

- A new day:

Outlook	Temp.	Humidity	Windy	Play
Sunny	66	90	true	?

Likelihood of "yes" = $2/9 \times 0.0340 \times 0.0221 \times 3/9 \times 9/14 = 0.000036$

Likelihood of "no" = $3/5 \times 0.0291 \times 0.0380 \times 3/5 \times 5/14 = 0.000136$

$P(\text{"yes"}) = 0.000036 / (0.000036 + 0.000136) = 20.9\%$

$P(\text{"no"}) = 0.000136 / (0.000036 + 0.000136) = 79.1\%$

- Missing values during training are not included in calculation of mean and standard deviation

Naïve Bayes: discussion

- Naïve Bayes works surprisingly well (even if independence assumption is clearly violated)
- Why? Because classification doesn't require accurate probability estimates *as long as maximum probability is assigned to correct class*
- However: adding too many redundant attributes will cause problems (e.g. identical attributes)
- Note also: many numeric attributes are not normally distributed (\rightarrow *kernel density estimators*)

Naïve Bayes Extensions

- Improvements:
 - select best attributes (e.g. with greedy search)
 - often works as well or better with just a fraction of all attributes
- Bayesian Networks

Summary

- **ZeroR** – not really a classifier, predicts majority class
- **OneR** – uses rules based on just one attribute
- **Naïve Bayes** – use all attributes and Bayes rules to estimate probability of the class given an instance.
- Simple methods frequently work well, but ...
 - Complex methods can be better (as we will see)