

Lab 6

Bayes classification

Useful information

Alex Becheru

irlab.becheru.net

irlab@becheru.net

Laboratory scope

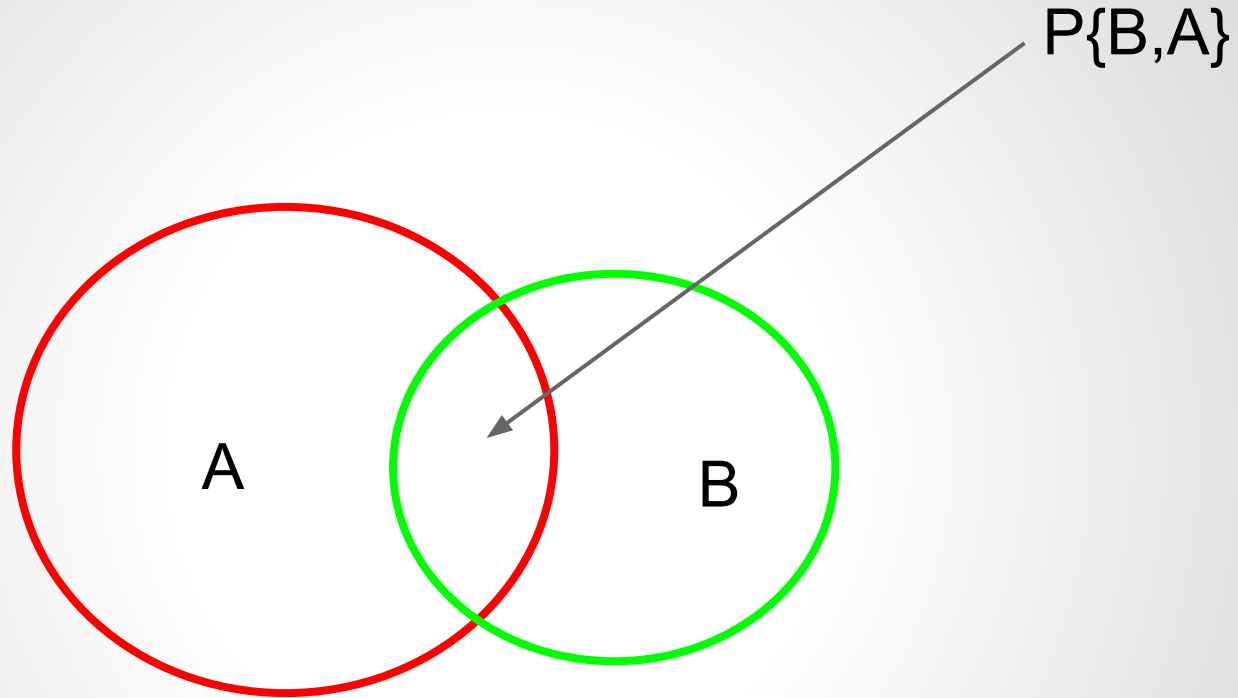
Naive Bayes Classification

Bayesian probabilistics

What is the probability of an event A to take place if it is conditioned by event B ?

$$P\{B|A\} = \frac{P\{A|B\} \cdot P\{B\}}{P\{A\}}.$$

$$\textit{posterior} = \frac{\textit{likelihood} \times \textit{prior probability}}{\textit{evidence}}.$$

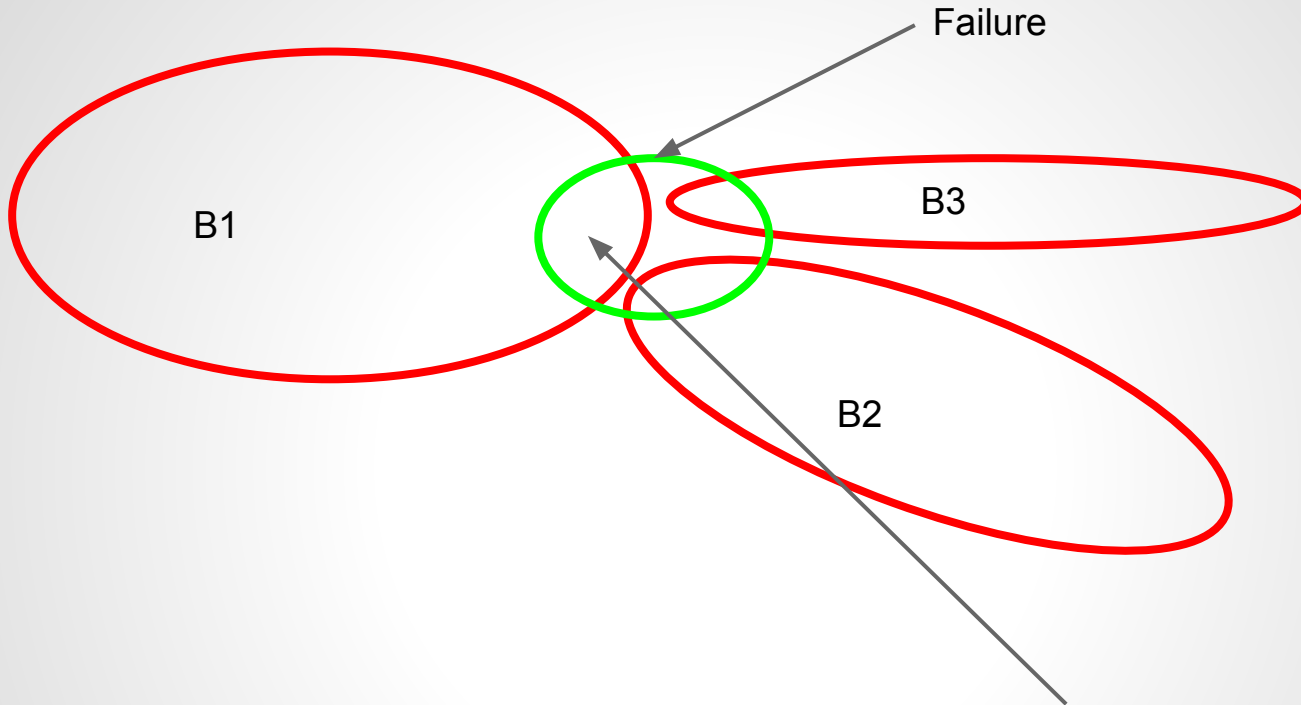


Simplification of formula

- There are 2 neighbors A and B that share a common land.
- We know that the common land is $\frac{1}{3}$ of B's territory: $P\{A|B\}$
- We know both A's and B's surface: $P\{A\}$, $P\{B\}$
- How much those the common land represent of A's territory (percent)? : $P\{B|A\}$?

Example 1

- A factory produces the same items using machines B1, B2 and B3.
- B1 produces 60% of items with 6% failure
- B2 produces 30% of items with 3% failure
- B3 produces 10% of items with 5% failure
- If i pick a failed product randomly, what is the probability to be produced from B1,B2 or B3?



Probability that an element made by B1 to be failed? = $P\{B1|F\}$
Probability that an element made by B2 to be failed? = $P\{B2|F\}$
Probability that an element made by B3 to be failed? = $P\{B3|F\}$

- $P\{B1|F\} = P\{F|B1\} * P\{B1\} / P\{F\}$ = number of items broken made with B1 / total number of broken items
- $P\{F|B1\} = 6\%$ of items made by B1
- $P\{B1\} = 60\%$ of total items made by all machines
- $P\{F|B1\} * P\{B1\} = 0.06 * 0.6 = 0.036\%$ of all elements made are broken and made with B1
- $P\{F\} = P\{F|B1\} + P\{F|B2\} + P\{F|B3\} = 0.06 * 0.6 + 0.03 * 0.3 + 0.05 * 0.1 = 0.05\%$ of all elements are broken
- $P\{B1|F\} = 0.036 / 0.05 = 72\%$ chance to pick an element made by B1 if we chose randomly from the broken elements

Ex 1.

What about $P\{B_2|F\}$ and $P\{B_3|F\}$?

Naive bayes classification

It is called naive or idiot because it supposes that there is no dependency between attributes.

Is that ok?

Usually attributes depend in some way on the other attributes!

The horsepower and a car's speed depend on each other, but naive bayes algorithm does not consider that.

The greatest advantage of Naive Bayes algorithm is that it needs a small training set!

Classification

Classification of instance I having attributes $Ia_1 \dots Ia_n$.

Classes: ClassA, ClassB

$$P\{\text{ClassA}|I\} = P\{Ia_1|\text{ClassA}\} * P\{Ia_2|\text{ClassA}\} * \dots * P\{Ia_n|\text{ClassA}\} * P\{\text{ClassA}\}$$

$$P\{\text{ClassB}|I\} = P\{Ia_1|\text{ClassB}\} * P\{Ia_2|\text{ClassB}\} * \dots * P\{Ia_n|\text{ClassB}\} * P\{\text{ClassB}\}$$

If $P\{\text{ClassA}|I\} > P\{\text{ClassB}|I\}$ then : I is classified as ClassA

else: I is classified as ClassB

Example 2

Salary	Has car?	Marital status	Happy
100k	Yes	No	Yes
100k	Yes	Yes	No
120k	Yes	Yes	Yes
60k	No	No	Yes
60k	No	Yes	No

A guy with a salary of 100k with a wife and no car is he happy?

$$P\{100k, \text{happy}=\text{yes}\} = \frac{1}{2}$$

$$P\{100k, \text{happy}=\text{no}\} = \frac{1}{2}$$

$$P\{\text{car}=\text{no}, \text{happy}=\text{yes}\} = \frac{1}{2}$$

$$P\{\text{car}=\text{no}, \text{happy}=\text{no}\} = \frac{1}{2}$$

$$P\{\text{wife}=\text{yes}, \text{happy}=\text{yes}\} = \frac{1}{3}$$

$$P\{\text{wife}=\text{yes}, \text{happy}=\text{no}\} = \frac{2}{3}$$

$$P\{\text{happy}=\text{yes}\} = \frac{3}{5}$$

$$P\{\text{happy}=\text{no}\} = \frac{2}{5}$$

$$P\{100k, \text{car}=\text{no}, \text{wife}=\text{yes} | \text{happy}=\text{yes}\} = \frac{1}{2} * \frac{1}{2} * \frac{1}{3} * \frac{3}{5} = \frac{3}{60}$$

$$P\{100k, \text{car}=\text{no}, \text{wife}=\text{yes} | \text{happy}=\text{no}\} = \frac{1}{2} * \frac{1}{2} * \frac{2}{3} * \frac{2}{5} = \frac{4}{60}$$

No he is not happy

Ex. 2

Salary	Has car?	Marital status	Happy
100k	Yes	No	Yes
100k	Yes	Yes	No
120k	Yes	Yes	No
60k	No	No	Yes

A guy with a salary of 100k with a wife and no car is he happy?

Problem

With the current method of calculation when

$P\{\text{married}=\text{yes}, \text{happy}=\text{yes}\}=0$ which makes

$P\{100\text{k}, \text{married}=\text{yes}, \text{car}=\text{no} | \text{happy}=\text{no}\}=0$

The occurrence of 0 in calculations might mess up the classification!

The solution was to calculate probabilities from frequencies!

Ex. 3

- open iris.arff
- classify using NaiveBayes classifier
- see the result

Ex. 4

- open iris.arff
- classify using the BayesNet classifier
- where does the difference in precision of classification come?

Ex. 5

- compare NaiveBayes, J48 and ZeroR on:
 - diabetes.arff
 - supermarket.arff
- which is the best algorithm?