

A spam classifier based on Bayes network

Alberto Franzin, Fabio Palese

Sistemi Intelligenti

January 15th, 2013

INTRODUCTION

Introduction

Bayesian networks

Definition

Naive Bayes

Naive Bayes for spam classification

SpamBayes

Implementation

Tests

Results and Conclusions

THE BAYESIAN APPROACH

- Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

defines the *a posteriori* probability of event A , knowing the event B has already occurred.

THE BAYESIAN APPROACH

- Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

defines the *a posteriori* probability of event A , knowing the event B has already occurred.

- In other words, we can estimate the probability of an hypothesis, given that we know the consequences.

THE BAYESIAN APPROACH

- Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

defines the *a posteriori* probability of event A , knowing the event B has already occurred.

- In other words, we can estimate the probability of an hypothesis, given that we know the consequences.
- This has led to two different interpretations of the theorem.

THE BAYESIAN APPROACH

- Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

THE BAYESIAN APPROACH

- Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

- $P(A|B)$ is the *a posteriori* probability

THE BAYESIAN APPROACH

- Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

- $P(A|B)$ is the *a posteriori* probability
- $P(B|A)$ is the *likelihood*

THE BAYESIAN APPROACH

- ▶ Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

- ▶ $P(A|B)$ is the *a posteriori* probability
- ▶ $P(B|A)$ is the *likelihood*
- ▶ $P(B|A)P(A)$ is the *prior* probability

THE BAYESIAN APPROACH

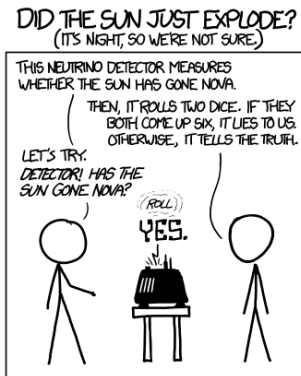
- Bayes rule:

$$P(A|B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B|A)P(A)}{P(B)}$$

- $P(A|B)$ is the *a posteriori* probability
- $P(B|A)$ is the *likelihood*
- $P(B|A)P(A)$ is the *prior* probability
- $P(B) = \sum_{a \in A} P(B|A = a)P(A = a)$ is the *total* probability

THE BAYESIAN APPROACH

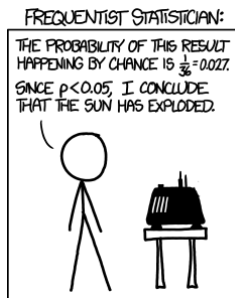
Frequentists vs. Bayesians



from <http://xkcd.com/1132>, see also
http://en.wikipedia.org/wiki/Sunrise_problem

THE BAYESIAN APPROACH

Frequentists vs. Bayesians



from <http://xkcd.com/1132>, see also
http://en.wikipedia.org/wiki/Sunrise_problem

The frequentist relies on the theoretical probability of the events.

THE BAYESIAN APPROACH

Frequentists vs. Bayesians



from <http://xkcd.com/1132>, see also
http://en.wikipedia.org/wiki/Sunrise_problem

The bayesian observes the past events occurred,
and adapts the probability accordingly.

WHAT IT IS

A Bayes network is a way to describe causal relationships between events.

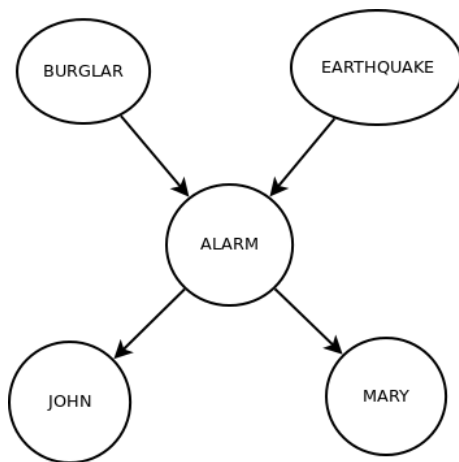
- ▶ Nodes = events
- ▶ (Directed) Edges = causal relationship

WHAT IT IS

A Bayes network is a way to describe causal relationships between events.

- ▶ Nodes = events
- ▶ (Directed) Edges = causal relationship
- ▶ Two nodes are connected by an edge: the child of an arc is influenced by its ancestor in a probabilistic way

AN EXAMPLE



CONDITIONAL INDEPENDENCE

- If

$$P(A|B, C) = P(A|B)$$

then we say that B and C are *conditionally independent*.

CONDITIONAL INDEPENDENCE

- If

$$P(A|B, C) = P(A|B)$$

then we say that B and C are *conditionally independent*.

- Note that *conditional independence* \neq *independence*

CONDITIONAL INDEPENDENCE

- ▶ If

$$P(A|B, C) = P(A|B)$$

then we say that B and C are *conditionally independent*.

- ▶ Note that *conditional independence* \neq *independence*
- ▶ Explaining away: if we know that one possible cause of the event has happened, this may *explain away* the event, being all the other causes less probable once we know the one that happened.

NAIVE BAYES

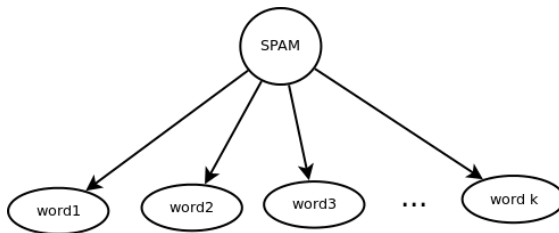
- ▶ Computing all the probabilities in a Bayesian network requires exponential time. We introduce the assumption of independence among variables.

NAIVE BAYES

- ▶ Computing all the probabilities in a Bayesian network requires exponential time. We introduce the assumption of independence among variables.
- ▶ It is called *naive*, since it's often unrealistic, but it yields good results.

NAIVE BAYES

- ▶ Computing all the probabilities in a Bayesian network requires exponential time. We introduce the assumption of independence among variables.
- ▶ It is called *naive*, since it's often unrealistic, but it yields good results.
- ▶ In spam classification:



NAIVE BAYES FOR SPAM CLASSIFICATION

Algorithm:

- ▶ Training

NAIVE BAYES FOR SPAM CLASSIFICATION

Algorithm:

- ▶ Training
- ▶ Validation

NAIVE BAYES FOR SPAM CLASSIFICATION

Algorithm:

- ▶ Training
- ▶ Validation
- ▶ Testing

NAIVE BAYES FOR SPAM CLASSIFICATION

Formulas:

- For each word:

$$P_{word|spam} = \frac{\text{\# occurrences of word in spam mails}}{\text{\# total occurrences of word}}$$

NAIVE BAYES FOR SPAM CLASSIFICATION

Formulas:

- For each word:

$$P_{word|spam} = \frac{\text{\# occurrences of word in spam mails}}{\text{\# total occurrences of word}}$$

- Final for spam is:

$$P_{spam} = \prod_{words \in mail} P_{spam|word}$$

NAIVE BAYES FOR SPAM CLASSIFICATION

Formulas:

- For each word:

$$P_{word|spam} = \frac{\# \text{ occurrences of word in spam mails}}{\# \text{ total occurrences of word}}$$

- Final for spam is:

$$P_{spam} = \prod_{words \in mail} P_{spam|word}$$

- same for ham

NAIVE BAYES FOR SPAM CLASSIFICATION

Formulas:

- For each word:

$$P_{word|spam} = \frac{\# \text{ occurrences of word in spam mails}}{\# \text{ total occurrences of word}}$$

- Final for spam is:

$$P_{spam} = \prod_{words \in mail} P_{spam|word}$$

- same for ham
- Outcome is the class that maximizes the probabilities.

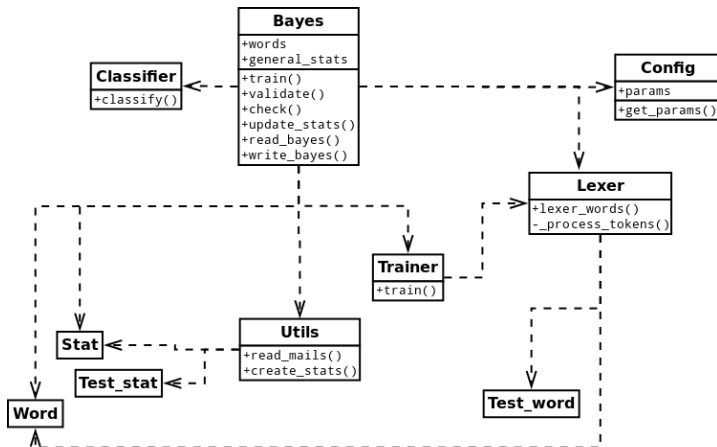
SPAMBAYES

- ▶ SpamBayes: applying NaiveBayes to spam classification
Python with Ply and BeautifulSoup
dataset: SpamAssassin archive
Code and documentation at <http://code.google.com/p/sist-int-2012project/>

SPAMBAYES

- ▶ SpamBayes: applying NaiveBayes to spam classification
Python with Ply and BeautifulSoup
dataset: SpamAssassin archive
Code and documentation at <http://code.google.com/p/sist-int-2012project/>
- ▶ spam: undesired mail
ham: valid mail

SPAMBAYES



NOTES ON IMPLEMENTATION

► Smoothing:

$$P_{word|spam} = \frac{\# \text{ occurrences of word in spam mails} + k}{\# \text{ total occurrences of the word} + |C| \times k}$$

NOTES ON IMPLEMENTATION

- Smoothing:

$$P_{word|spam} = \frac{\# \text{ occurrences of word in spam mails} + k}{\# \text{ total occurrences of the word} + |C| \times k}$$

- Calculations can be simplified: some words bring little contribution to the mail status

NOTES ON IMPLEMENTATION

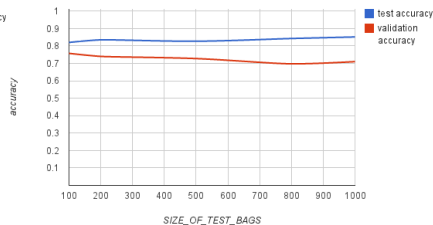
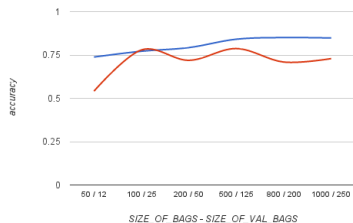
- Smoothing:

$$P_{word|spam} = \frac{\# \text{ occurrences of word in spam mails} + k}{\# \text{ total occurrences of the word} + |C| \times k}$$

- Calculations can be simplified: some words bring little contribution to the mail status
- Several parameters to tune: we describe the more relevant ones

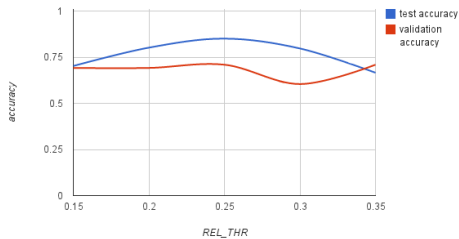
PARAMETERS

Size of training/validation/test sets:



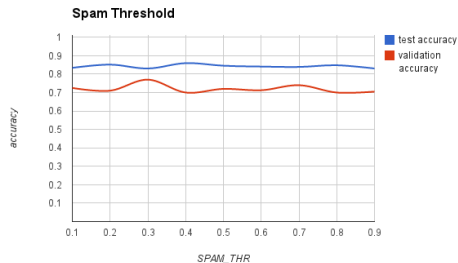
PARAMETERS

Relevance threshold:



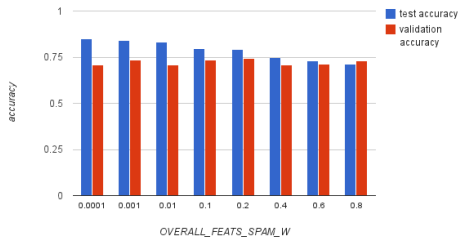
PARAMETERS

“Spamicity” threshold:



PARAMETERS

Feature statistic threshold:



RESULTS

- About the dataset:

RESULTS

- ▶ About the dataset:
 - ▶ General features

RESULTS

- ▶ About the dataset:
 - ▶ General features
 - ▶ Single words

RESULTS

- ▶ About the dataset:
 - ▶ General features
 - ▶ Single words
- ▶ About the classification:

RESULTS

- ▶ About the dataset:
 - ▶ General features
 - ▶ Single words
- ▶ About the classification:
 - ▶ Accuracy

RESULTS

- ▶ About the dataset:
 - ▶ General features
 - ▶ Single words
- ▶ About the classification:
 - ▶ Accuracy
 - ▶ How to improve?

FINAL

Any questions?