

How a robot can use basic probability theory to
classify things:

The basics of statistical classification theory

Georgia Tech

PROJECT: Text-classification: an application

We are going to classify text into topics. We assume a random model for generating documents.

- Given the topic the words are independent.
- We do not consider the order the words appear in sentence nor how often they appear. Only if they appear or not.

Text-classification: the word \times document matrix

We consider the document \times word matrix X which records which words appear in which document:

	dog	cat	car	home	run	be	nice
Document1	1	0	0	1	1	0	0
Document2	0	1	0	1	0	1	0
Document3	1	0	0	0	0	1	0
Document4	0	0	1	1	0	1	0
Document5	0	0	1	0	1	0	1
Document6	0	0	1	0	0	1	1

(12)

One topic for each document

we consider two topics:

$T0 = \text{animal life}, T1 = \text{home and car}$

<i>Topic</i>		dog	cat	car	home	run	be	nice
<i>T0</i>	Document 1	1	0	0	1	1	0	0
<i>T0</i>	Document 2	0	1	0	1	0	1	0
<i>T0</i>	Document 3	1	0	0	0	0	1	0
<i>T1</i>	Document 4	0	0	1	1	0	1	0
<i>T1</i>	Document 5	0	0	1	0	1	0	1
<i>T1</i>	Document 6	0	0	1	0	0	1	1

(13)

We assume independence of words given topic

First document :

document 1 = (*dog, run, home*)

has probability under topic 0:

$$P(\textit{dog, run, home} | \textit{topic 0}) = P(\textit{dog} | \textit{topic0}) \cdot P(\textit{run} | \textit{topic0}) \cdot P(\textit{home} | \textit{topic0})$$

and under topic 1:

$$P(\textit{dog, run, home} | \textit{topic 1}) = P(\textit{dog} | \textit{topic1}) \cdot P(\textit{run} | \textit{topic1}) \cdot P(\textit{home} | \textit{topic1})$$

Conditional probabilities of words given topic

Let the conditional probabilities under topic 0 be:

$$\vec{p} = (P(\text{dog}|T0), P(\text{cat}|T0), P(\text{car}|T0), P(\text{home}|T0), P(\text{run}|T0), P(\text{be}|T0), P(\text{nice}|T0))$$

Let the conditional probabilities under topic 1 be:

$$\vec{q} = (P(\text{dog}|T1), P(\text{cat}|T1), P(\text{car}|T1), P(\text{home}|T1), P(\text{run}|T1), P(\text{be}|T1), P(\text{nice}|T1))$$

Bayse classification

Given document (*dog, run, cat*) compare the probabillites and chose the one with higher probability as the topic. IF

$$\pi_0 \cdot P(dog|T0) \cdot P(run|T0) \cdot P(cat|T0) > \pi_1 \cdot P(dog|T1) \cdot P(run|T1) \cdot P(cat|T1)$$

we classify as topic 0. Here π_0 is probability of topic 0 and π_1 is probability of topic 1.

Example of estimation of condition probabilities

<i>Topic</i>		dog	cat	car	home	run	be	nice
<i>T0</i>	Document 1	1	0	0	1	1	0	0
<i>T0</i>	Document 2	0	1	0	1	0	1	0
<i>T0</i>	Document 3	1	0	0	0	0	1	0
<i>T1</i>	Document 4	0	0	1	1	0	1	0
<i>T1</i>	Document 5	0	0	1	0	1	0	1
<i>T1</i>	Document 6	0	0	1	0	0	1	1

(14)

Estimated conditional probabilities given T0:

$$\hat{p} = \left(\frac{2}{3}, \frac{1}{3}, 0, \frac{2}{3}, \frac{1}{3}, \frac{2}{3}, 0 \right)$$

Estimated conditional probabilities given T1:

$$\hat{q} = \left(0, 0, 1, \frac{1}{3}, \frac{1}{3}, \frac{2}{3}, \frac{2}{3} \right)$$

Summary text classification naive bayse

- Estimate the coditional probabilities \vec{p} and \vec{q} .
- Use the estimates $\hat{\vec{p}}$ and $\hat{\vec{q}}$ to estimate the probabilities of a new document which you which to classify.
- Chose the topic classficiation which has higher estimated probability.