

Machine Learning Crash Course

Expectation Maximization

Dr. Henrich Matzinger and Dr. Greg Mayer
Machine Learning and AI Crash Course
December 2019

School of Mathematics
Georgia Institute of Technology, Atlanta, GA

Text Classification

Recall the document classification problem. We have **training data** with known topics.

<i>Topic</i>		dog	cat	car	home	run	be	nice
<i>T0</i>	Document1	1	0	0	1	1	0	0
<i>T0</i>	Document2	0	1	0	1	0	1	0
<i>T0</i>	Document3	1	0	0	0	0	1	0
<i>T1</i>	Document4	0	0	1	1	0	1	0
<i>T1</i>	Document5	0	0	1	0	1	0	1
<i>T1</i>	Document6	0	0	1	0	0	1	1

We calculate the conditional probability estimates:

$$T0: \hat{p} = \left(\frac{2}{3}, \frac{1}{3}, 0, \frac{2}{3}, \frac{1}{3}, \frac{2}{3}, 0\right), \quad T1: \hat{q} = \left(0, 0, 1, \frac{1}{3}, \frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right)$$

We then use \hat{p} and \hat{q} to classify a new document.

Errors in Training Data Topics

- Assume that the topics are given to us but with 30% percent error.
- Can we still use the topics? How?

Example of partially messed up topic

Messed up Topic	True Topic			dog	cat	car	home	run	be	nice
T_0	T_0	Doc1		1	0	0	1	1	0	0
T_0	T_0	Doc2		0	1	0	1	0	1	0
T_0	T_0	Doc3		1	0	0	0	0	1	0
T_1	T_1	Doc4		0	0	1	1	0	1	0
T_1	T_1	Doc5		0	0	1	0	1	0	1
T_1	T_1	Doc6		0	0	1	0	0	1	1

Unknown estimated conditional probabilities:

$$T_0: \hat{p} = \left(\frac{2}{3}, \frac{1}{3}, 0, \frac{2}{3}, \frac{1}{3}, \frac{2}{3}, 0\right) \quad T_1: \hat{q} = \left(0, 0, 1, \frac{1}{3}, \frac{1}{3}, \frac{2}{3}, \frac{2}{3}\right)$$

What should we do?

First Step of Solution

Try an iterative approach.

Messed up Topic	True Topic			dog	cat	car	home	run	be	nice
$T0$	$T0$	Doc1		1	0	0	1	1	0	0
$T0$	$T0$	Doc2		0	1	0	1	0	1	0
$T0$	$T0$	Doc3		1	0	0	0	0	1	0
$T1$	$T1$	Doc4		0	0	1	1	0	1	0
$T1$	$T1$	Doc5		0	0	1	0	1	0	1
$T1$	$T1$	Doc6		0	0	1	0	0	1	1

Step 1: Use the messed up topic as if they were accurate.

$$T0: \hat{p}_I = \left(\frac{2}{3}, 0, \frac{1}{3}, \frac{1}{3}, \frac{2}{3}, \frac{1}{3}, \frac{1}{3}\right) \quad T1: \hat{q}_I = \left(0, \frac{1}{3}, \frac{2}{3}, \frac{2}{3}, 0, 1, \frac{1}{3}\right)$$

Step 2

Having estimated conditional probabilities with messed-up topics:

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

Classify every document in the data using \hat{p}_I and \hat{q}_I to create new estimated conditional probabilities:

$$\hat{p}_{II}, \hat{q}_{II}$$

Then iterate steps 1 and 2 until there is very little change in the topics.

Final Notes

- IF Naive bayse does not work with given true topics, this one will not work either \implies tricks are needed
- **Enormous potential:** if we have annotated only a small portion of the data