

# Text Classification

## A Get-To-Know Introduction

Weiwei Cheng

University of Marburg / Deutsche Bank

[www.chengweiwei.com](http://www.chengweiwei.com)



Deutsche Bank





- Top Stories
- World
- U.S.
- Business
- Sci/Tech
- Entertainment
- Sports
- Health
- Spotlight
- Most Popular

## World »

### [Iraqi Shoe Thrower Says He Was Tortured in Jail](#)

New York Times - [Marc Santora](#) - 17 minutes ago

BAGHDAD - Hours after his release from prison, the Iraqi journalist who hurled his shoes at former President George W. Bush said that he had been tortured while in jail, and his family said that he would flee Iraq, fearing for his life.

[Video: Iraqi shoe-throwing journalist released from prison](#) ITN  

[Bush 'shoe thrower' claims he was tortured in prison](#) CNN International

[Christian Science Monitor](#) - [Los Angeles Times](#) - [guardian.co.uk](#) - [BBC News](#)

[all 1,894 news articles »](#)  [Email this story](#)




Telegraph.co.uk

### [US may have killed al Qaeda target in Somalia, officials say](#)


CNN International - [Barbara Starr](#) - 38 minutes ago

(CNN) -- US special operations forces raid in Somalia on Monday may have killed a wanted al Qaeda terrorist, US officials said. The FBI has this photo of Saleh Ali Saleh Nabhan on its Web site.

[Video: US forces have "killed" leading Al Qaeda suspect](#) ITN  

[Somalia raid likely wins intel haul, stirs tension](#) Reuters

[The Associated Press](#) - [guardian.co.uk](#) - [Bloomberg](#) - [Wikipedia: Saleh Ali Saleh Nabhan](#)

[all 721 news articles »](#)  [Email this story](#)



Financial Times

### [Only Decisive Force Can Prevail in Afghanistan](#)

Wall Street Journal - [Lindsey Graham](#) - 43 minutes ago

Growing numbers of Americans are starting to doubt whether we should have troops in Afghanistan and whether the war there is even winnable.

[SCENARIOS: Obama's options in Afghan war](#) Reuters

[Joint Chiefs Chairman Mullen says more troops likely needed in ...](#) Los Angeles Times

[Bloomberg](#) - [The Associated Press](#) - [AFP](#) - [United Press International](#)

[all 809 news articles »](#)  [Email this story](#)



WA today

## Business »

### [Business Inventories in US Decline in July](#)

Bloomberg - [Timothy R. Homan](#) - 1 hour ago

Sept. 15 (Bloomberg) -- Inventories in the U.S. declined in July, a sign that the economy is slowing.

[Car Sales Lead Increase in US Retail Sales](#) Reuters

[NRF: Retail sales up in Aug. for first time in 10 months](#) Reuters

[Retail Sales in August](#) Reuters

## U.S. »

### [Police: Roseville home broken into, cat killed](#)

Chicago Tribune - 27 minutes ago

AP ROSEVILLE, Mich. - Police in Roseville say at least one burglar broke into a home in the suburb and killed the family cat by putting it in the washing machine.

[Burglar killed cat in washing machine](#) The Detroit News

[burglars ransack home, kill cat](#) Detroit Free Press

[WDIV](#) - [WLNS](#) - [MLive.com](#)

[all 8 news articles »](#)  [Email this story](#)

### [Report: Student stabbed to death at Florida high school](#)

CNN - 46 minutes ago

MIAMI, Florida (CNN) -- Police said Tuesday they are investigating a death at a high school in Coral Gables, Florida. Miami-Dade County police told CNN the death occurred at Coral Gables Senior High School but would not provide further information.

[Student killed at Coral Gables High School](#) MiamiHerald.com

[Student stabbed at Fla. school; 1 in custody](#) The Associated Press

[WESH.com](#) - [NBC2 News](#) - [KMBC.com](#) - [Miami New Times](#)

[all 391 news articles »](#)  [Email this story](#)

### [FBI warns police depts after NYC terror raid](#)

The Associated Press - [Tom Hays](#), [Devlin Barrett](#) - 29 minutes ago

NEW YORK - Counterterrorism officials are warning police departments around the country to be on the lookout for evidence of homemade bombs following raids on several New York City apartments in a hunt for explosives and possible links to al-Qaida ...

[NY Police: Terror Task Force Conducts Raids in Queens](#) Voice of America

[NY Homes Raided in Terror Probe](#) Washington Post

[New York Times](#) - [Christian Science Monitor](#) - [Los Angeles Times](#) - [CBS News](#)

[all 1,232 news articles »](#)  [Email this story](#)

## Sci/Tech »

### [Zune HD: No ipod Killer](#)

PC World - [Daniel Ionescu](#) - 38 minutes ago

Really think it can outdo Apple and its fleet of iPods with the Zune HD? Get a flaky refresh of the iPod touch last week, the Zune HD still doesn't stand up to the competition that well against Apple's popular army of music-playing ...

[Under the hood](#) CNET News

[Zune HD: No ipod Killer](#) News.com





Search

[Advanced Search](#)

## Advanced Search

Use this form to automatically construct your query. (Alternatively, you can type [search operators](#) directly into the search box.)

Find tweets based on...

Search

Words

All of these words

Attitudes

With positive attitude :)

With negative attitude :(

Asking a question ?



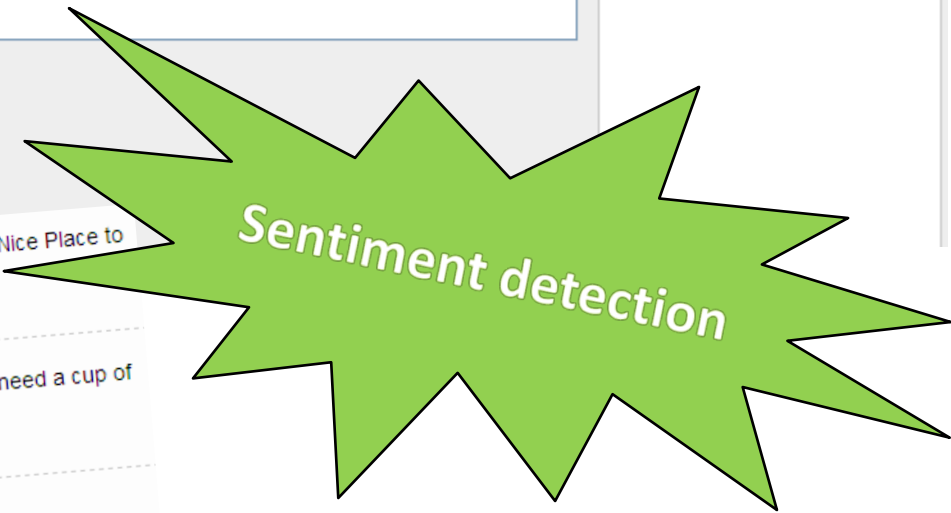
[bergcross121189](#): oh! how i wish I'm in **Germany**!!! :) Such a Nice Place to live..... :)  
35 minutes ago from web · [Reply](#) · [View Tweet](#)



[Susi1984](#): Good Morning from **Germany**. It's 08:18 am and I need a cup of coffee! :)  
about 1 hour ago from web · [Reply](#) · [View Tweet](#)



[denigab](#): Good Morning **Germany**. :) Today is a hard day, too...uff...[@Chaosqueen79](#) okay, honey, that's a deal :D  
about 1 hour ago from web · [Reply](#) · [View Tweet](#)



try it at <http://search.twitter.com/advanced>

# Formally...

## Training (or learning):

Input: a set of  $m$  labeled documents  $(x_1, y_1), \dots, (x_m, y_m)$

Output: a learned classifier  $f: x \rightarrow y$

## Testing (or predicting):

Input: a document  $x$

Output: a class  $y$  from some fixed set of labels  $y_1, \dots, y_K$

# Why we need computers?

If we have *enough* brains,  
we can solve *all* problems!

**Wrong!**

# Can human do this?

Human brain has limited a **capacity**.

*Can you remember all the books from the library?*

# Is the computer only for our “dirty work”?

Human **perception** is fuzzy.

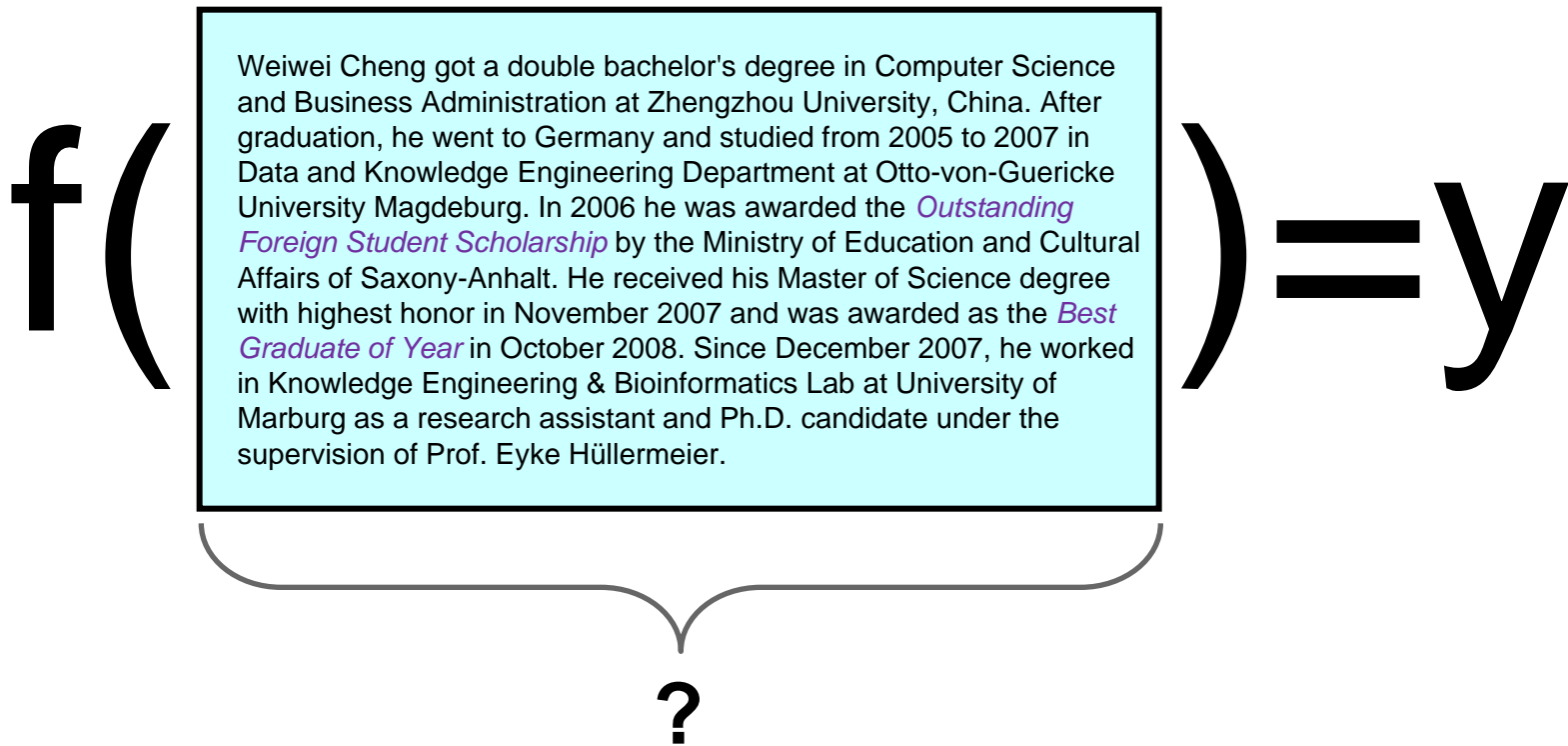
Human **intuition** is (more than often) wrong.

Human **reasoning** is by nature not robust against illusion.

*You like to watch magic, don't you?*

***No!***

# Examples



*Question:*

What is the ~~best~~ representation for the document?

**simplest useful**



# Representing document with bag-of-words

**x** → [ A lady was picking through the frozen turkeys at the supermarket, but couldn't find one big enough for her family. She asked a stock boy, "Do these turkeys get any bigger?" The stock boy replied, "No madam, they're dead." ]

→ [ ladi pick frozen turkei supermarket find big her famili she ask stock boi turkei big stock boi repli madam dead ]

→ [ 166 193 78 ... .. 210 188 60 ]

Word	Index
a	1
:	:
frozen	78
:	:
repli	210
:	:

... ignores the ordering of the words.

# Multinomial Naïve Bayes classifier

- A probabilistic learning method;
- Based on Bayes' theorem;
- State-of-the-art;
- Simple, in terms of implementation and use.

*Great and simple idea  
with fancy name*

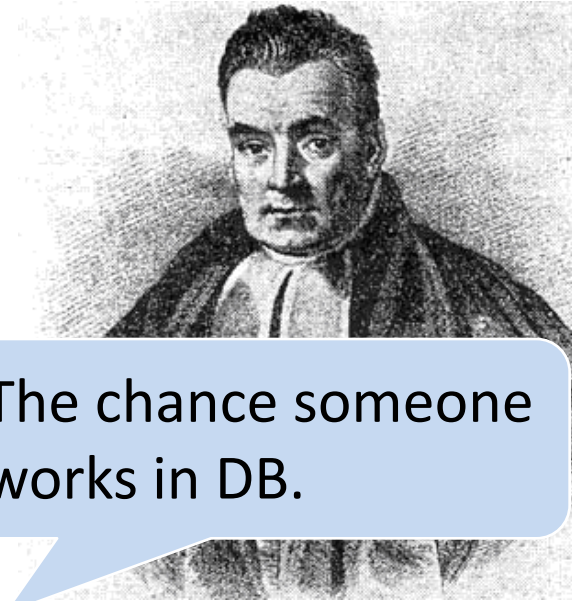
# Some basic probability

Notation	Example
$P(A)$	$P(\textit{someone works in DB})$
$P(B)$	$P(\textit{someone is smart})$
$P(A\&B)$	$P(\textit{someone works in DB and she/he is smart})$
$P(A\cup B)$	$P(\textit{someone works in DB or she/he is smart})$
$P(A B)$	$P(\textit{someone works in DB given she/he is smart})$

$$P(A\&B) = P(B)P(A|B)$$

If A and B are **independent**, e.g.,  
A: *I wear yellow shoes today,*  
B: *Today it rains,*  
then  $P(A\&B) = P(B)P(A)$ .

# Bayes' theorem



The chance that a smart guy works in DB.

The chance that a guy in DB is smart.

The chance someone works in DB.

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

The chance someone is smart.

A: *someone works in DB*

B: *someone is smart*

## Quiz - *drug test*

There is a drug test. It will correctly identify a drug user as testing positive **99%** of the time, and will correctly identify a non-user as testing negative **99%** of the time. Let's assume some company decides to test its employees for drug use, and it is known that **0.5%** of the population actually use the drug. What is the probability that, given a positive drug test, an employee is actually a drug user?

$$\begin{aligned}P(D|+) &= \frac{P(+|D)P(D)}{P(+)} \\ &= \frac{P(+|D)P(D)}{P(+|D)P(D) + P(+|N)P(N)} && \sim 33\% \\ &= \frac{0.99 \times 0.005}{0.99 \times 0.005 + 0.01 \times 0.995} \\ &= 0.3322.\end{aligned}$$

$\mathbf{x}$ : the document (e.g. *an email message*)  
 $y$ : the class (e.g. *spam* or *not spam*)

We want to know which class can lead to the highest

$$P(y|\mathbf{x}) .$$

The probability of the email's *class* (spam/not spam) given its *text*.

By **Bayes' theorem**

$$P(y|\mathbf{x}) = \frac{P(\mathbf{x}|y)P(y)}{P(\mathbf{x})} \propto P(\mathbf{x}|y)P(y)$$

With the **bag-of-words** representation

$$P(y)P(\mathbf{x}|y) = P(y)P(x_1|y)P(x_2|y)\dots P(x_n|y)$$

To choose a *class* that maximize this.

The probability of class  $y$ , e.g., the chance an email being spam.

The probability of word  $x_2$  found in class  $y$ , e.g., the chance of word "pharmacy" seen in a spam email.

# Take-away message

Text classification finds the patterns in documents, reduces the human effort.

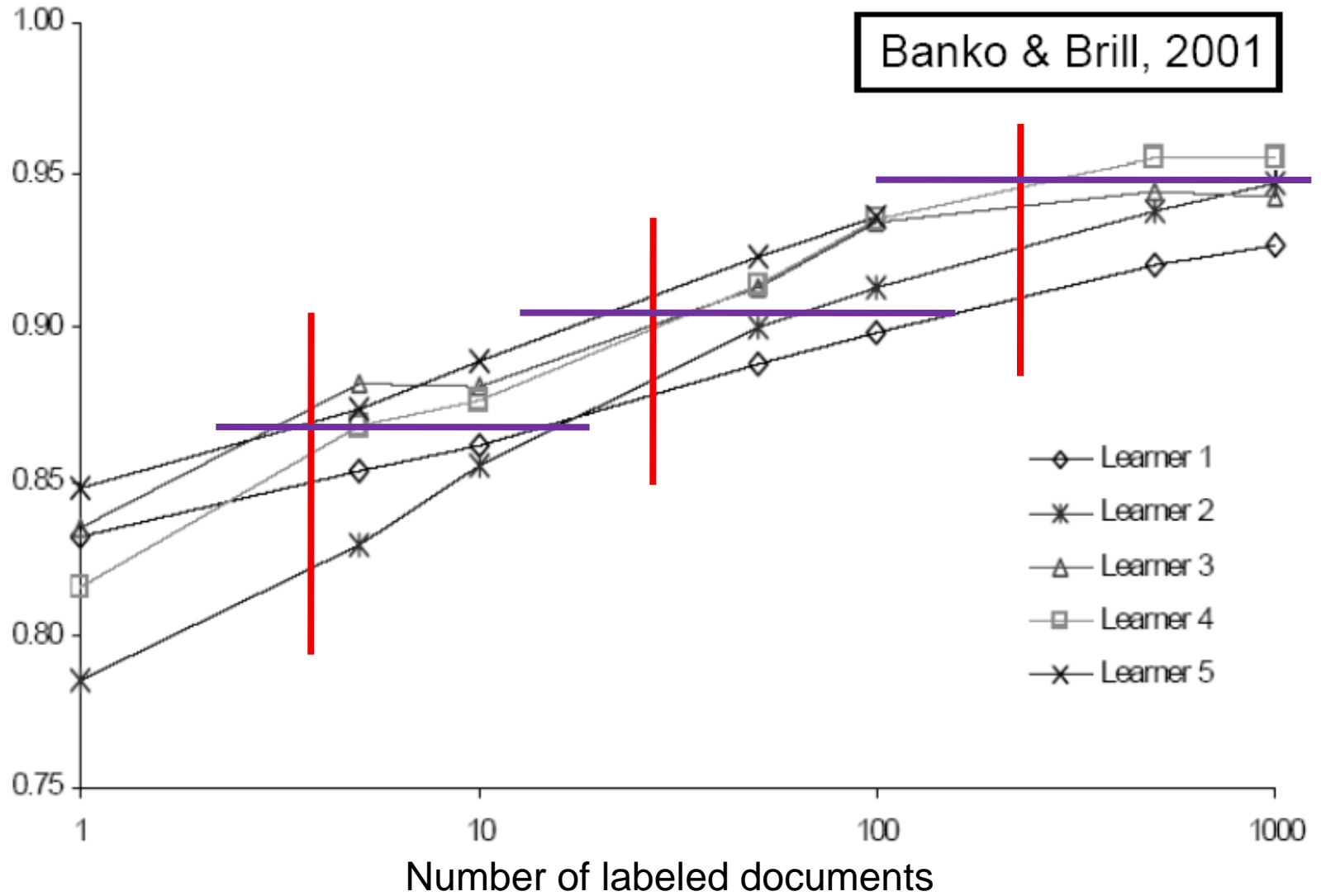
How do we represent the text in the computer?

*bag-of-words*

Which classifier?

*Naive Bayes*. Key idea: *Bayes' theorem*

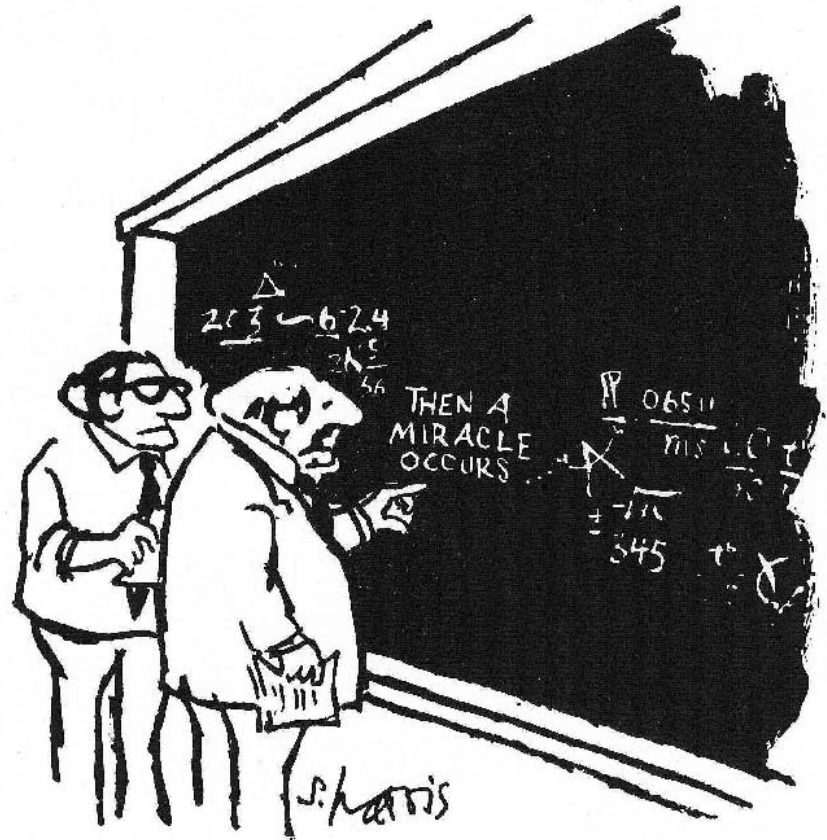
Accuracy





That is all from my side 😊  
Thank you all for being here!

Time for discussion!



"I think you should be more explicit here in step two."

[www.chengweiwei.com](http://www.chengweiwei.com)