

Bayes Theorem simply explained

With applications in Spam classifier and Autocorrect

Hung Tu Dinh

Jan 2018

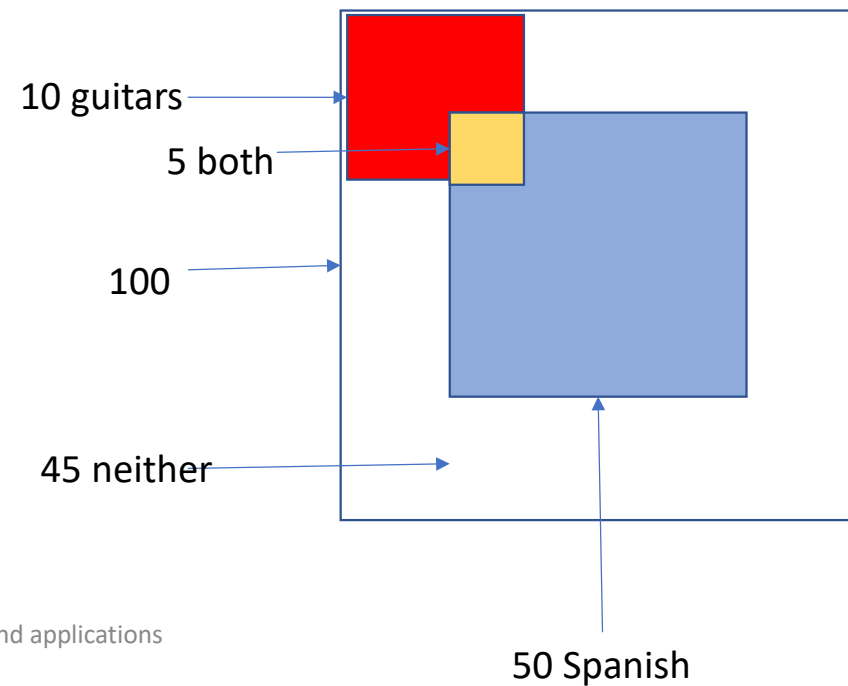
Agenda

- Part 1: Basic concepts of conditional probability and Bayes equation
- Part 2: How does spam filter work?
- Part 3: How does Auto-correct work?

Set up

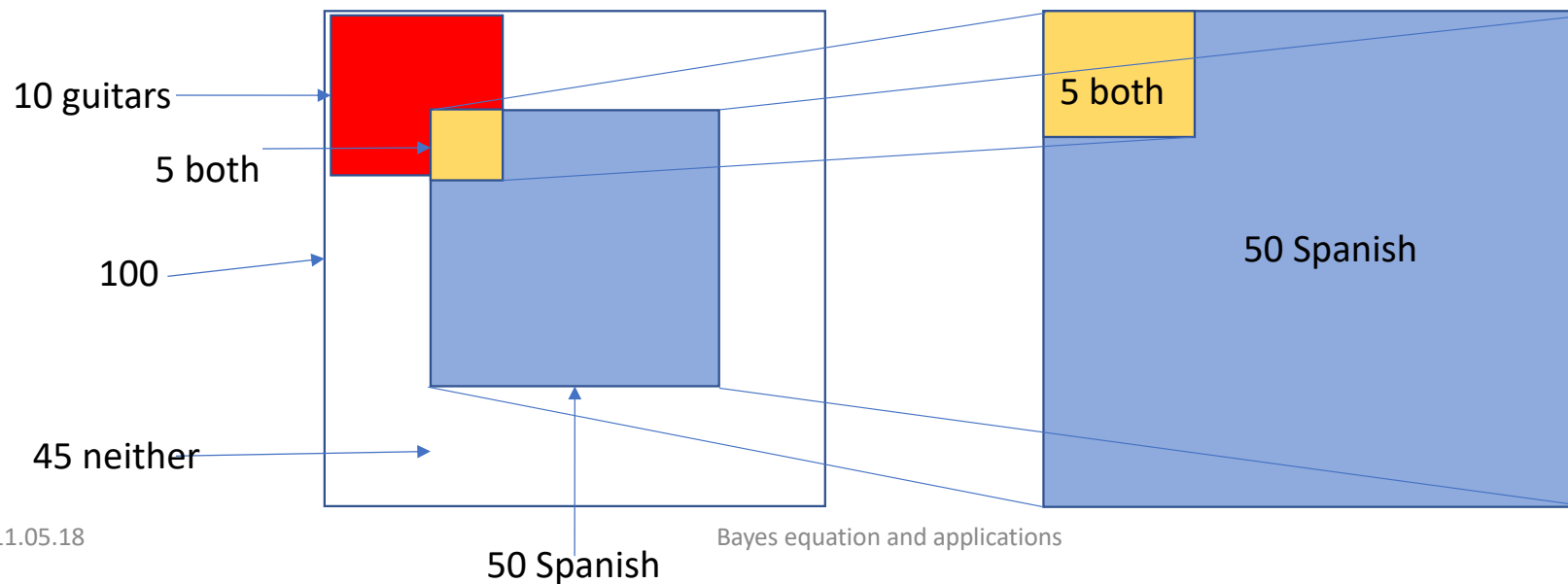
- In a school of 100 students
 - some take guitar in Music class
 - some take Spanish in Language class

Group	Value	Probability
Total	100	1
Guitar	10	0.1
Spanish	50	0.5
Both	5	0.05
Only guitar no Spanish	? 5	?
Neither	? 45	?



Conditional Probability

- $P(\text{Guitar \& Spanish}) = 5/100 = 0.05$
- Among the 50 Spanish learners, only 5 (or 10%) also play guitar
- $P(\text{Guitar \& Spanish} \mid \text{Spanish}) = 5/50 = 0.10$
- $P(\text{Guitar \& Spanish} \mid \text{Guitar}) = ? \quad 5/10=0.5$



Informal way of saying probability

- Probability is “chance” of a variable takes a value
- Music = {guitar, piano, violin}, Language = {Spanish, English, French}
- $P(\text{music=guitar})=0.1$
- $P(\text{lang=Spanish})=0.5$
- $P(\text{music=guitar}|\text{lang=Spanish}) = \frac{\text{Num}(\text{music=guitar \& lang=Spanish})}{\text{Num}(\text{lang=Spanish})} = \frac{\text{Num}(\text{music=guitar \& lang=Spanish}) / \text{Total}}{\text{Num}(\text{lang=Spanish}) / \text{Total}}$
- $P(\text{music=guitar}|\text{lang=Spanish}) = \frac{P(\text{music=guitar \& lang=Spanish})}{P(\text{lang=Spanish})}$

Rearrange the equation

- $P(\text{music=guitar} | \text{lang=Spanish}) = \frac{P(\text{music=guitar} \& \text{lang=Spanish})}{P(\text{lang=Spanish})}$
- Cross multiplication
$$P(\text{music=guitar} \& \text{lang=Spanish}) = P(\text{lang=Spanish}) \times P(\text{music=guitar} | \text{lang=Spanish})$$
$$= 0.5 \times 0.1 = 0.05 \text{ (also } = 5/100=0.05\text{)}.$$
- Intuitive:
 - Among the whole school, 50% learn Spanish.
 - Among the Spanish learners, 10% also learn guitar.
 - \Rightarrow Among the whole school, percentage of students doing both Spanish and guitar is $0.5 \times 0.1 = 0.05$.
- Likewise, switch the order of logic between language and music
- $P(\text{music=guitar} \& \text{lang=Spanish}) = P(\text{lang=Spanish} | \text{music=guitar}) \times P(\text{music} = \text{guitar}) = 0.5 \times 0.1 = 0.05$

Bayes Equation

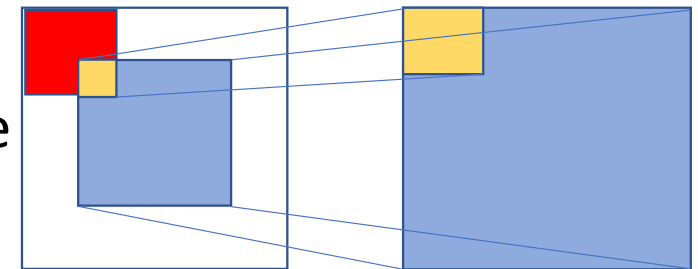
- Using symbol to make the equation more succinct and general
- $P(Y = y|X = x)P(X = x) = P(X = x|Y = y)P(Y = y) = P(X = x \cap Y = y)$
- $P(Y = y|X = x) = \frac{P(X = x|Y = y)P(Y=y)}{P(X=x)}$

Why the equation is needed?

- Because some probabilities is easier to compute than other.
- For example: Spam classifier
 - Objective: to estimate the probability of an email being a spam $P(\text{spam}=\text{True})$
 - Not straightforward to get $P(\text{spam}=\text{True})$
 - we look for other simpler probabilities
 - $$P(\text{spam} = \text{True} | \text{word}_1 = \text{money}) = \frac{P(\text{word}_1 = \text{money} | \text{spam} = \text{True}) \times P(\text{spam}=\text{True})}{P(\text{word}_1=\text{money})}$$
 - More straightforward to
 - Get $P(\text{spam}=\text{True})$ by counting spam emails/all emails available.
 - Among the spam, get the probability of the words e.g. chance of “money”
- More detailed of spam classifier in next part

Summary Part 1

- Probability is chance of a variable takes a (range) value
 - $P(\text{music=guitar})$, $P(\text{language=Spanish})$
- Probability is associate with a population/sample
 - Among the whole school or only Spanish-learners
- Some probabilities are easier to compute
 - Use Bayes equation to compute the difficult one from the easier
 - $$P(Y = y|X = x) = \frac{P(X = x|Y = y)P(Y=y)}{P(X=x)}$$



Part 2

- Frequentist vs Bayesian
 - Keyword: Updating belief
- Spam filtering application
 - Problem description & how it was solved
 - Propose method

Two interpretations of probability

	Frequentist	Bayesian (Belief)
Example	Flip a coin for 100 times, on average, 50 heads and 50 tails	Weather forecast: quantify the belief that tomorrow is a sunny day
Interpretation	$P=1$ means the events always happen	$P=1$ means the belief is 100% certain
Situation	Common in scientific experiments	Belief changes before and after an evidence

Example: Vietnam U23 team in the tournament of 16 teams

- Before the tournament, with no particular evidence, prior belief of winning $\sim 1/16$
- After some early games, with more evidence, update the belief of winning $\sim 1/2$

Spam filter is not only for email



Markiplier 1 day ago

Thank you all so much for watching! We had a blast making this and I hope you're enjoying it as much as we did!

3041

[View all 500 replies](#) ▼



Markiplier 1 day ago

► SteamWallet ► PSNCodes ► XboxCodes ► iTunesCodes ► GooglePlayCodes ► ClashOf ClansGEMS
► AmazonCodes ► RIOTPoints ► FIFAWorldCoins ► MinecraftPremium ► iPhone 6 ►

<https://plus.google.com/115684174143487307725/posts/iB4PnvbjkrF>

So Guys I Really Hope You Enjoyed Video Guys ►► And Presents

[Read more](#)

38



Markiplier 1 day ago

► SteamWallet ► PSNCodes ► XboxCodes ► iTunesCodes ► GooglePlayCodes ► ClashOf ClansGEMS
► AmazonCodes ► RIOTPoints ► FIFAWorldCoins ► MinecraftPremium ► iPhone 6 ►

<https://plus.google.com/115684174143487307725/posts/iB4PnvbjkrF>

Dear Friends I Really Hope You Enjoyed Video Guys ►► And Presents

[Read more](#)

32

How was it solved?

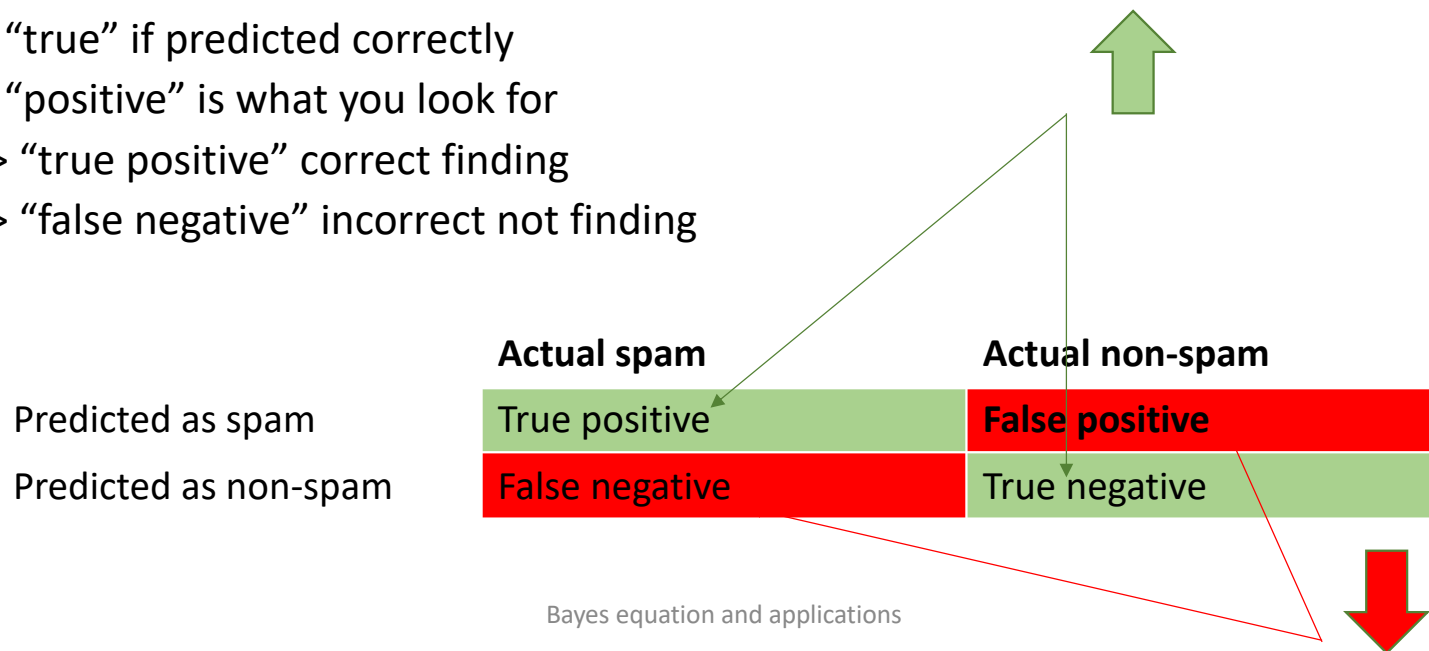
- Fighters: manually identify features to detect spam
 - Common words in spam: money, viagra, discount, deals, ...
- Spammers circumvent



Source: <https://flowingdata.com/2011/10/13/visualizing-yahoo-email-processing-in-real-time/>

Some concepts

- “Spam” ~ undesired # “Scam” ~ fraudulent
- Measure of effectiveness
 - True positive vs false positive:
 - “true” if predicted correctly
 - “positive” is what you look for
 - => “true positive” correct finding
 - => “false negative” incorrect not finding



Calculate spam score

- Find probability of an email being a spam, given a word in the email

$$P(\text{spam} = \text{True} | \text{word}_1 = \text{money}) = \frac{P(\text{word}_1 = \text{money} | \text{spam} = \text{True})}{P(\text{word}_1 = \text{money})} \times P(\text{spam} = \text{True})$$

- Combine probabilities for all words in the email

- Evidence 1: word “money” give $P_1 = P(\text{spam} = \text{True} | \text{word}_1 = \text{money}) = 0.9$
- Evidence 1: word “hi” give $P_2 = P(\text{spam} = \text{True} | \text{word}_1 = \text{hi}) = 0.5$
- ...
- Evidence n: word “hot” give $P_n = P(\text{spam} = \text{True} | \text{word}_n = \text{hot}) = 0.8$

$$P_{\text{combine}} = \frac{P_1 \times \dots \times P_n}{(P_1 \times \dots \times P_n) + (1 - P_1) \times \dots \times (1 - P_n)} \quad \text{simplified} \quad \frac{abc}{abc + (1-a)(1-b)(1-c)}$$

Implementation

$$P(\text{spam} = \text{True} | \text{word}_1 = \text{money}) = \frac{P(\text{word}_1 = \text{money} | \text{spam} = \text{True})}{P(\text{word}_1 = \text{money})} \times P(\text{spam} = \text{True})$$

- Train

Predict

1. Data sets: corpus of spam and corpus of non-spam $P(\text{spam} = \text{True})$
 2. Tokenize (split into words) & count frequency of words in each corpus $P(\text{word}_1 = \text{money} | \text{spam} = \text{True})$ & $P(\text{word}_1 = \text{money} | \text{spam} = \text{False})$
 3. Hash table: $\text{dict}[\text{"money"}] = \frac{P(\text{word}_1 = \text{money} | \text{spam} = \text{True})}{P(\text{word}_1 = \text{money} | \text{spam} = \text{True} \& \text{False})} \times P(\text{spam} = \text{True})$
 4. For each test email:
 1. Tokenize into words
 2. For each word, check in hash table to find $P(\text{spam} = \text{True} | \text{word}_1 = \text{money})$
 3. Combine probabilities from all words $\frac{abc}{abc + (1-a)(1-b)(1-c)}$

Updating the belief

- $P(\text{spam} = \text{True} | \text{word}_1 = \text{money}) = \frac{P(\text{word}_1 = \text{money} | \text{spam} = \text{True})}{P(\text{word}_1 = \text{money})} \times P(\text{spam} = \text{True})$
- No evidence: global rate $P(\text{spam} = \text{True}) = 0.5$.
- With only 1 evidence of word “money” the belief is $P_1 = 1.8 \times 0.5 = 0.9$
- With another evidence of “hi”, $P_2 = 0.5$ the total belief is updated $\frac{0.9 \times 0.5}{0.9 \times 0.5 + (1 - 0.9) \times (1 - 0.5)} = 0.9$
- With another evidence of “hot”, $P_3 = 0.8$ the total belief is updated $\frac{0.9 \times 0.8}{0.9 \times 0.5 + 0.1 \times 0.2} = 0.97$

$$\frac{abc}{abc + (1 - a)(1 - b)(1 - c)}$$

DEMO TIME

- Check the GitHub repository
 - <https://github.com/browning/comment-troll-classifier>

Strength of Bayesian filter

- Identify the features by itself (which words are good - bad)
- Automatically update the probability (“belief”) of spamming words
 - Robust against tricks of spammers
- Take into account both good words & suspicious words
- Customize for individual users (different training sets)

Summary Part 2

- Two interpretations
 - Frequentist: how many times an event occurs
 - Bayesian: how certain is the belief that an event occurring
- Bayes formula
 - Given a hypothesis H and evidence E
 - $P(H|E) = \frac{P(E|H)}{P(E)} \times P(H)$
 - Update belief about H after learning about the evidence E
- Spam filtering
 - Get evidence from each words using Bayes formula
 - Combine evidence from all words
 - Good words and bad words both contribute

Armchair philosophy

- *With more evidences, your belief about the world is updated.*
- *Among many possible evidences, choose wisely.*
- *Beware that we may reinforce our prior belief because of bias.*

Reference

- Brilliant explanation <https://brilliant.org/wiki/bayes-theorem/>
- Some interesting articles from Paul Graham @ “Hackers & Painters”
 - <http://www.paulgraham.com/spam.html> (August 2002)
 - <http://www.paulgraham.com/falsepositives.html>
 - <http://www.paulgraham.com/naivebayes.html>
 - <http://www.paulgraham.com/better.html> (January 2003)