
INFO 2950: Intro to Data Science

Lecture 20
2023-11-06

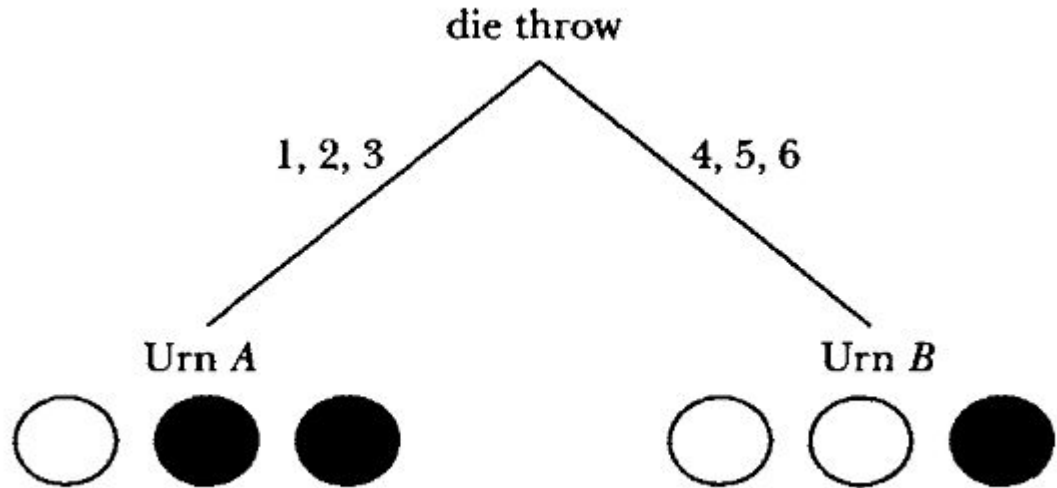
Agenda

1. Bayes Review
2. Text Analysis
3. Log probability
4. Naive Bayes Classifier

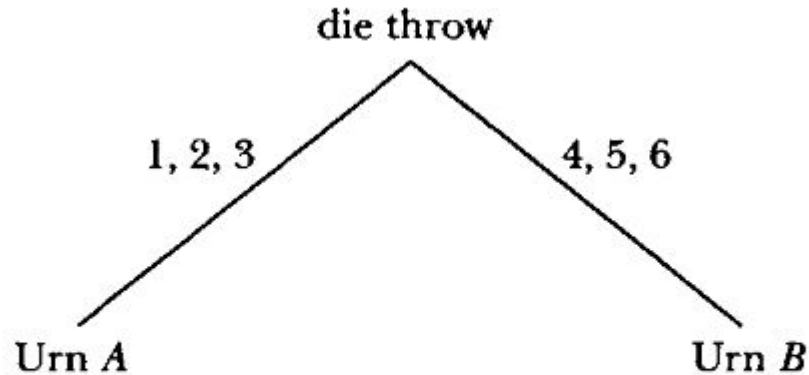
What is Bayes' rule for?

- Lots of teaching examples involve drawing “balls from urns” because it’s an easy way to explain the concept
- But, there are lots of places you might use Bayes’ rule in data science jobs, e.g.: estimating probabilities about spam filters detecting spam emails, given they are spam (or not spam)
- This is why lots of data science interviews test for things like understanding Bayes’ rule!

Balls in Urns!



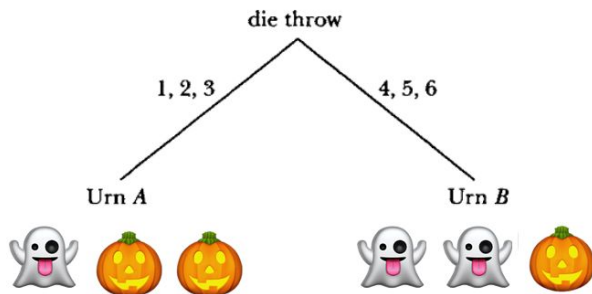
We secretly throw a die!



We won't show you the die roll

- But, we draw a single element from the “urn” that was decided based on the die roll
- What is the probability that what we drew came from Urn A?

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



How did you think about the probability?

- Calculation using Bayes' rule?
- Intuition without Bayes' rule?

Interview question alert!

What is the probability that you roll two fair dice and the sum of their faces is 7?

@Allison inspired by u I've been asking this as an introductory interview question



Everyone ive asked so far has taken at least like 15 secs to answer and done it by brute force



Interview question alert!

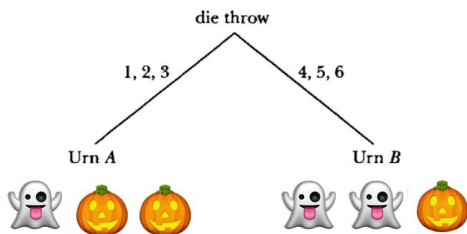
Brute force method:

- $[1,6],[2,5],[3,4],[4,3],[5,2],[6,1] \rightarrow 6$ ways to roll a 7
- $6 \times 6 = 36$ total ways to roll two dice
- $6/36 = \frac{1}{6}$ probability of rolling a 7

Non-brute-force method:

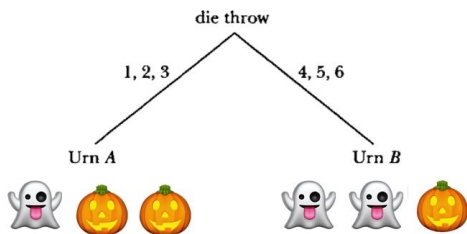
Roll one die. No matter what it lands on, there is exactly one roll by the second die such that the sum equals 7





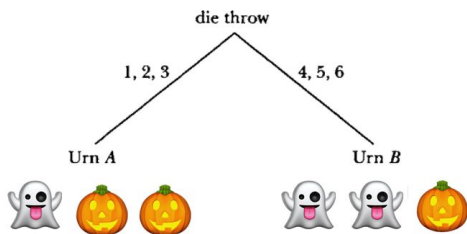
If the first draw was a ghost

- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$



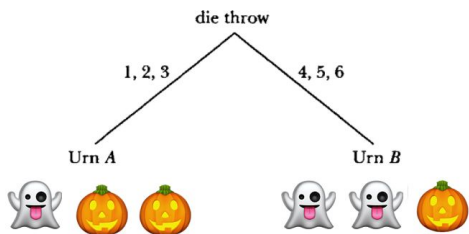
If the first draw was a ghost

- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
Conditional Joint / Marginal



If the first draw was a ghost

- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
 - $P(\text{ghost}, A) = \text{joint probability}$
 $= \frac{1}{2} * \frac{1}{3} = 1 / 6$

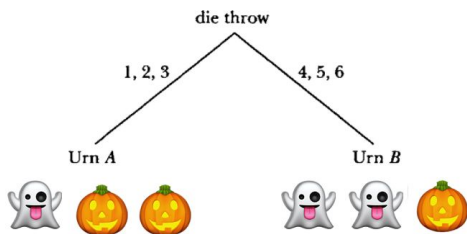


If the first draw was a ghost

- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
 - $P(\text{ghost}, A) = \text{joint probability}$
 $= \frac{1}{2} * \frac{1}{3} = 1 / 6$

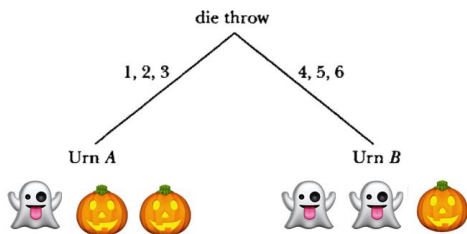
Prob of rolling
1, 2, 3

Within Urn A,
getting a ghost



If the first draw was a ghost

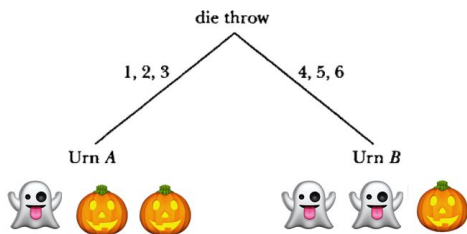
- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
 - $P(\text{ghost}, A) = \text{joint probability} = 1 / 6$
 - $P(\text{ghost}) = 1/2$



If the first draw was a ghost

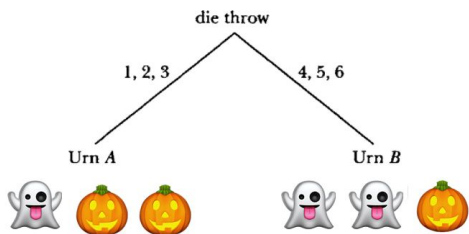
- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
 - $P(\text{ghost}, A) = \text{joint probability} = 1 / 6$
 - $P(\text{ghost}) = 1/2$

Across the urns, there are 6 items to draw, and 3 of them are ghosts



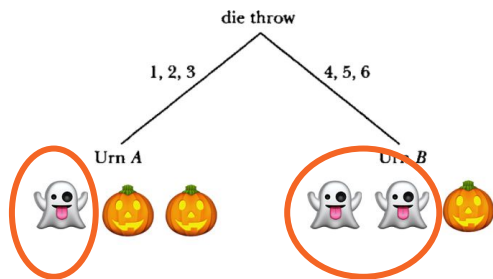
If the first draw was a ghost

- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
 - $P(\text{ghost}, A) = \text{joint probability} = 1 / 6$
 - $P(\text{ghost}) = 1/2$
 - $P(A|\text{ghost}) = \frac{1}{6} / \frac{1}{2} = \frac{1}{3}$



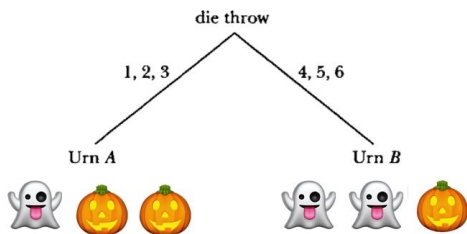
If the first draw was a ghost

- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
 - $P(A|\text{ghost}) = \frac{1}{6} / \frac{1}{2} = \frac{1}{3}$
- Intuition without Bayes' rule
 - ?



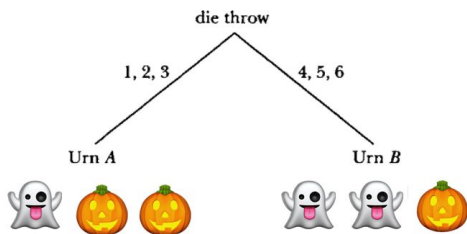
If the first draw was a ghost

- Calculation using Bayes' rule
 - $P(A|\text{ghost}) = P(\text{ghost}, A) / P(\text{ghost})$
 - $P(A|\text{ghost}) = \frac{1}{6} / \frac{1}{2} = \frac{1}{3}$
- Intuition without Bayes' rule
 - All 6 items are equally likely to be drawn before the die is thrown! (This is called our *prior*)
 - $\frac{1}{3}$ of the ghosts overall are in Urn A



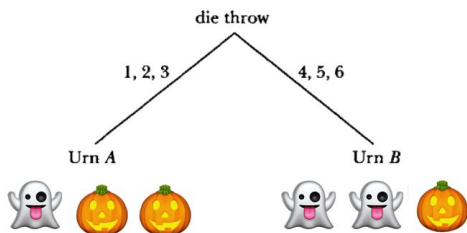
If the first draw was a pumpkin

- Calculation using Bayes' rule
 - $P(A|\text{pumpkin}) = P(\text{pumpkin}, A) / P(\text{pumpkin})$
 - **$P(\text{pumpkin}, A)$** = joint probability = ?
 - **$P(\text{pumpkin})$** = ?



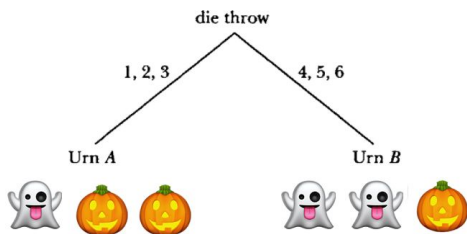
If the first draw was a pumpkin

- Calculation using Bayes' rule
 - $P(A|\text{pumpkin}) = P(\text{pumpkin}, A) / P(\text{pumpkin})$
 - $P(\text{pumpkin}, A) = 2/6 = \mathbf{1/3}$
 - $P(\text{pumpkin}) = \mathbf{1/2}$



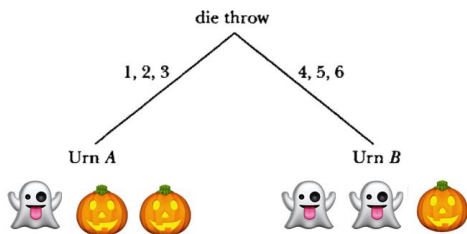
If the first draw was a pumpkin

- Calculation using Bayes' rule
 - $P(A|\text{pumpkin}) = P(\text{pumpkin}, A) / P(\text{pumpkin})$
 - $P(\text{pumpkin}, A) = 2/6 = \mathbf{1/3}$
 - $P(\text{pumpkin}) = \mathbf{1/2}$
 - Joint probability: $\Pr(\text{rolling a } 1,2,3) * \Pr(\text{drawing a pumpkin from urn A}) = \frac{1}{2} * \frac{2}{3} = 2/6 = \mathbf{1/3}$
 - Marginal: $\Pr(\text{pumpkin}) = 3 \text{ pumpkins} / 6 \text{ items across both urns} = \mathbf{1/2}$



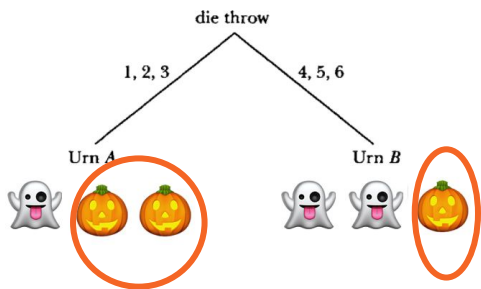
If the first draw was a pumpkin

- Calculation using Bayes' rule
 - $P(A|\text{pumpkin}) = P(\text{pumpkin}, A) / P(\text{pumpkin})$
 - $P(\text{pumpkin}, A) = 2/6 = \mathbf{1/3}$
 - $P(\text{pumpkin}) = \mathbf{1/2}$
 - $P(A|\text{pumpkin}) = 1/3 / 1/2 = \mathbf{2/3}$



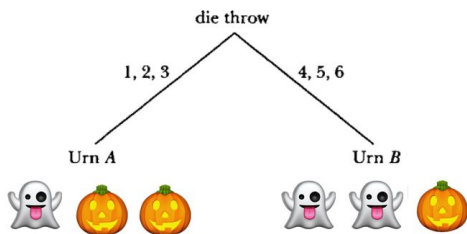
If the first draw was a pumpkin

- Calculation using Bayes' rule
 - $P(A|\text{pumpkin}) = P(\text{pumpkin}, A) / P(\text{pumpkin})$
 - $P(A|\text{pumpkin}) = \frac{1}{3} / \frac{1}{2} = \frac{2}{3}$
- Intuition without Bayes' rule
 - ?



If the first draw was a pumpkin

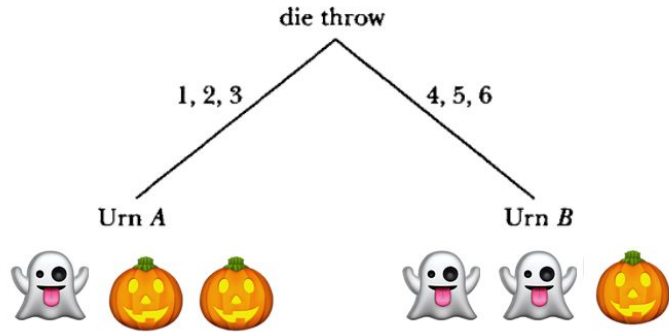
- Calculation using Bayes' rule
 - $P(A|\text{pumpkin}) = P(\text{pumpkin}, A) / P(\text{pumpkin})$
 - $P(A|\text{pumpkin}) = \frac{1}{3} / \frac{1}{2} = \frac{2}{3}$
- Intuition without Bayes' rule
 - All 6 items are equally likely to be drawn before the die is thrown! (This is called our *prior*)
 - $\frac{2}{3}$ of the pumpkins are in Urn A



Let's assume our first draw was a pumpkin

- Now we replace the pumpkin in the same urn we drew it out of (i.e., we drew *with replacement*)
- And now, we take a second draw from the **same urn** as before

Let's assume our first draw was a pumpkin



- Now we replace the pumpkin in the same urn we drew it out of (i.e., we drew *with replacement*)
- And now, we take a **second** draw from the **same urn** as before
- **Question: what is the probability that what we drew came from Urn A?**

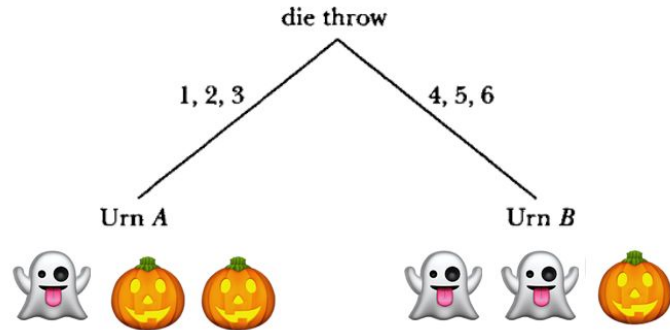
How did you think about the probability?

- Calculation using Bayes' rule?
- Intuition without Bayes' rule?

How did you think about the probability?

- Calculation using Bayes' rule?
- Intuition without Bayes' rule?

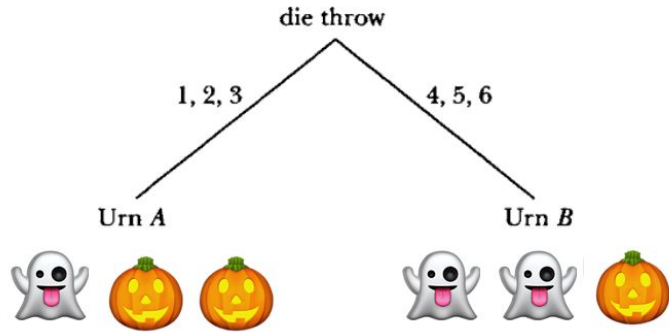
Bayes' rule

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


Bayes' rule

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

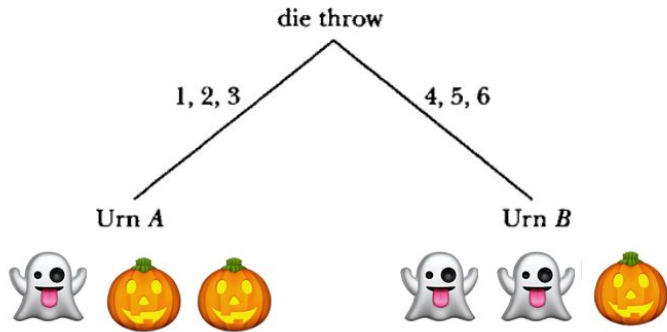
- First draw from urns was a pumpkin. It is returned to original urn **with replacement**
- And now, we take a second draw from the **same urn**, and get a pumpkin
- What is the **probability we drew from Urn A**?



Bayes' rule

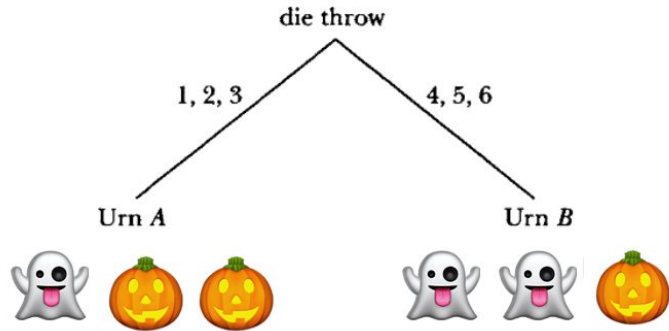
$$P(A | B) = \frac{P(B | A) \cdot P(A)}{P(B)}$$

- First draw from urns was a pumpkin. It is returned to original urn **with replacement**
- And now, we take a second draw from the **same urn**, and get a pumpkin
- What is the **probability we drew from Urn A**?
- How do we define A and B if we use Bayes' rule?



Bayes' rule

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

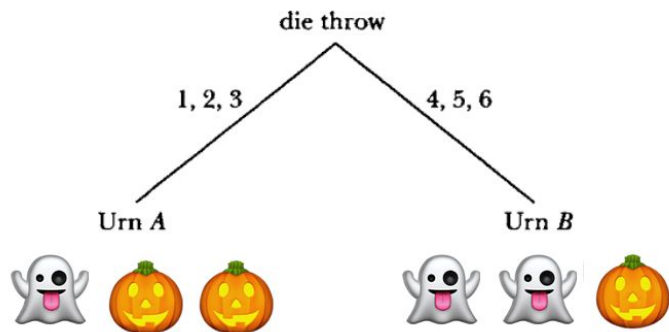


A = drawing from urn A

B = 1st draw 🎃 and 2nd draw
🎃

Bayes' rule: solve!

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



A = drawing from urn A

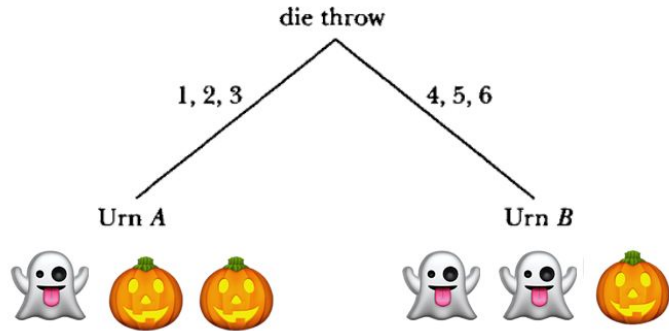
B = 1st draw 🎃 and 2nd draw





Bayes' rule

$$P(A, B) = \frac{2}{3} * \frac{2}{3} =$$

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



A = drawing from urn A

B = 1st draw  and 2nd draw 

Bayes' rule

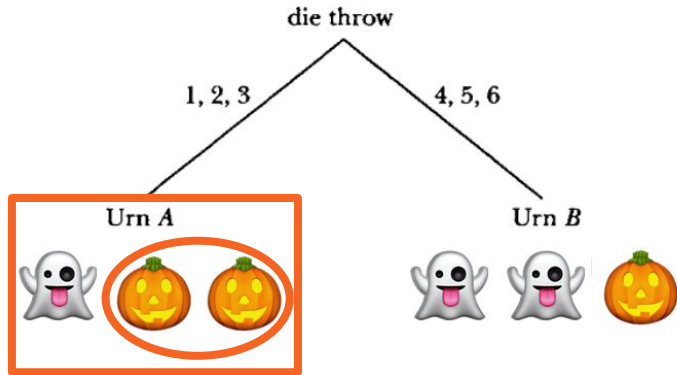
$$P(A, B) = \frac{2}{3} * \frac{2}{3} =$$

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

$P(\text{1st draw pumpkin, 2nd draw pumpkin, from Urn A}) = \frac{2}{3} * \frac{2}{3} = \frac{4}{9}$

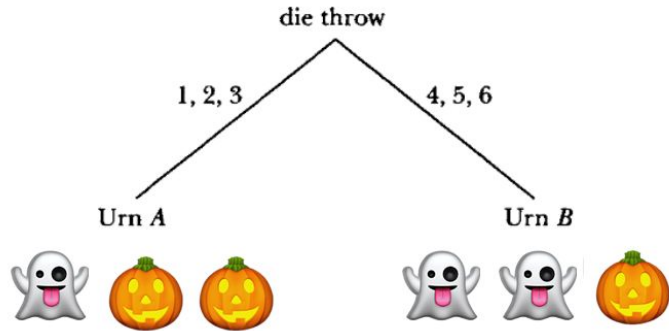
A = drawing from urn A

B = 1st draw 🎃 and 2nd draw
🎃





Bayes' rule

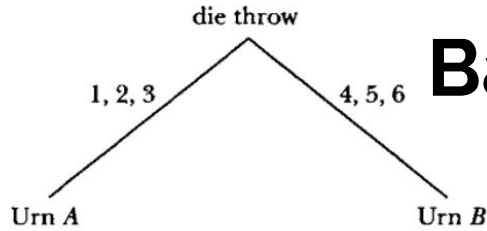
$$P(A | B) = \frac{P(B | A) \cdot P(A)}{P(B)}$$
$$= \frac{2}{3} * \frac{2}{3} + \frac{1}{3} * \frac{1}{3}$$



A = drawing from urn A

B = 1st draw  and 2nd draw 

Bayes' rule



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

$$= \frac{2}{3} * \frac{2}{3} + \frac{1}{3} * \frac{1}{3}$$
$$= P(B,A) + P(B, \text{not-A})$$

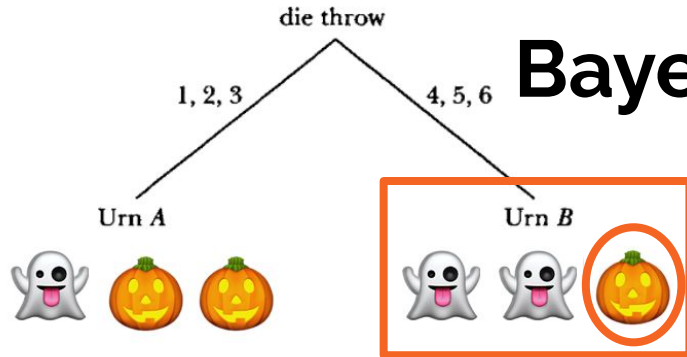
- $P(\text{1st draw pumpkin, 2nd draw pumpkin, Urn A}) = \frac{2}{3} * \frac{2}{3} = \frac{4}{9}$

A = drawing from urn A

B = 1st draw  and 2nd draw



Bayes' rule



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

$$= \frac{2}{3} * \frac{2}{3} + \frac{1}{3} * \frac{1}{3}$$
$$= P(B, A) + P(B, \text{not-A})$$

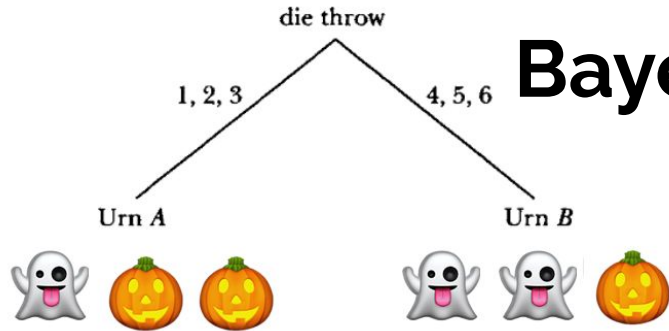
- $P(\text{1st draw pumpkin, 2nd draw pumpkin, Urn A}) = \frac{2}{3} * \frac{2}{3} = \frac{4}{9}$
- $P(\text{1st draw pumpkin, 2nd draw pumpkin, Urn B}) = \frac{1}{3} * \frac{1}{3} = \frac{1}{9}$

A = drawing from urn A

B = 1st draw  and 2nd draw



Bayes' rule



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

$$= \frac{2}{3} * \frac{2}{3} + \frac{1}{3} * \frac{1}{3}$$
$$= P(B, A) + P(B, \text{not-A})$$

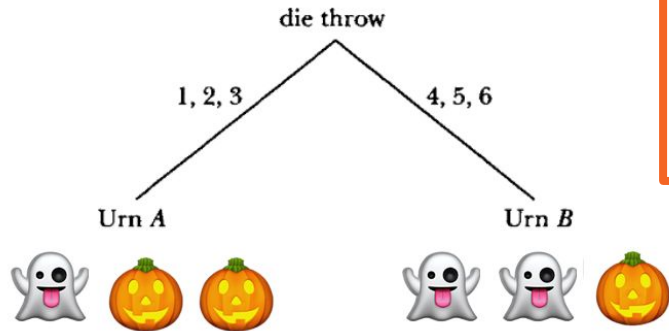
- $P(\text{1st draw pumpkin, 2nd draw pumpkin, Urn A}) = \frac{2}{3} * \frac{2}{3} = \frac{4}{9}$
- $P(\text{1st draw pumpkin, 2nd draw pumpkin, Urn B}) = \frac{1}{3} * \frac{1}{3} = \frac{1}{9}$
- $P(\text{1st draw pumpkin, 2nd draw pumpkin}) = P(\text{1st draw pumpkin, 2nd draw ghost, urn A}) + P(\text{1st draw pumpkin, 2nd draw ghost, urn B}) = \frac{5}{9}$

A = drawing from urn A

B = 1st draw  and 2nd draw



Bayes' rule



$$P(A | B) = \frac{P(B | A) \cdot P(A)}{P(B)}$$
$$= (4/9) / (5/9)$$
$$= 4/5$$

$$P(A, B) = \frac{2}{3} * \frac{2}{3} =$$
$$P(B | A) \cdot P(A)$$

$$P(B)$$
$$= \frac{2}{3} * \frac{2}{3} + \frac{1}{3} * \frac{1}{3}$$

A = drawing from urn A

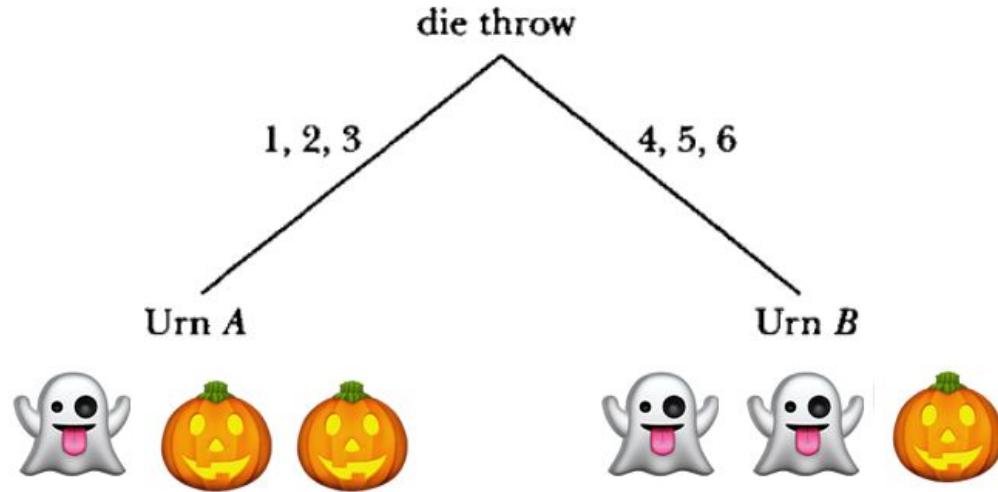
B = 1st draw 🎃 and 2nd draw



How did you think about the probability?

- Calculation using Bayes' rule?
- Intuition without Bayes' rule?

Can we express the “posterior belief” after the 1st roll?

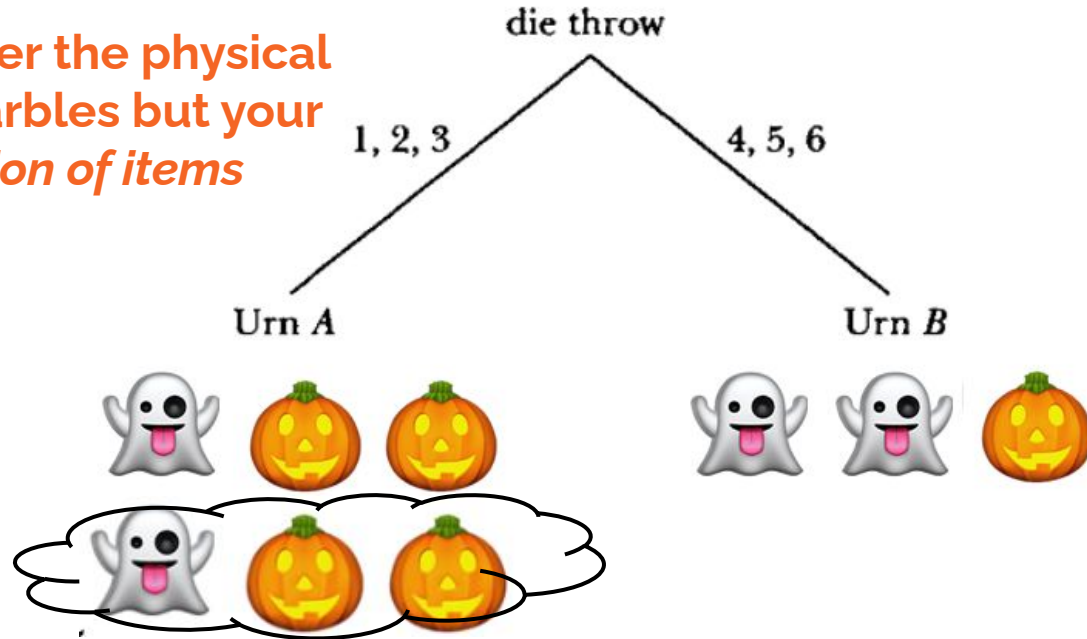


My beliefs before the 2nd draw are that the probability that Urn A is used is $\frac{2}{3}$ (i.e. it's twice as likely that Urn A is used as Urn B).

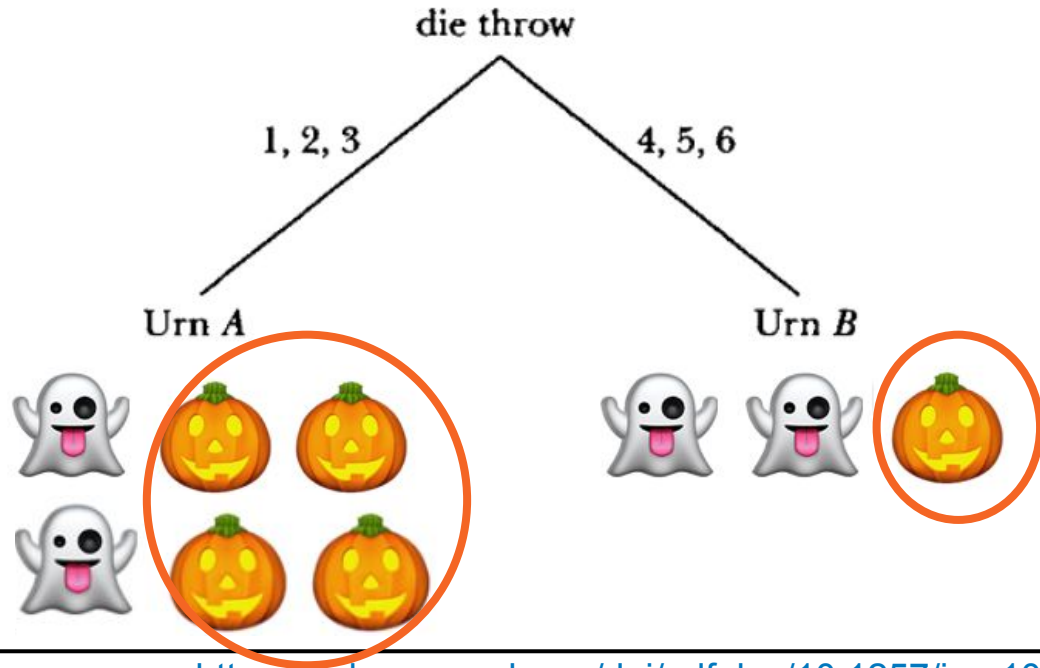
Can we express the “posterior belief” after the 1st roll?

This figure is no longer the physical representation of marbles but your *belief of the distribution of items*

My beliefs before the 2nd draw are that the probability that Urn A is used is $\frac{2}{3}$ (i.e. it's twice as likely that Urn A is used as Urn B).

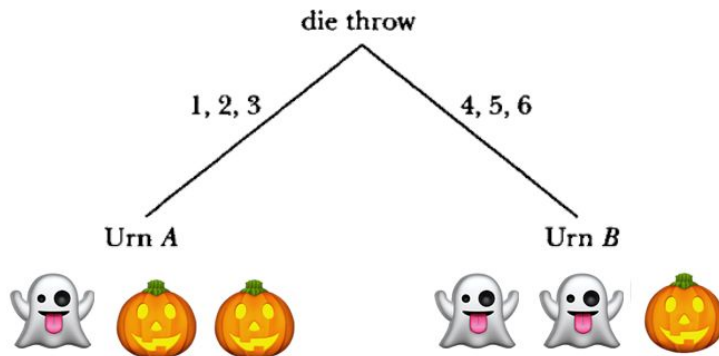


2nd draw a pumpkin: $\frac{4}{6}$ prob from A

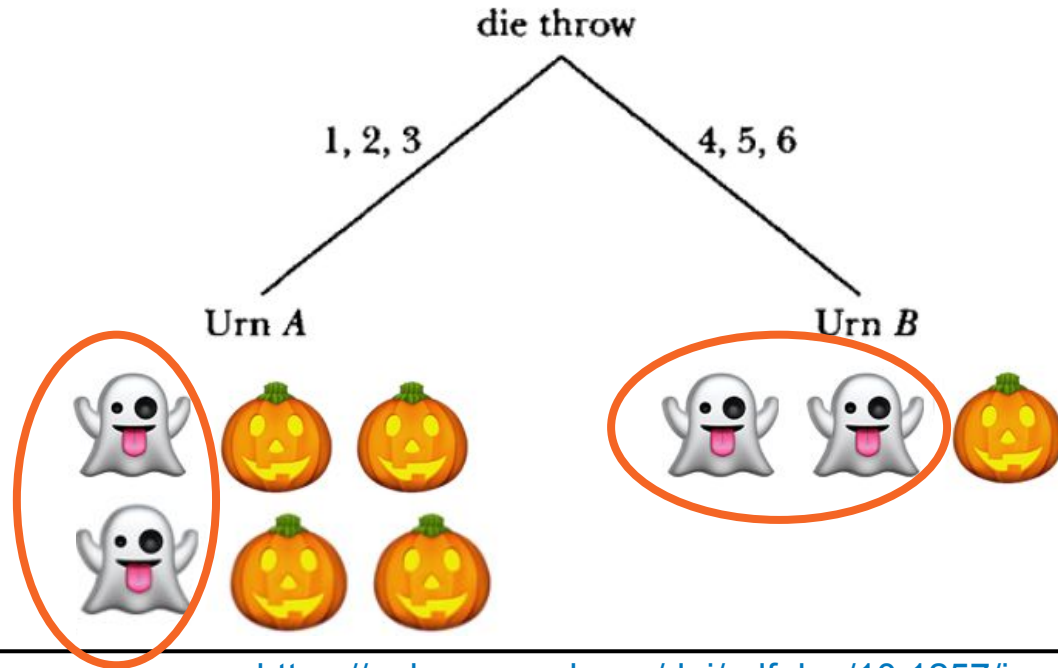


Practice problem for home

- First draw is pumpkin, second draw (with replacement) is ghost. Now, what is the probability that we drew from Urn A?



2nd draw a ghost: $\frac{1}{2}$ prob from A



If the second draw was a ghost

- $P(A | \text{1st draw pumpkin, 2nd draw ghost}) = P(\text{1st draw pumpkin, 2nd draw ghost, A}) / P(\text{1st draw pumpkin, 2nd draw ghost})$ (Bayes' Rule)
- $P(\text{1st draw pumpkin, 2nd draw ghost, A}) = \frac{2}{3} * \frac{1}{3} = \frac{2}{9}$
- $P(\text{1st draw pumpkin, 2nd draw ghost, B}) = \frac{1}{3} * \frac{2}{3} = \frac{2}{9}$
- $P(\text{1st draw pumpkin, 2nd draw ghost}) = P(\text{1st draw pumpkin, 2nd draw ghost, A}) + P(\text{1st draw pumpkin, 2nd draw ghost, B}) = \frac{4}{9}$ (Marginalizing!)
- $P(A | \text{1st draw pumpkin, 2nd draw ghost}) = (\frac{2}{9}) / (\frac{4}{9}) = \frac{1}{2}$

1 min break & attendance



tinyurl.com/yhajjfvk

Now, let's talk about text!

(we promise this will all connect back to Bayes!)

Goal: classify a sentence as positive or negative

A) "I've long been searching for the best camarones diablo and have found that gem here in their camarones endiablados offering."

B) "We walked in and nobody greeted us at the entrance and we decided to walk over to one of the open tables."

What do you guess (+ or -)?

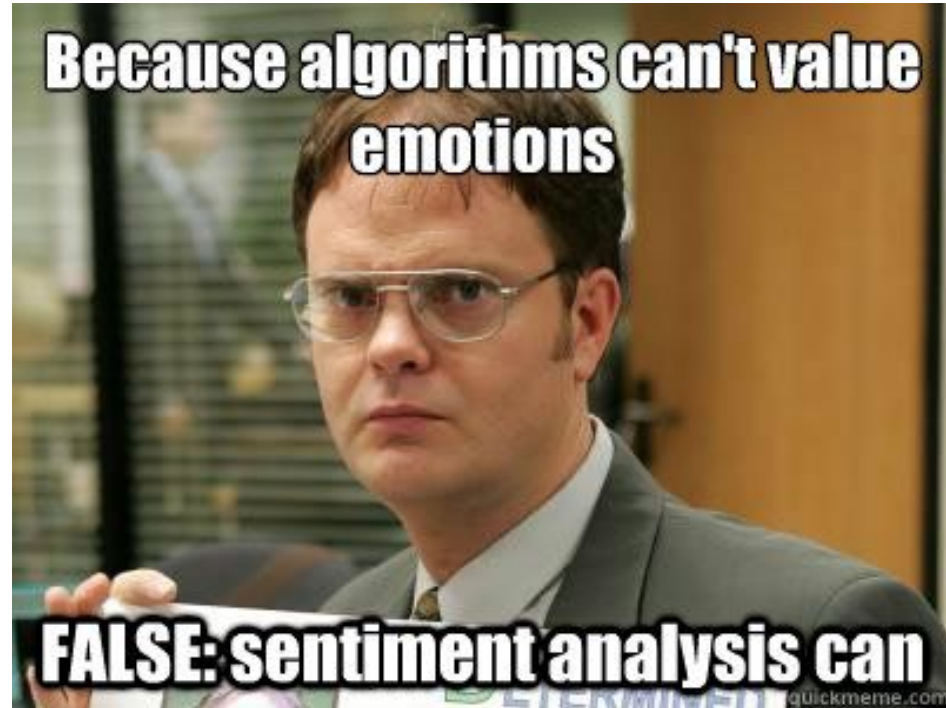
Goal: classify a sentence as positive or negative

+

A) "I've long been searching for the best camarones diabla and have found that gem here in their camarones endiablados offering." ★★★★★

-

B) "We walked in and nobody greeted us at the entrance and we decided to walk over to one of the open tables." ★



Apply Bayes' rule

$P(\star\star\star\star\star \mid \text{"great food, but service stank"})$

Apply Bayes' rule

$P(\star\star\star\star\star \mid \text{"great food, but service stank"}) =$

$P(\text{"great food, but service stank"} \mid \star\star\star\star\star) P(\star\star\star\star\star)$

$P(\text{"great food, but service stank"})$

Apply Bayes' rule

$P(\star\star\star\star\star \mid \text{"great food, but service stank"}) =$

$P(\text{"great food, but service stank"} \mid \star\star\star\star\star) P(\star\star\star\star\star)$

$P(\text{"great food, but service stank"})$

Colorless green ideas sleep furiously

Article [Talk](#)

From Wikipedia, the free encyclopedia

"Green ideas" redirects here. For the book series, see [Green Ideas](#).

Colorless green ideas sleep furiously was composed by [Noam Chomsky](#) in his 1957 book *[Syntactic Structures](#)* as an example of a [sentence](#) that is [grammatically well-formed](#), but [semantically nonsensical](#). The sentence was originally used in his 1955 thesis *[The Logical Structure of Linguistic Theory](#)* and in his 1956 paper "Three Models for the Description of Language".^{[1]:116}

But it must be recognized
that the notion of
"probability of a sentence"
is an entirely useless one,
under any known
interpretation of this term.



Colorless green ideas learn furiously: Chomsky and the two cultures of statistical learning

Peter Norvig

First published: 09 August 2012 | <https://doi.org/10.1111/j.1740-9713.2012.00590.x> | Citations: 4

☰ SECTIONS



PDF



TOOLS



SHARE

Abstract

Language recognition programs use massive databases of words, and statistical correlations between those words, to translate or to recognise speech. But correlation is not causation. Do these statistical data-dredgings give any insight into how language works? Or are they a mere big-number trick, useful but adding nothing to understanding? One who holds the latter view is the theorist of language Noam Chomsky. **Peter Norvig** disagrees.



Clearly, it is inaccurate to say that statistical models (and probabilistic models) have achieved limited success; rather they have achieved a dominant (although not exclusive) position.

Assigning probability to a sentence

$P(\text{"great food, but service stank"} \mid \star\star\star\star\star)$

Assigning probability to a sentence

$P(\text{"great", "food", "but", "service", "stank"} \mid \star \star \star \star \star)$

Split a string into *tokens*

"the food was great but the service was bad"

Split a string into *tokens*

"the food was great but the service was bad"

["the", "food", "was", "great", "but", "the", "service", "was", "bad"]

We tokenize the string into 9 word tokens

Split a string into *tokens*

"the food was great but the service was bad"

["**the**", "food", "**was**", "great", "but", "**the**", "service", "**was**", "bad"]

{'the': 2, 'was': 2, 'food': 1, 'great': 1, 'but': 1, 'service': 1, 'bad': 1}

The string contains 7 distinct word types

Assigning probability to a sentence

$$\begin{aligned} P(\text{"great food, but service stank"} \mid \star\star\star\star\star) = \\ P(\text{"great"} \mid \star\star\star\star\star) \times \\ P(\text{"food"} \mid \star\star\star\star\star) \times \\ P(\text{"but"} \mid \star\star\star\star\star) \times \\ P(\text{"service"} \mid \star\star\star\star\star) \times \\ P(\text{"stank"} \mid \star\star\star\star\star) \end{aligned}$$

Assigning probability to a sentence

Approximate
the probability
of a sentence
by multiplying
the probability
of each word

$$\begin{aligned} P(\text{"great food, but service stank"} \mid \star\star\star\star\star) = \\ P(\text{"great"} \mid \star\star\star\star\star) \times \\ P(\text{"food"} \mid \star\star\star\star\star) \times \\ P(\text{"but"} \mid \star\star\star\star\star) \times \\ P(\text{"service"} \mid \star\star\star\star\star) \times \\ P(\text{"stank"} \mid \star\star\star\star\star) \end{aligned}$$

Assigning probability to a sentence

Naïve
assumption

or

"bag of
words"
assumption

$$\begin{aligned} P(\text{"great food, but service stank"} \mid \star\star\star\star\star) = \\ P(\text{"great"} \mid \star\star\star\star\star) \times \\ P(\text{"food"} \mid \star\star\star\star\star) \times \\ P(\text{"but"} \mid \star\star\star\star\star) \times \\ P(\text{"service"} \mid \star\star\star\star\star) \times \\ P(\text{"stank"} \mid \star\star\star\star\star) \end{aligned}$$

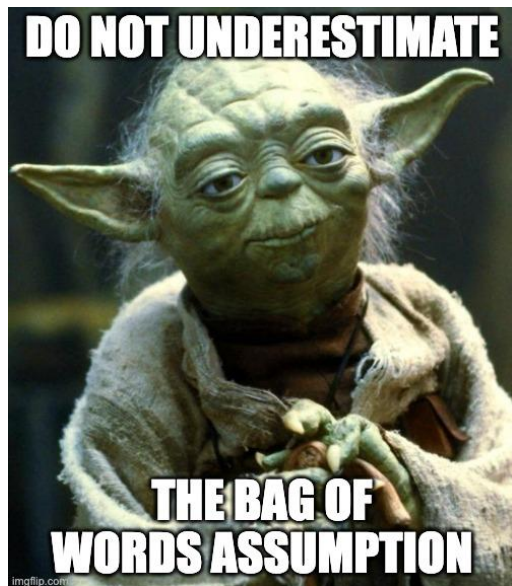
Bag of words: order doesn't matter

$P(\text{"great food, but service stank"} \mid \star\star\star\star\star) =$
 $P(\text{"great service, but food stank"} \mid \star\star\star\star\star)$

Bag of words: order doesn't matter

$P(\text{"great food, but service stank"} \mid \star\star\star\star\star) =$
 $P(\text{"great service, but food stank"} \mid \star\star\star\star\star) =$
 $P(\text{"stank service food but great"} \mid \star\star\star\star\star)$

Bag of words: order doesn't matter

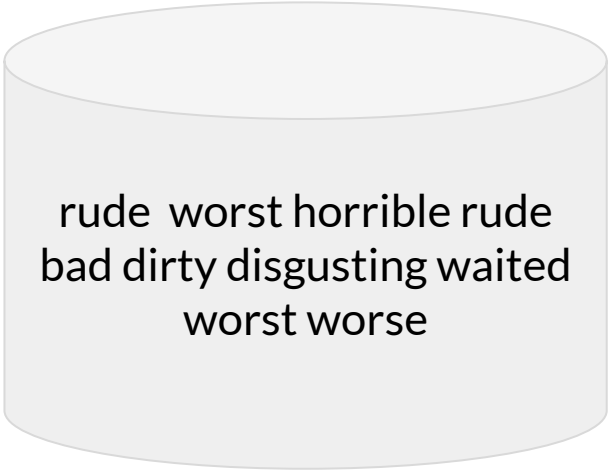


$P(\text{"great food, but service stank"} \mid \star\star\star\star\star) =$
 $P(\text{"great service, but food stank"} \mid \star\star\star\star\star) =$
 $P(\text{"stank service food but great"} \mid \star\star\star\star\star)$


Which source are we sampling?



Which source are we sampling?



rude worst horrible rude
bad dirty disgusting waited
worst worse



fantastic best perfect
incredible gem BEST perfect
delicious LOVE

Multiplying Probabilities

- If $P(E)$ is the probability of event E occurring, and $P(F)$ is the probability of event F occurring:
- E and F are independent events
- What is the probability that both E and F occur?

Multiplying Probabilities

- If $P(E)$ is the probability of event E occurring, and $P(F)$ is the probability of event F occurring:
- E and F are independent events
- What is the probability that both E and F occur?
 - $P(E \text{ and } F) = P(E, F) = P(E) * P(F)$

Probability

- If $P(E)$ is the probability of event E occurring, then we know:

$$\boxed{} \leq P(E) \leq \boxed{}$$

Probability

- If $P(E)$ is the probability of event E occurring, then we know:

$$0 \leq P(E) \leq 1$$

Log Probability

In **probability theory** and **computer science**, a **log probability** is simply a **logarithm** of a **probability**.

Log Probability

- If $P(E)$ is the probability of event E occurring, then we know:

$$0 \leq P(E) \leq 1$$

$$\boxed{} \leq \log P(E) \leq \boxed{}$$

Log Probability

- If $P(E)$ is the probability of event E occurring, then we know:

$$\begin{aligned} 0 &\leq P(E) \leq 1 \\ -\infty &\leq \log P(E) \leq 0 \end{aligned}$$

Log Probability

- If $P(E)$ is the probability of event E occurring, and $P(F)$ is the probability of event F occurring:

$$\log(P(E) \cdot P(F)) = \boxed{} + \boxed{}$$

Log Probability

- If $P(E)$ is the probability of event E occurring, and $P(F)$ is the probability of event F occurring:

$$\log(P(E) \cdot P(F)) = \log P(E) + \log P(F)$$

Log Probability

- If we already know the values of $P(E)$ and $P(F)$, why would we want to take the logarithm?

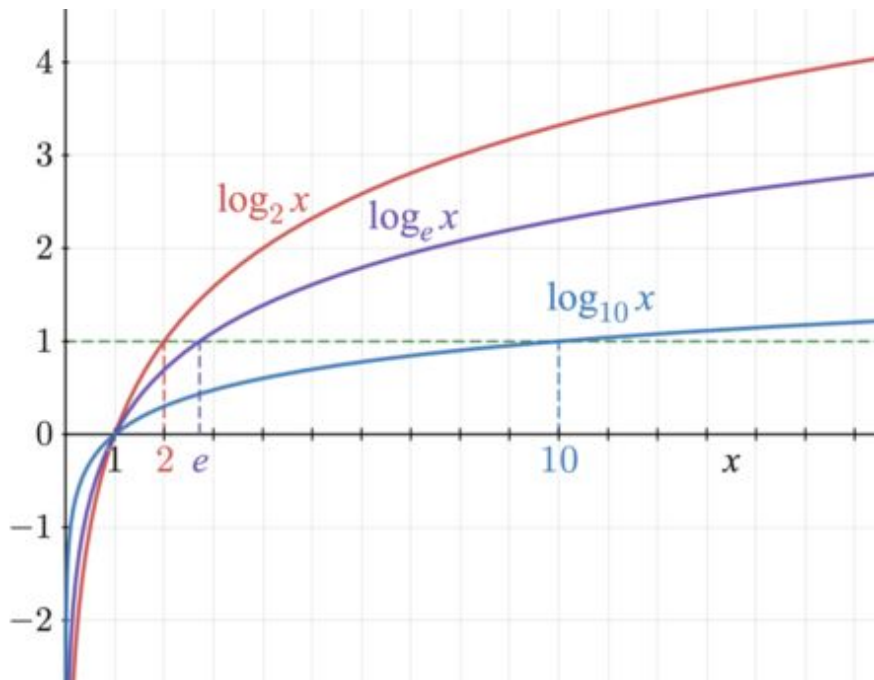
$$\log(P(E) \cdot P(F)) = \log P(E) + \log P(F)$$

Log Probability

One in	Probability	Log_{10}	Log_e
10	0.1	-1	-2.3
100	0.01	-2	-4.6
1,000	0.001	-3	-6.9
10,000	0.0001	-4	-9.2
100,000	0.00001	-5	-11.5

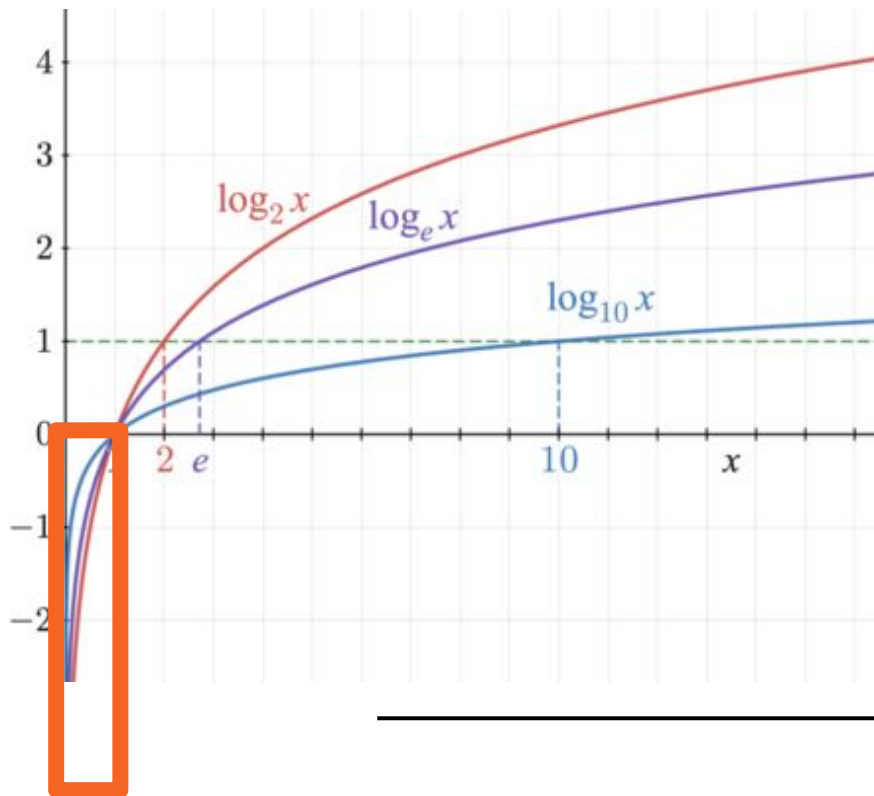
Logs are much more understandable when probabilities are small!

Log Probability



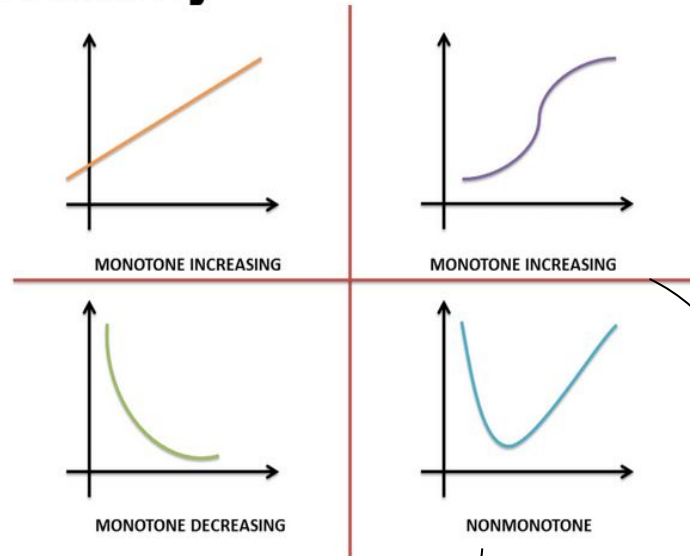
- Before, we've talked about log transforms as “squishing” big numbers

Log Probability



- Before, we've talked about log transforms as “squishing” big numbers
- But now we're constrained to $\log(P(E))$ where $0 < P(E) < 1$, where we're actually *not* squishing numbers

Monotonicity



<https://discover.hubpages.com/education/Differences-Between-Spearman-Rank-Correlation-and-Pearson-Correlation> 30

Log Probability

- Is the logarithm function (base e)
 - A monotonic function?
 - Increasing? Decreasing?

Log Probability

- Is the logarithm function (base e)
 - A monotonic function? **Yes**
 - Increasing? Decreasing? **Increasing**

Log Probability

- Why is it convenient that \ln (log base e) is a monotonic function?
 - The input point at which $a*b$ is maximized is the same input point at which $\log(a*b)$ is maximized

Log Probability

- Why is it convenient that \ln (log base e) is a monotonic function?
 - The input point at which $a*b$ is maximized is the same input point at which $\log(a*b)$ is maximized
 - This is called the **argmax** because we're sorting by a *function* of the point, not the point itself

Log Probability

- The key is in what we derived before:

$$\log(P(E) \cdot P(F)) = \log P(E) + \log P(F)$$

Log Probability

- The key is in what we derived before:

$$\log(P(E) \cdot P(F)) = \log P(E) + \log P(F)$$

$$\log \prod_i P(E_i) = \sum_i \log P(E_i)$$

Log Probability

- The key is in what we derived before:

What happens if you multiply a bunch of tiny numbers close to 0?

$$\log \prod_i P(E_i) = \sum_i \log P(E_i)$$

Log Probability

- The key is in what we derived before:

What happens if you multiply a bunch of tiny numbers close to 0? Unless you have a really special computer, you'll get 0

$$\log \prod_i P(E_i) = \sum_i \log P(E_i)$$

Log Probability

- The key is in what we derived before:

The smallest decimal Python can express is 2.225e-308

$$\log \prod_i P(E_i) = \sum_i \log P(E_i)$$

Log Probability

- The key is in what we derived before:

Do we run into small number issues if we do sums instead?

$$\log \prod_i P(E_i) = \sum_i \log P(E_i)$$

Log Probability

- The key is in what we derived before:

No issues! $\log(2.225e-308) = -307.652$, easy to store in Python

$$\log \prod_i P(E_i) = \sum_i \log P(E_i)$$

Simulate small
probabilities

Log Probability

```
>>> import numpy as np
```

```
>>> r = np.random.random_sample((100,)) * 0.000001
```

Log Probability

```
>>> import numpy as np
```

```
>>> r = np.random.random_sample((100,)) * 0.000001
```

What will we get?

```
>>> np.product(r)
```

Log Probability

```
>>> import numpy as np
```

```
>>> r = np.random.random_sample((100,)) * 0.000001
```

```
>>> np.product(r)
```

```
0.0
```

Not the actual
product – too
small for Python!

Log Probability

```
>>> import numpy as np
```

```
>>> r = np.random.random_sample((100,)) * 0.000001
```

```
>>> np.product(r)
```

```
0.0
```

```
>>> log_r = [np.log(x) for x in r]
```

Let's try taking
logs now

Log Probability

```
>>> import numpy as np
```

```
>>> r = np.random.random_sample((100,)) * 0.000001
```

```
>>> np.product(r)
```

```
0.0
```

```
>>> log_r = [np.log(x) for x in r]
```

```
>>> np.sum(log_r)
```

Do we expect a
weird number now?

Log Probability

```
>>> import numpy as np
```

```
>>> r = np.random.random_sample((100,)) * 0.000001
```

```
>>> np.product(r)
```

```
0.0
```

```
>>> log_r = [np.log(x) for x in r]
```

```
>>> np.sum(log_r)
```

```
-1472.245511811776
```

No Python
issues here!

Log Probability Takeaways

- If you're using a computer to help you calculate the product of probabilities, you should nearly always just use log probabilities instead
 - They are harder for humans to interpret...
 - but will allow for more accurate computation
 - and allow for finding the same maximizing points due to monotonicity

1 min break (stare at this table)

One in	Probability	Log_{10}	Log_e
10	0.1	-1	-2.3
100	0.01	-2	-4.6
1,000	0.001	-3	-6.9
10,000	0.0001	-4	-9.2
100,000	0.00001	-5	-11.5

Log prob quiz!

What probability has $\log_e -2.3$?

One in _____

What is the \log_e of **One in 100,000**?

Log prob quiz!

What probability has $\log_e -2.3$?

One in 10

What is the \log_e of One in 100,000?


-11.5

Classifying text: sports or not?

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

Classifying text: sports or not?

Training data



Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

Classifying text: sports or not?

Training data

Want to
classify new
data

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports
"A very close game"	

Classifying text: sports or not?

In math: $P(\text{____} | \text{____})?$

Want to
classify new
data



Text	Tag
"A very close game"	Sports / Not Sports?

Classifying text: sports or not?

In math: $P(\text{Sports} \mid \text{"A very close game"})$

Want to
classify new
data



Text	Tag
"A very close game"	Sports / Not Sports?

Classifying text: sports or not?

In math: $P(\text{Sports} \mid \text{"A very close game"})$

Notice: we aren't outputting a binary estimate of Sports or Not Sports – we're outputting an estimate of the probability that our tag is Sports

Want to
classify new
data



Text	Tag
"A very close game"	Sports / Not Sports?

Classifying text: sports or not?

In math: $P(\text{Sports} \mid \text{"A very close game"})$


"Sports" here is a binary variable: it can take either value 0 or 1. We can, for example, write a subset probability: $P(\text{Sports}=1 \mid \text{"A very close game"})$

Want to
classify new
data



Text	Tag
"A very close game"	Sports / Not Sports?

Classifying text: sports or not?



Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

How do we use **training data** to calculate $P(\text{Sports} \mid \text{"A very close game"})$?

Classifying text: sports or not?

First, let's apply Bayes' theorem:

$$P(\text{Sports} \mid \text{"A very close game"}) = \frac{P(\text{"A very close game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A very close game"})}$$

Classifying text: sports or not?

First, let's apply Bayes' theorem:

$$P(\text{Sports} \mid \text{"A very close game"}) = \frac{P(\text{"A very close game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A very close game"})}$$

Remember we're being naive today... let's
assume all the words are *independent* (Yoda)

Classifying text: sports or not?

Assuming words are independent:

$$P(\text{Sports} \mid \text{"A", "very", "close", "game"}) = \frac{P(\text{"A", "very", "close", "game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

Classifying text: sports or not?

$$P(\text{Sports} \mid \text{"A", "very", "close", game"}) = \frac{P(\text{"A", "very", "close", "game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

Classifying text: sports or not?

We can reduce further using the “chain rule”

$$P(\text{Sports} \mid \text{“A”, “very”, “close”, game}) = \frac{P(\text{“A”, “very”, “close”, “game”} \mid \text{Sports}) * P(\text{Sports})}{P(\text{“A”, “very”, “close”, “game”})}$$

$$P(y \mid x_1, \dots, x_n) = \frac{P(x_1 \mid y) P(x_2 \mid y) \dots P(x_n \mid y) P(y)}{P(x_1) P(x_2) \dots P(x_n)}$$

Classifying text: sports or not?

$P(\text{Sports} \mid \text{"A very close game"}) =$

$$\frac{P(\text{"A"} \mid \text{Sports}) * P(\text{"very"} \mid \text{Sports}) * P(\text{"close"} \mid \text{Sports}) * P(\text{"game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

Classifying text: sports or not?

$P(\text{Sports} \mid \text{"A very close game"}) =$

$$\frac{P(\text{"A"} \mid \text{Sports}) * P(\text{"very"} \mid \text{Sports}) * P(\text{"close"} \mid \text{Sports}) * P(\text{"game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

We can reduce this slightly more!

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

These are the same; big Pi means multiplying (the same way big Sigma means summing)

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

alpha here means “proportional to” (e.g., instead “=” which means “equals to”)

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

This is just the numerator of the above expression! What happened to the denominator?

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

The goal of calculating $P(y|x_1, \dots, x_n)$ is to find the argmax of it!

$$y = \operatorname{argmax}_y P(y) \prod_{i=1}^n P(x_i|y)$$

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

$$y = \operatorname{argmax}_y P(y) \prod_{i=1}^n P(x_i|y)$$

Classifying text: sports or not?

$$P(y|x_1, \dots, x_n) = \frac{P(x_1|y)P(x_2|y)\dots P(x_n|y)P(y)}{P(x_1)P(x_2)\dots P(x_n)}$$

$$P(y|x_1, \dots, x_n) \propto P(y) \prod_{i=1}^n P(x_i|y)$$

It turns out the denominator isn't affected by y , so it'll just be a constant that's irrelevant for calculating the argmax!

$$y = \operatorname{argmax}_y P(y) \prod_{i=1}^n P(x_i|y)$$

Classifying text: sports or not?

What does this mean in words? Let's switch back to our example...

$$y = \operatorname{argmax}_y P(y) \prod_{i=1}^n P(x_i|y)$$

Classifying text: sports or not?

$P(\text{Sports} \mid \text{"A very close game"}) =$

$$\frac{P(\text{"A"} \mid \text{Sports}) * P(\text{"very"} \mid \text{Sports}) * P(\text{"close"} \mid \text{Sports}) * P(\text{"game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

Classifying text: sports or not?

$P(\text{Sports} \mid \text{"A very close game"}) =$

$$\frac{P(\text{"A"} \mid \text{Sports}) * P(\text{"very"} \mid \text{Sports}) * P(\text{"close"} \mid \text{Sports}) * P(\text{"game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

We have $y = \text{Sports}$, and $x = \text{"A very close game"}$

We want to pick the value of y ($\text{Sports} = 1$ or $\text{Sports} = 0$) such that the probability above is maximized.

Classifying text: sports or not?

$P(\text{Sports} \mid \text{"A very close game"}) =$

$$\frac{P(\text{"A"} \mid \text{Sports}) * P(\text{"very"} \mid \text{Sports}) * P(\text{"close"} \mid \text{Sports}) * P(\text{"game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

$P(\text{"A", "very", "close", "game"})$ is *the same* for Sports and NotSports.

We want to pick the value of y (Sports = 1 or Sports = 0) such that the probability above is maximized.

Classifying text: sports or not?

$P(\text{Sports} \mid \text{"A very close game"}) =$

$$\frac{P(\text{"A"} \mid \text{Sports}) * P(\text{"very"} \mid \text{Sports}) * P(\text{"close"} \mid \text{Sports}) * P(\text{"game"} \mid \text{Sports}) * P(\text{Sports})}{P(\text{"A", "very", "close", "game"})}$$

$\alpha P(\text{"A"} \mid \text{Sports}) * P(\text{"very"} \mid \text{Sports}) * P(\text{"close"} \mid \text{Sports}) * P(\text{"game"} \mid \text{Sports}) * P(\text{Sports})$
alpha-looking thing means “proportional to”

Classifying text: sports or not?

To calculate the probability $P(\text{Sports} \mid \text{"A very close game"})$, we need the following ingredients:

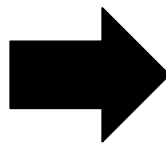
- $P(\text{"A"} \mid \text{Sports})$
- $P(\text{"very"} \mid \text{Sports})$
- $P(\text{"close"} \mid \text{Sports})$
- $P(\text{"game"} \mid \text{Sports})$
- $P(\text{Sports})$

How do we get $P(\text{"word"} \mid \text{Sports})$??

Generate word frequencies!

Training Data

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports



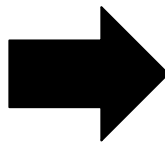
Joint Count Data

Word	# in Sports	# in Not sports
"a"		
"very"		
"close"		
"game"		

Generate word frequencies!

Training Data

Text	Tag
" A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
" A clean but forgettable game"	Sports
"It was a close election"	Not sports



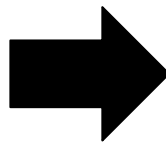
Joint Count Data

Word	# in Sports	# in Not sports
"a"	2	1
"very"		
"close"		
"game"		

Generate word frequencies!

Training Data

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
" Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports



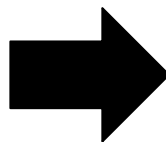
Joint Count Data

Word	# in Sports	# in Not sports
"a"	2	1
"very"	1	0
"close"		
"game"		

Generate word frequencies!

Training Data

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports



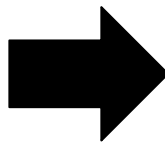
Joint Count Data

Word	# in Sports	# in Not sports
"a"	2	1
"very"	1	0
"close"	?	?
"game"		

Generate word frequencies!

Training Data

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports



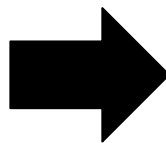
Joint Count Data

Word	# in Sports	# in Not sports
"a"	2	1
"very"	1	0
"close"	0	1
"game"		

Generate word frequencies!

Training Data

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports



Joint Count Data

Word	# in Sports	# in Not sports
"a"	2	1
"very"	1	0
"close"	0	1
"game"	2	0

Classifying text: sports or not?

For $P(\text{Sports} = 1 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Sports})$
- $P(\text{"very"} \mid \text{Sports})$
- $P(\text{"close"} \mid \text{Sports})$
- $P(\text{"game"} \mid \text{Sports})$
- $P(\text{Sports})$

For $P(\text{Sports} = 0 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Not sports})$
- $P(\text{"very"} \mid \text{Not sports})$
- $P(\text{"close"} \mid \text{Not sports})$
- $P(\text{"game"} \mid \text{Not sports})$
- $P(\text{Not sports})$

Classifying text: sports or not?

For $P(\text{Sports} = 1 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Sports})$
- $P(\text{"very"} \mid \text{Sports})$
- $P(\text{"close"} \mid \text{Sports})$
- $P(\text{"game"} \mid \text{Sports})$
- $P(\text{Sports})$

For $P(\text{Sports} = 0 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Not sports})$
- $P(\text{"very"} \mid \text{Not sports})$
- $P(\text{"close"} \mid \text{Not sports})$
- $P(\text{"game"} \mid \text{Not sports})$
- $P(\text{Not sports})$

Goal: compare which probability is bigger!

Classifying text: sports or not?

For $P(\text{Sports} = 1 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Sports})$
- $P(\text{"very"} \mid \text{Sports})$
- $P(\text{"close"} \mid \text{Sports})$
- $P(\text{"game"} \mid \text{Sports})$
- $P(\text{Sports})$

For $P(\text{Sports} = 0 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Not sports})$
- $P(\text{"very"} \mid \text{Not sports})$
- $P(\text{"close"} \mid \text{Not sports})$
- $P(\text{"game"} \mid \text{Not sports})$
- $P(\text{Not sports})$

Our joint count table gave us these numbers

Classifying text: sports or not?

For $P(\text{Sports} = 1 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Sports})$
- $P(\text{"very"} \mid \text{Sports})$
- $P(\text{"close"} \mid \text{Sports})$
- $P(\text{"game"} \mid \text{Sports})$
- $P(\text{Sports})$

For $P(\text{Sports} = 0 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Not sports})$
- $P(\text{"very"} \mid \text{Not sports})$
- $P(\text{"close"} \mid \text{Not sports})$
- $P(\text{"game"} \mid \text{Not sports})$
- $P(\text{Not sports})$

But we still don't have these probabilities!

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

$P(\text{Sports}) = ?$

$P(\text{Not sports}) = ?$

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

$$P(\text{Sports}) = 3/5$$

$$P(\text{Not sports}) = 2/5$$

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

$$P(\text{Sports}) = 3/5$$

$$P(\text{Not sports}) = 2/5$$

$$\# \text{ words in Sports} = 11$$

$$\# \text{ words in Not sports} = 9$$

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

$$P(\text{Sports}) = 3/5$$

$$P(\text{Not sports}) = 2/5$$

words in Sports = 11

words in Not sports = 9

Word	# in Sports	# in Not sports	P(Word Sports)	P(Word Not sports)
"a"	2	1	?	?
"very"	1	0	?	?
"close"	0	1	?	?
"game"	2	0	?	?

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

$$P(\text{Sports}) = 3/5$$

$$P(\text{Not sports}) = 2/5$$

$$\# \text{ words in Sports} = 11$$

$$\# \text{ words in Not sports} = 9$$

Word	# in Sports	# in Not sports	P(Word Sports)	P(Word Not sports)
"a"	2	1	2/11	1/9
"very"	1	0	1/11	0/9
"close"	0	1	0/11	1/9
"game"	2	0	2/11	0/9

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

We want to multiply...

- $P(\text{"A"}|\text{Sports}) = 2/11$
- $P(\text{"very"}|\text{Sports}) = 1/11$
- $P(\text{"close"}|\text{Sports}) = 0/11$
- $P(\text{"game"}|\text{Sports}) = 2/11$
- $P(\text{Sports}) = 3/5$

Word	# in Sports	# in Not sports	$P(\text{Word} \text{Sports})$	$P(\text{Word} \text{Not sports})$
"a"	2	1	2/11	1/9
"very"	1	0	1/11	0/9
"close"	0	1	0/11	1/9
"game"	2	0	2/11	0/9

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

We'll just get a 0 probability since "close" never appeared in the Sports=1 training data!

That's not fair, so we have to correct for it somehow...

Word	# in Sports	# in Not sports	P(Word Sports)	P(Word Not sports)
"a"	2	1	2/11	1/9
"very"	1	0	1/11	0/9
"close"	0	1	0/11	1/9
"game"	2	0	2/11	0/9

Laplace Smoothing: add 1 to every count so nothing is ever zero!

Word	# in Sports	# in Not sports	P(Word Sports)	P(Word Not sports)
“a”	2	1	2/11	1/9
“very”	1	0	1/11	0/9
“close”	0	1	0/11	1/9
“game”	2	0	2/11	0/9

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

With **Laplace Smoothing**:

Every probability *numerator* increases by 1

Word	# in Sports	# in Not sports	P(Word Sports)	P(Word Not sports)
"a"	2	1	3/25	2/23
"very"	1	0	2/25	1/23
"close"	0	1	1/25	2/23
"game"	2	0	3/25	1/23

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

With **Laplace Smoothing**:

Every probability *denominator* increases by the *distinct word count* across all the training data

Word	# in Sports	# in Not sports	P(Word Sports)	P(Word Not sports)
"a"	2	1	3/25	2/23
"very"	1	0	2/25	1/23
"close"	0	1	1/25	2/23
"game"	2	0	3/25	1/23

Text	Tag
"A great game"	Sports
"The election was over"	Not sports
"Very clean match"	Sports
"A clean but forgettable game"	Sports
"It was a close election"	Not sports

With **Laplace Smoothing**:

Distinct word count overall = 14
 # words in Sports = $11 + 14 = 25$
 # words in Not sports = $9 + 14 = 23$

Word	# in Sports	# in Not sports	P(Word Sports)	P(Word Not sports)
"a"	2	1	$3/25$	$2/23$
"very"	1	0	$2/25$	$1/23$
"close"	0	1	$1/25$	$2/23$
"game"	2	0	$3/25$	$1/23$

$$P(\text{Sports}) = \frac{3}{5}$$

$$P(\text{Not sports}) = \frac{2}{5}$$

Word	P(Word Sports)	P(Word Not sports)
"a"	3/25	2/23
"very"	2/25	1/23
"close"	1/25	2/23
"game"	3/25	1/23

$P(\text{Sports}=1 \mid \text{"A very close game"})$ is the multiplication of:

- $P(\text{"A"}|\text{Sports}) = 3/25$
- $P(\text{"very"}|\text{Sports}) = 2/25$
- $P(\text{"close"}|\text{Sports}) = 1/25$
- $P(\text{"game"}|\text{Sports}) = 3/25$
- $P(\text{Sports}) = 3/5$

$$P(\text{Sports}) = \frac{3}{5}$$

$$P(\text{Not sports}) = \frac{2}{5}$$

Word	P(Word Sports)	P(Word Not sports)
"a"	3/25	2/23
"very"	2/25	1/23
"close"	1/25	2/23
"game"	3/25	1/23

$P(\text{Sports}=1 \mid \text{"A very close game"})$ is the multiplication of:

- $P(\text{"A"}|\text{Sports}) = 3/25$
- $P(\text{"very"}|\text{Sports}) = 2/25$
- $P(\text{"close"}|\text{Sports}) = 1/25$
- $P(\text{"game"}|\text{Sports}) = 3/25$
- $P(\text{Sports}) = 3/5$

We got lucky this time because this is calculable in Python still...
0.000027648

How else can we deal with this P?

$$P(\text{Sports}) = \frac{3}{5}$$

$$P(\text{Not sports}) = \frac{2}{5}$$

Word	P(Word Sports)	P(Word Not sports)
“a”	3/25	2/23
“very”	2/25	1/23
“close”	1/25	2/23
“game”	3/25	1/23

$\text{Log}[P(\text{Sports}=1 \mid \text{“A very close game”})]$ is the sum of:

- $\text{Log}[P(\text{“A”}|\text{Sports})] = \text{log}(3/25)$
- $\text{Log}[P(\text{“very”}|\text{Sports})] = \text{log}(2/25)$
- $\text{Log}[P(\text{“close”}|\text{Sports})] = \text{log}(1/25)$
- $\text{Log}[P(\text{“game”}|\text{Sports})] = \text{log}(3/25)$
- $\text{Log}[P(\text{Sports})] = \text{log}(3/5)$

This gives us -10.5

Log prob quiz!

What probability has $\log_e -2.3$?

One in _____

What is the \log_e of **One in 100,000**?

Log prob quiz!

What probability has $\log_e -2.3$?

One in 10

What is the \log_e of One in 100,000?

-11.5

$$P(\text{Sports}) = \frac{3}{5}$$

$$P(\text{Not sports}) = \frac{2}{5}$$

Word	P(Word Sports)	P(Word Not sports)
"a"	3/25	2/23
"very"	2/25	1/23
"close"	1/25	2/23
"game"	3/25	1/23

$\text{Log}[P(\text{Sports}=1 \mid \text{"A very close game"})]$ is the sum of:

- $\text{Log}[P(\text{"A"}|\text{Sports})] = \text{log}(3/25)$
- $\text{Log}[P(\text{"very"}|\text{Sports})] = \text{log}(2/25)$
- $\text{Log}[P(\text{"close"}|\text{Sports})] = \text{log}(1/25)$
- $\text{Log}[P(\text{"game"}|\text{Sports})] = \text{log}(3/25)$
- $\text{Log}[P(\text{Sports})] = \text{log}(3/5)$

This gives us -10.5

Classifying text: sports or not?

For $P(\text{Sports} = 1 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Sports})$
- $P(\text{"very"} \mid \text{Sports})$
- $P(\text{"close"} \mid \text{Sports})$
- $P(\text{"game"} \mid \text{Sports})$
- $P(\text{Sports})$

For $P(\text{Sports} = 0 \mid \text{"A very close game"})$, multiply:

- $P(\text{"A"} \mid \text{Not sports})$
- $P(\text{"very"} \mid \text{Not sports})$
- $P(\text{"close"} \mid \text{Not sports})$
- $P(\text{"game"} \mid \text{Not sports})$
- $P(\text{Not sports})$

Goal: compare which probability is bigger!

Classifying text: sports or not?

$P(\text{Sports} = 1 \mid \text{"A very close game"})$

$\propto 0.000027648$

⋮

$P(\text{Sports} = 0 \mid \text{"A very close game"})$

$\propto 0.00000572$

Goal: compare which probability is bigger!

Classifying text: sports or not?

$P(\text{Sports} = 1 \mid \text{"A very close game"})$

$\propto 0.000027648$

⋮

$P(\text{Sports} = 0 \mid \text{"A very close game"})$

$\propto 0.00000572$

Normalize by $0.000027648 + 0.00000572 = 0.000033368$

$= 0.000027648 / 0.000033368 \approx 83\%$

⋮

$= 0.00000572 / 0.000033368 \approx 17\%$

If we want to compare this without 'proportional to'...

Classifying text: sports or not?

$P(\text{Sports} = 1 \mid \text{"A very close game"})$

$\propto 0.000027648$

$\text{Log}[P(\text{Sports} = 1 \mid \text{"A very close game"})]$

$= -10.5$

$P(\text{Sports} = 0 \mid \text{"A very close game"})$

$\propto 0.00000572$

$\text{Log}[P(\text{Sports} = 0 \mid \text{"A very close game"})]$

$= -12.07$

Goal: compare which probability is bigger!

Classifying text: sports or not?

$P(\text{Sports} = 1 \mid \text{"A very close game"})$

$\propto 0.000027648$

$\text{Log}[P(\text{Sports} = 1 \mid \text{"A very close game"})]$

$= -10.5$ Better than 1 in 100k

$P(\text{Sports} = 0 \mid \text{"A very close game"})$

$\propto 0.00000572$

$\text{Log}[P(\text{Sports} = 0 \mid \text{"A very close game"})]$

$= -12.07$

Goal: compare which probability is bigger!

Classifying text: sports or not?

$P(\text{Sports} = 1 \mid \text{"A very close game"})$

$\propto 0.000027648$

$\text{Log}[P(\text{Sports} = 1 \mid \text{"A very close game"})]$

$= -10.5$

$P(\text{Sports} = 0 \mid \text{"A very close game"})$

$\propto 0.00000572$

$\text{Log}[P(\text{Sports} = 0 \mid \text{"A very close game"})]$

$= -12.07$ Worse than 1 in 100k

Goal: compare which probability is bigger!

Classifying text: sports or not?

$P(\text{Sports} = 1 \mid \text{"A very close game"})$

$= 0.000027648$

$\text{Log}[P(\text{Sports} = 1 \mid \text{"A very close game"})]$

$= -10.5$

$P(\text{Sports} = 0 \mid \text{"A very close game"})$

$= 0.00000572$

$\text{Log}[P(\text{Sports} = 0 \mid \text{"A very close game"})]$

$= -12.07$

It's more likely that "A very close game" is about Sports!

Probability ratios

$$\log \left(\frac{P(Sports = 1|text)}{P(Sports = 0|text)} \right) =$$

$$\log P(Sports = 1|text) - \log P(Sports = 0|text)$$

Probability ratios

$$\log \left(\frac{P(Sports = 1|text)}{P(Sports = 0|text)} \right) =$$

$$\begin{aligned} & \log P(Sports = 1|text) - \log P(Sports = 0|text) \\ & = -10.5 - (-12.07) = 1.57 \end{aligned}$$

Probability ratios

$$\log \left(\frac{P(\textit{Sports} = 1 | \textit{text})}{P(\textit{Sports} = 0 | \textit{text})} \right) =$$

$$\begin{aligned} & \log P(\textit{Sports} = 1 | \textit{text}) - \log P(\textit{Sports} = 0 | \textit{text}) \\ &= -10.5 - (-12.07) = 1.57 \end{aligned}$$

Exponentiating gives us how many times more likely this text is Sports is vs. Not sports

Probability ratios

$$\log \left(\frac{P(\textit{Sports} = 1 | \textit{text})}{P(\textit{Sports} = 0 | \textit{text})} \right) =$$

$$\begin{aligned} & \log P(\textit{Sports} = 1 | \textit{text}) - \log P(\textit{Sports} = 0 | \textit{text}) \\ &= -10.5 - (-12.07) = 1.57 \end{aligned}$$

Exponentiating gives us how many times more likely this text is Sports is vs. Not sports: $e^{1.57} = 4.8$ times more likely

Classifying text: sports or not?

- This process was called **Naive Bayes**
 - Get frequencies of your data units (e.g. word)
 - Calculate probabilities assuming independence of data units
 - Use Bayes' rule to determine what classification has the highest probability

Naive Bayes

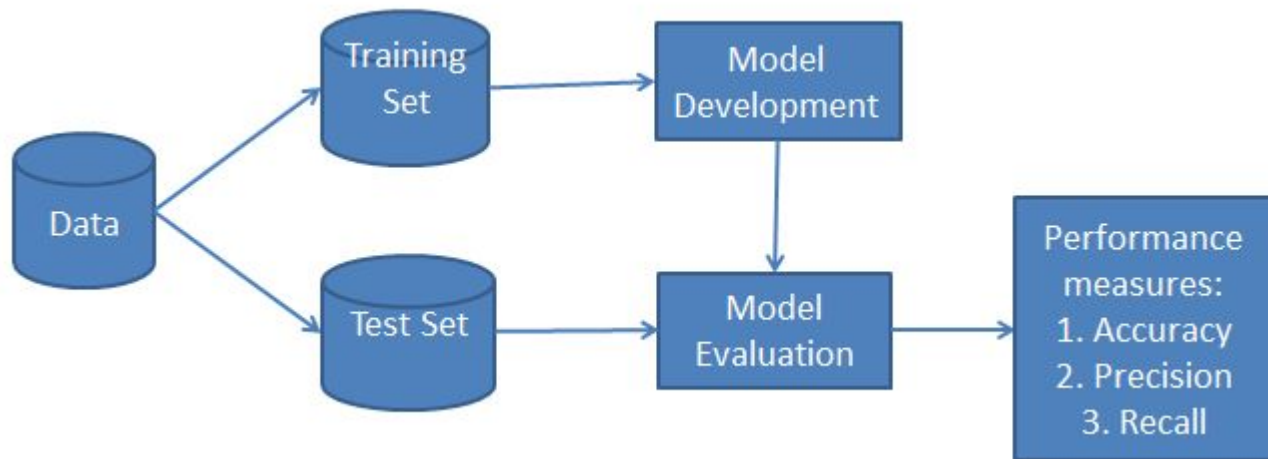
- **Naive Bayes allows us to classify our outcomes**
 - These can be binary (Sports / not sports) or multi-category (sport A / sport B / elections)
 - We calculate probabilities based on the frequencies of these categories in our data for each independent input x (whether x_i 's are words or other df inputs)

Gaussian Naive Bayes

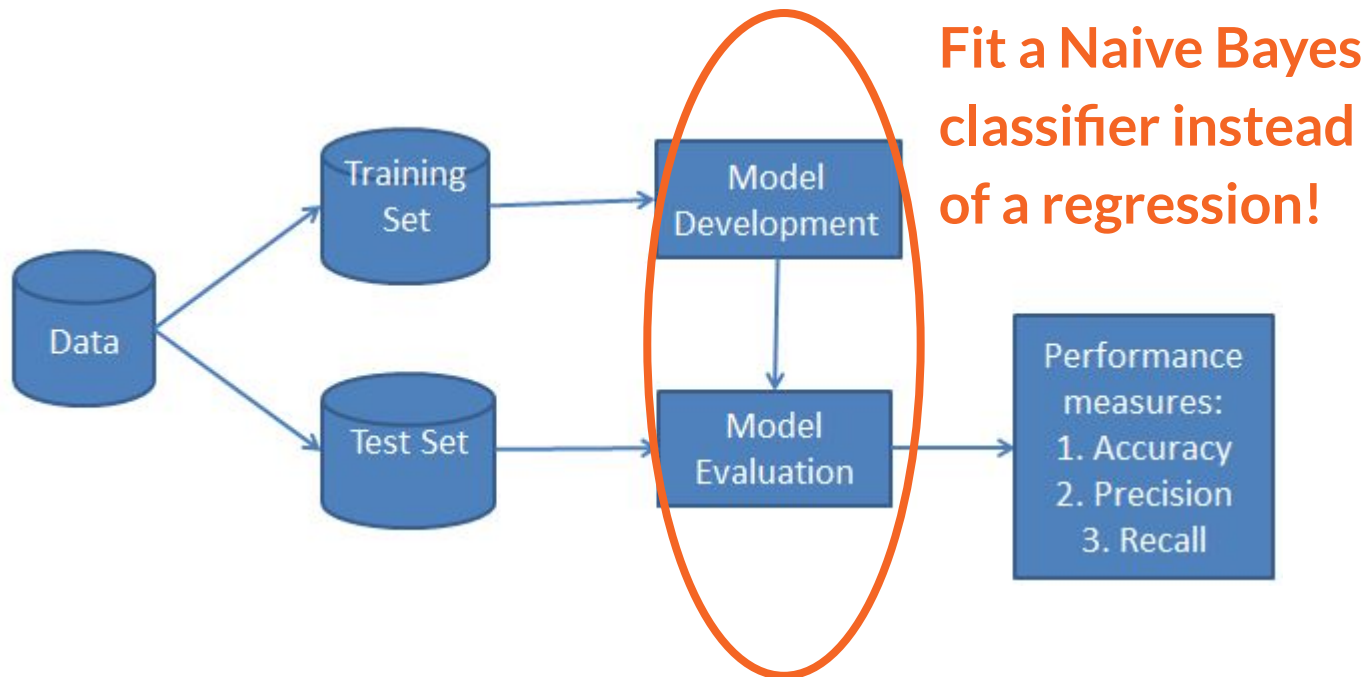
- Naive Bayes allows us to classify our outcomes
 - These can be binary (Sports / not sports) or multi-category (sport A / sport B / elections)

○ We calculate probabilities based on the normal probability density function ~~frequencies of these categories~~ in our data for each independent input x (whether x_i 's are words or other df inputs)

Naive Bayes... as a model?



Naive Bayes... as a model?



Naive Bayes using scikit learn

```
from sklearn.naive_bayes import GaussianNB

model = GaussianNB()

model.fit(X_train,Y_train)

model.predict(X_test)
```