

Bayes' Theorem



Remember Language ID?

- Let $p(X)$ = probability of text X in English
- Let $q(X)$ = probability of text X in Polish
- Which probability is higher?

– (we'd also like bias toward English since it's more likely *a priori* – ignore that for now)

Let's revisit this

"Horses and Lukasiewicz are on the curriculum."

$p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots)$

Bayes' Theorem



- $p(A | B) = p(B | A) * p(A) / p(B)$
- Easy to check by removing syntactic sugar
- **Use 1:** Converts $p(B | A)$ to $p(A | B)$
- **Use 2:** Updates $p(A)$ to $p(A | B)$
- Stare at it so you'll recognize it later

Language ID

- Given a sentence x , I suggested comparing its prob in different languages:
 - $p(\text{SENT}=x \mid \text{LANG}=\text{english})$ (i.e., $p_{\text{english}}(\text{SENT}=x)$)
 - $p(\text{SENT}=x \mid \text{LANG}=\text{polish})$ (i.e., $p_{\text{polish}}(\text{SENT}=x)$)
 - $p(\text{SENT}=x \mid \text{LANG}=\text{xhosa})$ (i.e., $p_{\text{xhosa}}(\text{SENT}=x)$)
- But surely for language ID we should compare
 - $p(\text{LANG}=\text{english} \mid \text{SENT}=x)$
 - $p(\text{LANG}=\text{polish} \mid \text{SENT}=x)$
 - $p(\text{LANG}=\text{xhosa} \mid \text{SENT}=x)$

Language ID

- For language ID we should compare

- $p(\text{LANG}=\text{english} \mid \text{SENT}=x)$
 - $p(\text{LANG}=\text{polish} \mid \text{SENT}=x)$
 - $p(\text{LANG}=\text{xhosa} \mid \text{SENT}=x)$
- a posteriori*

- For ease, multiply by $p(\text{SENT}=x)$ and compare

- $p(\text{LANG}=\text{english}, \text{SENT}=x)$
 - $p(\text{LANG}=\text{polish}, \text{SENT}=x)$
 - $p(\text{LANG}=\text{xhosa}, \text{SENT}=x)$
- sum of these is a way to find $p(\text{SENT}=x)$; can divide back by that to get posterior probs

- Must know prior probabilities; then rewrite as

- $p(\text{LANG}=\text{english})$
 - $p(\text{LANG}=\text{polish})$
 - $p(\text{LANG}=\text{xhosa})$
- a priori*
- * $p(\text{SENT}=x \mid \text{LANG}=\text{english})$
 - * $p(\text{SENT}=x \mid \text{LANG}=\text{polish})$
 - * $p(\text{SENT}=x \mid \text{LANG}=\text{xhosa})$
- likelihood (what we had before)*

Let's try it!

"First we pick a random LANG,
then we roll a random SENT
with the LANG dice."

best	0.7	$p(\text{LANG}=\text{english})$	*	$p(\text{SENT}=x \mid \text{LANG}=\text{english})$	0.00001
	0.2	$p(\text{LANG}=\text{polish})$	*	$p(\text{SENT}=x \mid \text{LANG}=\text{polish})$	0.00004
	0.1	$p(\text{LANG}=\text{xhosa})$	*	$p(\text{SENT}=x \mid \text{LANG}=\text{xhosa})$	0.00005 best
	↑	prior prob		likelihood	↑

from a very simple
model: a single die
whose sides are the
languages of the world

from a set of
trigram dice
(actually 3 sets,
one per language)

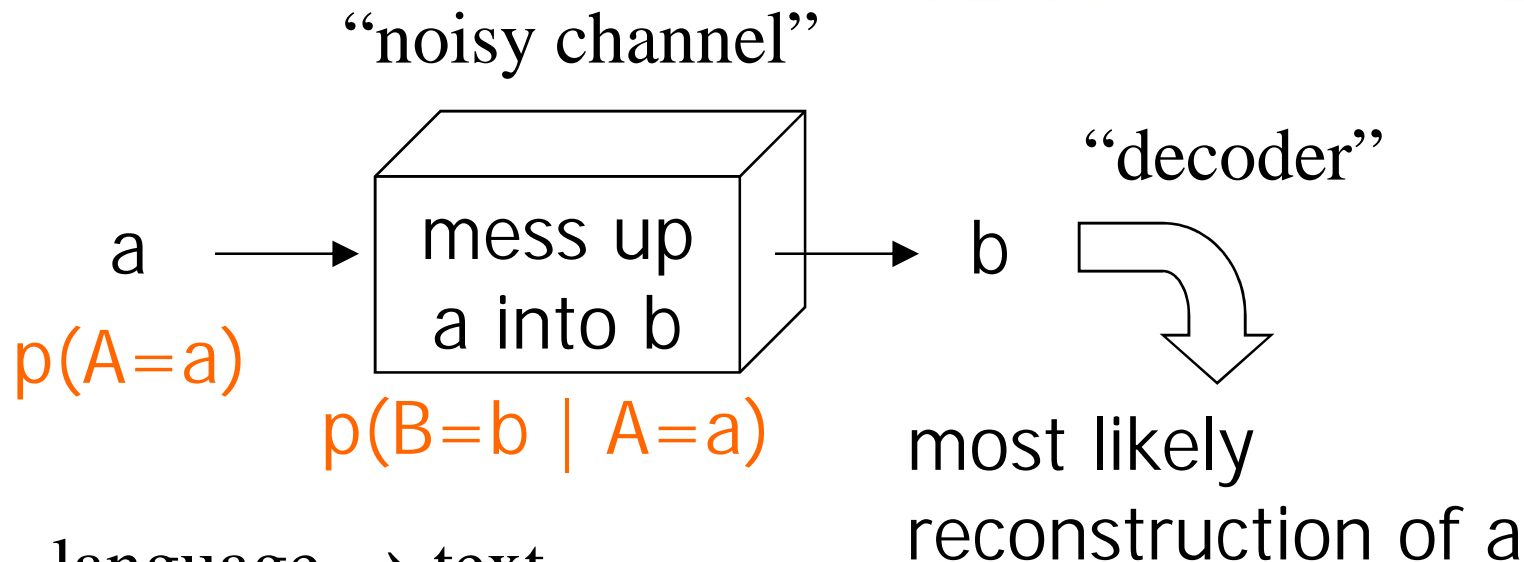
=	$p(\text{LANG}=\text{english}, \text{SENT}=x)$	0.000007	
=	$p(\text{LANG}=\text{polish}, \text{SENT}=x)$	0.000008	best compromise
=	$p(\text{LANG}=\text{xhosa}, \text{SENT}=x)$	0.000005	
	joint probability		
<hr/>			
	$p(\text{SENT}=x)$	0.000020	total over all ways of getting $\text{SENT}=x$
	probability of evidence		

Let's try it!

"First we pick a random LANG,
then we roll a random SENT
with the LANG dice."

	=	$p(\text{LANG}=\text{english}, \text{SENT}=x)$	0.000007	
...	=	$p(\text{LANG}=\text{polish}, \text{SENT}=x)$	0.000008	best compromise
	=	$p(\text{LANG}=\text{xhosa}, \text{SENT}=x)$	0.000005	
		joint probability		
add up		$p(\text{SENT}=x)$	0.000020	total probability of getting $\text{SENT}=x$ one way or another!
normalize (divide by a constant so they'll sum to 1)		$p(\text{LANG}=\text{english} \mid \text{SENT}=x)$	$0.000007/0.000020 = 7/20$	
		$p(\text{LANG}=\text{polish} \mid \text{SENT}=x)$	$0.000008/0.000020 = 8/20$	best
		$p(\text{LANG}=\text{xhosa} \mid \text{SENT}=x)$	$0.000005/0.000020 = 5/20$	
		posterior probability		given the evidence $\text{SENT}=x$, the possible languages sum to 1

General Case (“noisy channel”)



language → text
text → speech
spelled → misspelled
English → French

maximize $p(A=a | B=b)$

$$= p(A=a) p(B=b | A=a) / p(B=b)$$
$$= p(A=a) p(B=b | A=a) / \sum_{a'} p(A=a') p(B=b | A=a')$$

Language ID

- For language ID we should compare

- $p(\text{LANG}=\text{english} \mid \text{SENT}=x)$
- $p(\text{LANG}=\text{polish} \mid \text{SENT}=x)$
- $p(\text{LANG}=\text{xhosa} \mid \text{SENT}=x)$

a posteriori

- For ease, multiply by $p(\text{SENT}=x)$ and compare

- $p(\text{LANG}=\text{english}, \text{SENT}=x)$
- $p(\text{LANG}=\text{polish}, \text{SENT}=x)$
- $p(\text{LANG}=\text{xhosa}, \text{SENT}=x)$

- which we find as follows (we need prior probs!):

- $p(\text{LANG}=\text{english})$
- $p(\text{LANG}=\text{polish})$
- $p(\text{LANG}=\text{xhosa})$

a priori

- * $p(\text{SENT}=x \mid \text{LANG}=\text{english})$
- * $p(\text{SENT}=x \mid \text{LANG}=\text{polish})$
- * $p(\text{SENT}=x \mid \text{LANG}=\text{xhosa})$

likelihood

General Case (“noisy channel”)

- Want most likely A to have generated evidence B

- $p(A = a1 \mid B = b)$
- $p(A = a2 \mid B = b)$
- $p(A = a3 \mid B = b)$

a posteriori

- For ease, multiply by $p(B=b)$ and compare

- $p(A = a1, B = b)$
- $p(A = a2, B = b)$
- $p(A = a3, B = b)$

- which we find as follows (we need prior probs!):

- $p(A = a1)$
- $p(A = a2)$
- $p(A = a3)$

a priori

- * $p(B = b \mid A = a1)$
- * $p(B = b \mid A = a2)$
- * $p(B = b \mid A = a3)$

likelihood

Speech Recognition

- For baby speech recognition we should compare
 - $p(\text{MEANING}=\text{gimme} \mid \text{SOUND}=\text{uhh})$
 - $p(\text{MEANING}=\text{changeme} \mid \text{SOUND}=\text{uhh})$ *a posteriori*
 - $p(\text{MEANING}=\text{loveme} \mid \text{SOUND}=\text{uhh})$
- For ease, multiply by $p(\text{SOUND}=\text{uhh})$ & compare
 - $p(\text{MEANING}=\text{gimme}, \text{SOUND}=\text{uhh})$
 - $p(\text{MEANING}=\text{changeme}, \text{SOUND}=\text{uhh})$
 - $p(\text{MEANING}=\text{loveme}, \text{SOUND}=\text{uhh})$
- which we find as follows (we need prior probs!):

■ $p(\text{MEAN}=\text{gimme})$	*	$p(\text{SOUND}=\text{uhh} \mid \text{MEAN}=\text{gimme})$
■ $p(\text{MEAN}=\text{changeme})$	*	$p(\text{SOUND}=\text{uhh} \mid \text{MEAN}=\text{changeme})$
■ $p(\text{MEAN}=\text{loveme})$	*	$p(\text{SOUND}=\text{uhh} \mid \text{MEAN}=\text{loveme})$
<i>a priori</i>		<i>likelihood</i>

A simpler view? Odds Ratios

- What A values are probable, given that B=b?

- Bayes' Theorem says:

- $p(A=a1 | B=b) = p(A=a1) * p(B=b | A=a1) / p(B=b)$

- $p(A=a2 | B=b) = p(A=a2) * p(B=b | A=a2) / p(B=b)$

- Therefore

- $$\frac{p(A=a1 | B=b)}{p(A=a2 | B=b)} = \frac{p(A=a1)}{p(A=a2)} * \frac{p(B=b | A=a1)}{p(B=b | A=a2)}$$

$$\text{posterior odds ratio} = \text{prior odds ratio} * \text{likelihood ratio}$$

A simpler view? Odds Ratios

$$\frac{p(A=a1 | B=b)}{p(A=a2 | B=b)} = \frac{p(A=a1)}{p(A=a2)} * \frac{p(B=b | A=a1)}{p(B=b | A=a2)}$$

posterior odds ratio
prior odds ratio
likelihood odds ratio

0.7	p(LANG=english)	0.00001	
0.2	p(LANG=polish)	0.00004	
0.1	p(LANG=xhosa)	0.00005	
	prior	likelihood	

- A priori, English is 7 times as probable as Xhosa (7:1 odds)
- But likelihood of English is only 1/5 as large (1:5 odds)
- So a posteriori, English now $7 * 1/5 = 1.2$ times as probable (7:5 odds)
 - That is: $p(\text{English}) = 7/12$, $p(\text{Xhosa}) = 5/12$ if no other options

Growing evidence eventually overwhelms the prior

- We were expecting Polish text but actually it's English
- What happens as we read more & more words?
 - The prior odds ratio stays the same
 - But the likelihood odds ratio becomes extreme (much bigger or much smaller than 1, depending on which hypothesis is correct)
 - Suppose each trigram is 1.001 times more probable under the English model than the Polish model
 - Then after 700 trigrams, the likelihood ratio is > 2 in favor of English ($1.001^{700} > 2$)
 - And after 7000 trigrams, the likelihood ratio is $> 2^{10}$ in favor of English!
 - As long as the prior $p(\text{English}) > 0$, eventually we come to believe it's English a posteriori. We get surer and surer with more evidence.

Life or Death!

Does Epitaph have hoof-and-mouth disease?
He tested positive – oh no!
False positive rate only 5%

- $p(\text{hoof}) = 0.001$ so $p(\neg\text{hoof}) = 0.999$

- $p(\text{positive test} \mid \neg\text{hoof}) = 0.05$ “false pos”

- $p(\text{negative test} \mid \text{hoof}) = 0$ “false neg”

so $p(\text{positive test} \mid \text{hoof}) = 1 - 0 = 1$

- What is $p(\text{hoof} \mid \text{positive test})$?

- Consider the hoof: $\neg\text{hoof}$ odds ratio

- Prior odds ratio **1:999** (improbable!)

- Likelihood ratio at most $1:0.05$, or equivalently **20:1**

- So posterior odds ratio at most $20:999$, or about **1:50**

- That is, $p(\text{hoof} \mid \text{positive test})$ at most about $1/51$