

Goals of this lecture

- Probability notation like p(X | Y):
 - What does this expression mean?
 - How can I manipulate it?
 - How can I estimate its value in practice?
- Probability models:
 - What is one?
 - Can we build one for language ID?
 - How do I know if my model is any good?





What does that really mean? p(Paul Revere wins | weather's clear) = 0.9 Past performance? Revere's won 90% of races with clear weather Hypothetical performance? If he ran the race in many parallel universes ... Subjective strength of belief? Would pay up to 90 cents for chance to win \$1 Output of some computable formula? Ok, but then which formulas should we trust? p(X | Y) versus q(X | Y)









p measures total probability of a set of events.

Commas denote conjunction

p(Paul Revere wins, Valentine places, Epitaph shows | weather's clear)

Commas denote conjunction

- p(Paul Revere wins, Valentine places, Epitaph shows | weather's clear)
- p(Paul Revere wins | weather's clear, ground is dry, jockey getting over sprain, Epitaph also in race, Epitaph was recently bought by Gonzalez, race is on May 17, ...)
 - what happens as we add conjuncts to right of bar ?
 probability could increase or decrease
 probability gets more relevant to our case (less *bias*)
 probability *estimate* gets less reliable (more *variance*) # times Revere wins AND weather clear AND ... it's May 17 # times weather clear AND ... it's May 17

Simplifying Right Side: Backing Off

p(Paul Revere wins | weather's clear, ground is dry, jockey getting over sprain, Epitaph also in race, Epitaph was recently bought by Gonzalez, race is on May 17, ...)

not exactly what we want but at least we can get a reasonable estimate of it! (i.e., more bias but less variance) try to *keep* the conditions that we suspect will have the most influence on whether Paul Revere wins

Simplifying Right Side: Backing Off

p(Paul Revere wins, Valentine places, Epitaphshows | weather's clear)

Facto	oring Left Side: The Cha	in Rule		
p(Rever	re, Valentine, Epitaph weather's clear)	RVEW/W		
= p(Rever	re Valentine, Epitaph , weather's clear)	= RVEW/VEW		
^ *	p(Valentine Epitaph, weather's clear)	* VEW/EW		
/	* p(Epitaph weather's clear)	* EW/W		
	True because numerators cancel against denominators			
If this prob is unchanged by backoff, we say Revere was				
CONDITIONALLY INDEPENDENT of Valentine and Epitaph				
(conditioned on the weather's being clear). Often we just				
ASSUME conditional independence to get the nice product above.				

Remember Language ID?

- "Horses and Lukasiewicz are on the curriculum."
- Is this English or Polish or what?
- We had some notion of using n-gram models ...
- Is it "good" (= likely) English?
- Is it "good" (= likely) Polish?
- Space of events will be not races but character sequences (x₁, x₂, x₃, ...) where x_n = EOS

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)

"Horses and Lukasiewicz are on the curriculum." $p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, ...)$

Apply the Chain Rule		
$p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, x_6=s$)	
= p(x ₁ =h)	4470/	52108
* p(x ₂ =o x ₁ =h)	395/	4470
* p(x ₃ =r x ₁ =h, x ₂ =o)	5/	395
* $p(x_4=s x_1=h, x_2=o, x_3=r)$	3/	5
* $p(x_5=e x_1=h, x_2=o, x_3=r, x_4=s)$	3/	3
* p(x ₆ =s x ₁ =h, x ₂ =o, x ₃ =r, x ₄ =s, x ₅ =	e) 0/	3
* = 0		
	counts f	from
	Brown	corpus

Back Off On Right Side





Another Independence	e Assumption
p(x ₁ = h , x ₂ = o , x ₃ = r , x ₄ = s , x ₅ =	=e, x ₆ =s,)
$\approx p(x_1 = h)$	4470/52108
* p(x ₂ =o x ₁ =h)	395/ 4470
* p(x _i =r x _{i-2} =h, x _{i-1} =o)	1417/14765
* p(x _i =s x _{i-2} =0, x _{i-1} =r	() 1573/26412
* p(x _i =e x _{i-2} =r,	, x _{i-1} = s) 1610/12253
* p(x _i =s	x _{i-2} = s , x _{i-1} = 20) 4/21250
* = 5.4e-7 *	
	counts from
	Brown corpus
600.465 Intro to NI R Eispor	19

Simplify the Notation	
p (x ₁ = h , x ₂ = o , x ₃ = r , x ₄ = s , x ₅	5=e, x ₆ =s,)
$\approx p(x_1 = h)$	4470/52108
* p(x ₂ =o x ₁ =h)	395/ 4470
* p(r h, o)	1417/14765
* p(s o, r)	1573/26412
* p(e r, s)	1610/12253
* p(s s, e)	2044/21250
*	
	counts from
	Brown corpus











What is "X" in "p(X)"?

- Element of some implicit "event space" • e.g., race, sentence, text ...
- Suppose an event is a sequence of letters: p(horses)
- But we rewrote p(horses) as $p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, ...)$ $\approx p(x_1=h) * p(x_2=o | x_1=h) * ...$
- What does this variable=value notation mean?

Random Variables: What is "variable" in "p(variable= Answer: variable is really a function of Event • $p(x_1=h) * p(x_2=o | x_1=h) * ...$ • Event is a sequence of letters • X_2 is the second letter in the sequence • p(number of heads=2) or just p(H=2) • Event is a sequence of 3 coin flips • H is the number of heads • p(weather's clear=true) or just p(weather's clear) • Event is a race

• weather's clear is true or false



- p(weather's clear=true) or just p(weather's clear) • Event is a race
 - weather's clear (Event) is true or false

Random Variables: What is "variable" in "p(variable=

- p(number of heads=2) or just p(H=2)
 - Event is a sequence of 3 coin flips
 - H is the number of heads in the event
 - So p(H=2)
 - = p(H(Event)=2) picks out the set of events with 2 heads $= p({HHT,HTH,THH})$

TTT

THT

ТТН НТТ

THH HHT HHH

HTH

= p(HHT)+p(HTH)+p(THH)



Random Variables:

What is "variable" in "p(variable)"?

- p(x₁=h) * p(x₂=o | x₁=h) * ...
 - Event is a sequence of letters
 - X_2 is the second letter in the sequence
 - So $p(x_2=0)$
 - = $p(x_2(Event)=0)$ picks out the *set* of events with ...
 - = Σ p(Event) over all events whose second letter ...
 - = p(horses) + p(boffo) + p(xoyzkklp) + ...



A Different Model

• Exploit fact that horses is a common word

$p(W_1 = horses)$

where word vector W is a function of the event (the sentence) just as character vector X is.

- $= p(W_i = horses \mid i=1)$ $\approx p(W_i = horses) = 7.2e-5$

independence assumption says that sentence-initial words w₁ are just like all other words w_i (gives us more data to use)

Much larger than previous estimate of 5.4e-7 – why?

Advantages, disadvantages?

Improving the New Model: Weaken the Indep. Assumption • Don't totally cross off i=1 since it's not irrelevant: - Yes, horses is common, but less so at start of sentence since most sentences start with determiners. $$\begin{split} p(W_1 = \text{horses}) &= \Sigma_t \, p(W_1 = \text{horses}, \ \text{T}_1 = \text{t}) \\ &= \Sigma_t \, p(W_1 = \text{horses} \mid \text{T}_1 = \text{t}) * p(\text{T}_1 = \text{t}) \\ &= \Sigma_t \, p(W_i = \text{horses} \mid \text{T}_i = \text{t}, i = 1) * p(\text{T}_1 = \text{t}) \end{split}$$ $\approx \Sigma_{t} p(W_{i} = horses | T_{i} = t) * p(T_{1} = t)$ $p(W_i = horses | T_i = PlNoun) * p(T_1 = PlNoun)$ (if first factor is 0 for any other part of speech) ≈ (72 / 55912) * (977 / 52108) = 2.4e-5

Which Model is Better?

- Model 1 predict each letter X_i from previous 2 letters X_{i-2}, X_{i-1}
- Model 2 predict each word W_i by its part of speech T_i, having predicted T_i from i
- Models make different independence assumptions that reflect different intuitions
- Which intuition is better???

Measure Performance!

- Which model does better on language ID? - Administer test where you know the right answers
 - Seal up test data until the test happens
 - · Simulates real-world conditions where new data comes along that you didn't have access to when choosing or training model
 - In practice, split off a test set as soon as you obtain the data, and never look at it
 - Need *enough* test data to get statistical significance
- For a different task (e.g., speech transcription instead of language ID), use that task to evaluate the models

Cross-Entropy ("xent")

- Another common measure of model quality
 - Task-independent
 - Continuous so slight improvements show up here even if they don't change # of right answers on task
- Just measure probability of (enough) test data
 - Higher prob means model better predicts the future
 - There's a limit to how well you can predict random stuff
 Limit depends on "how random" the dataset is (easier to predict weather than headlines, especially in Arizona)

Cross-Entropy ("xent")

- Want prob of test data to be high: p(h | BOS, BOS) * p(o | BOS, h) * p(r | h, o) * p(s | o, r) ... 1/8 * 1/8 * 1/8 * 1/8 * 1/16 ...
- high prob → low xent by 3 cosmetic improvements:
 Take logarithm (base 2) to prevent underflow:
 - $\log (1/8 * 1/8 * 1/8 * 1/16 ...) = \log 1/8 + \log 1$
 - Negate to get a positive value in *bits* 3+3+3+4+...
 - Divide by length of text to get bits per letter or bits per word
 - Want this to be small (equivalent to wanting good compression!)
 Lower limit is called *entropy* obtained in principle as cross-entropy of best possible model on an infinite amount of test data
 - Or use *perplexity* = 2 to the xent (9.5 choices instead of 3.25 bits)