

How to Use Probabilities

The Crash Course

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?

3 Kinds of Statistics

- **descriptive:** mean Hopkins SAT (or median)
- **confirmatory:** statistically significant?
- **predictive:** wanna bet?

this course – why?

Notation for Greenhorns



0.9

probability model

$$p(\text{Paul Revere wins} | \text{weather's clear}) = 0.9$$

What does that really mean?

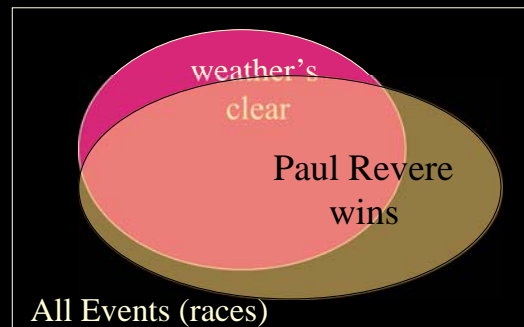
$$p(\text{Paul Revere wins} | \text{weather's clear}) = 0.9$$

- 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?

$$p(X | Y) \text{ versus } q(X | Y)$$

p is a function on event sets

$$p(\text{win} | \text{clear}) \equiv p(\text{win, clear}) / p(\text{clear})$$



p is a function on event sets

$$p(\text{win} \mid \text{clear}) \equiv p(\text{win}, \text{clear}) / p(\text{clear})$$

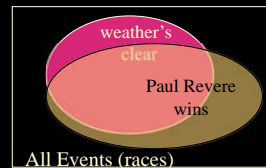
syntactic sugar logical conjunction of predicates predicate selecting races where weather's clear



p measures total probability of a set of events.

Required Properties of p (axioms) ^{most of the}

- $p(\emptyset) = 0$ $p(\text{all events}) = 1$
- $p(X) \leq p(Y)$ for any $X \subseteq Y$
- $p(X) + p(Y) = p(X \cup Y)$ provided $X \cap Y = \emptyset$
 e.g., $p(\text{win} \& \text{clear}) + p(\text{win} \& \neg \text{clear}) = p(\text{win})$



p measures total probability of a set of events.

Commas denote conjunction

$p(\text{Paul Revere wins, Valentine places, Epitaph shows} \mid \text{weather's clear})$

what happens as we add conjuncts to left of bar ?

- probability can only decrease
- numerator of historical estimate likely to go to zero:
 $\frac{\# \text{ times Revere wins AND Val places... AND weather's clear}}{\# \text{ times weather's clear}}$

Commas denote conjunction

$p(\text{Paul Revere wins, Valentine places, Epitaph shows} \mid \text{weather's clear})$

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

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*)
 $\frac{\# \text{ times Revere wins AND weather clear AND ... it's May 17}}{\# \text{ times weather clear AND ... it's May 17}}$

Simplifying Right Side: Backing Off

$p(\text{Paul Revere wins} \mid \text{weather's clear, } \cancel{\text{ground is dry, jockey getting over sprain}}, \text{ Epitaph also in race, } \cancel{\text{Epitaph was recently bought by Gonzalez, race is on May 17, ...}})$

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

$p(\text{Paul Revere wins, } \cancel{\text{Valentine places, Epitaph shows}} \mid \text{weather's clear})$

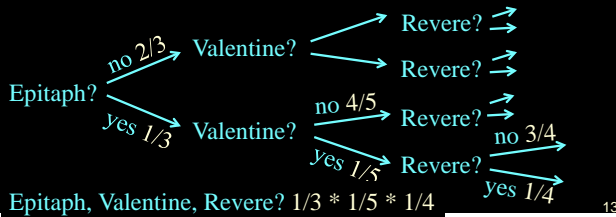
NOT ALLOWED!

but we can do something similar to help ...

Factoring Left Side: The Chain Rule

$$\begin{aligned}
 & p(\text{Revere, Valentine, Epitaph} \mid \text{weather's clear}) = \text{RVEW/W} \\
 = & p(\text{Revere} \mid \text{Valentine, Epitaph, weather's clear}) = \text{RVEW/VEW} \\
 & * p(\text{Valentine} \mid \text{Epitaph, weather's clear}) = \text{VEW/EW} \\
 & * p(\text{Epitaph} \mid \text{weather's clear}) = \text{EW/W}
 \end{aligned}$$

True because numerators cancel against denominators
 Makes perfect sense when read from bottom to top



Epitaph, Valentine, Revere? $1/3 * 1/5 * 1/4$ 13

Factoring Left Side: The Chain Rule

$$\begin{aligned}
 & p(\text{Revere, Valentine, Epitaph} \mid \text{weather's clear}) = \text{RVEW/W} \\
 = & p(\text{Revere} \mid \text{Valentine, Epitaph, weather's clear}) = \text{RVEW/VEW} \\
 & * p(\text{Valentine} \mid \text{Epitaph, weather's clear}) = \text{VEW/EW} \\
 & * p(\text{Epitaph} \mid \text{weather's clear}) = \text{EW/W}
 \end{aligned}$$

True because numerators cancel against denominators
 Makes perfect sense when read from bottom to top
 Moves material to right of bar so it can be ignored

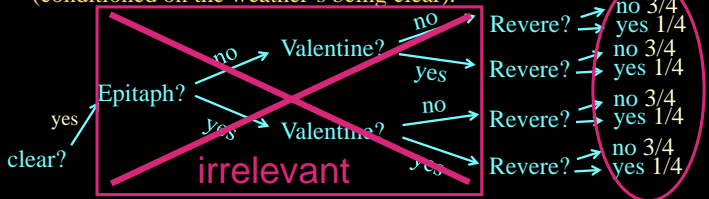
If this prob is unchanged by backoff, we say Revere was **CONDITIONALLY INDEPENDENT** of Valentine and Epitaph (conditioned on the weather's being clear). Often we just **ASSUME** conditional independence to get the nice product above.

Factoring Left Side: The Chain Rule

$$p(\text{Revere} \mid \text{Valentine, Epitaph, weather's clear})$$

conditional independence lets us use backed-off data from all four of these cases to estimate their shared probabilities

If this prob is unchanged by backoff, we say Revere was **CONDITIONALLY INDEPENDENT** of Valentine and Epitaph (conditioned on the weather's being clear).



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 (x_1, x_2, x_3, \dots) where $x_n = \text{EOS}$

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."

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

Apply the Chain Rule

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 = & p(x_1=h) && 4470/52108 \\
 & * p(x_2=o \mid x_1=h) && 395/4470 \\
 & * p(x_3=r \mid x_1=h, x_2=o) && 5/395 \\
 & * p(x_4=s \mid x_1=h, x_2=o, x_3=r) && 3/5 \\
 & * p(x_5=e \mid x_1=h, x_2=o, x_3=r, x_4=s) && 3/3 \\
 & * p(x_6=s \mid x_1=h, x_2=o, x_3=r, x_4=s, x_5=e) && 0/3 \\
 & * \dots = 0
 \end{aligned}$$

counts from Brown corpus

Back Off On Right Side

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx p(x_1=h) && 4470/52108 \\
 & * p(x_2=o | x_1=h) && 395/ 4470 \\
 & * p(x_3=r | x_1=h, x_2=o) && 5/ 395 \\
 & * p(x_4=s | x_2=o, x_3=r) && 12/ 919 \\
 & * p(x_5=e | x_3=r, x_4=s) && 12/ 126 \\
 & * p(x_6=s | x_4=s, x_5=e) && 3/ 485 \\
 & * \dots = 7.3e-10 * \dots
 \end{aligned}$$

counts from
Brown corpus

Change the Notation

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx p(x_1=h) && 4470/52108 \\
 & * p(x_2=o | x_1=h) && 395/ 4470 \\
 & * p(x_3=r | x_{i-2}=h, x_{i-1}=o, i=3) && 5/ 395 \\
 & * p(x_4=s | x_{i-2}=o, x_{i-1}=r, i=4) && 12/ 919 \\
 & * p(x_5=e | x_{i-2}=r, x_{i-1}=s, i=5) && 12/ 126 \\
 & * p(x_6=s | x_{i-2}=s, x_{i-1}=e, i=6) && 3/ 485 \\
 & * \dots = 7.3e-10 * \dots
 \end{aligned}$$

counts from
Brown corpus

Another Independence Assumption

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx p(x_1=h) && 4470/52108 \\
 & * p(x_2=o | x_1=h) && 395/ 4470 \\
 & * p(x_3=r | x_{i-2}=h, x_{i-1}=o) && 1417/14765 \\
 & * p(x_4=s | x_{i-2}=o, x_{i-1}=r) && 1573/26412 \\
 & * p(x_5=e | x_{i-2}=r, x_{i-1}=s) && 1610/12253 \\
 & * p(x_6=s | x_{i-2}=s, x_{i-1}=e) && 2044/21250 \\
 & * \dots = 5.4e-7 * \dots
 \end{aligned}$$

counts from
Brown corpus

Simplify the Notation

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx p(x_1=h) && 4470/52108 \\
 & * p(x_2=o | x_1=h) && 395/ 4470 \\
 & * p(r | h, o) && 1417/14765 \\
 & * p(s | o, r) && 1573/26412 \\
 & * p(e | r, s) && 1610/12253 \\
 & * p(s | s, e) && 2044/21250 \\
 & * \dots
 \end{aligned}$$

counts from
Brown corpus

Simplify the Notation

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx p(h | \text{BOS, BOS}) && 4470/52108 \\
 & * p(o | \text{BOS, h}) && 395/ 4470 \\
 & * p(r | h, o) && 1417/14765 \\
 & * p(s | o, r) && 1573/26412 \\
 & * p(e | r, s) && 1610/12253 \\
 & * p(s | s, e) && 2044/21250 \\
 & * \dots && \text{These basic probabilities} \\
 & && \text{are used to define } p(\text{horses})
 \end{aligned}$$

the parameters
of our old
trigram generator!
Same assumptions
about language.

values of
those
parameters,
as naively
estimated
from Brown
corpus.

counts from
Brown corpus

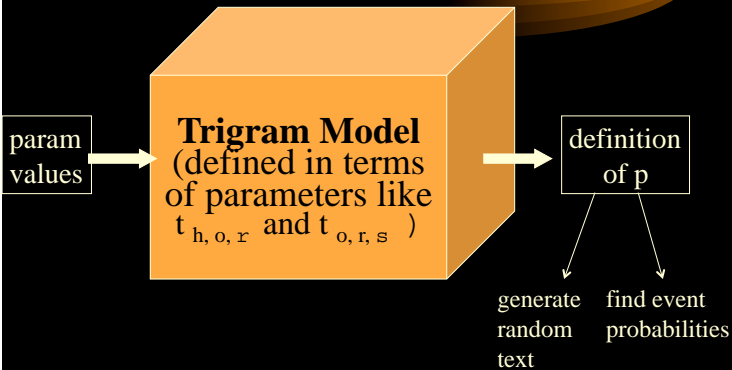
Simplify the Notation

$$\begin{aligned}
 & p(x_1=h, x_2=o, x_3=r, x_4=s, x_5=e, x_6=s, \dots) \\
 & \approx t_{\text{BOS, BOS, h}} && 4470/52108 \\
 & * t_{\text{BOS, h, o}} && 395/ 4470 \\
 & * t_{h, o, r} && 1417/14765 \\
 & * t_{o, r, s} && 1573/26412 \\
 & * t_{r, s, e} && 1610/12253 \\
 & * t_{s, e, s} && 2044/21250 \\
 & * \dots && \text{This notation emphasizes that} \\
 & && \text{they're just real variables} \\
 & && \text{whose value must be estimated}
 \end{aligned}$$

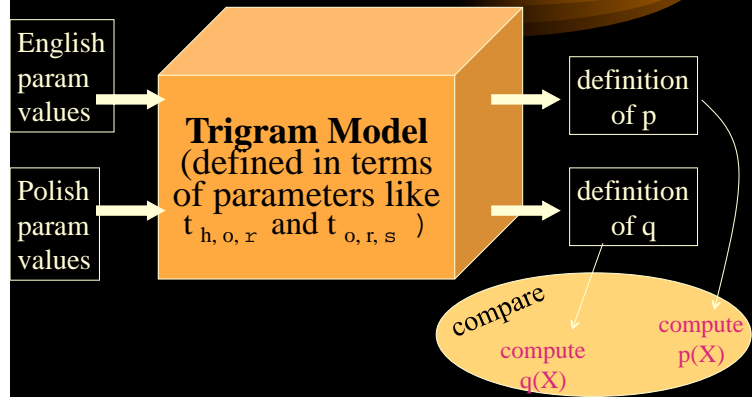
values of
those
parameters,
as naively
estimated
from Brown
corpus.

counts from
Brown corpus

Definition: Probability Model

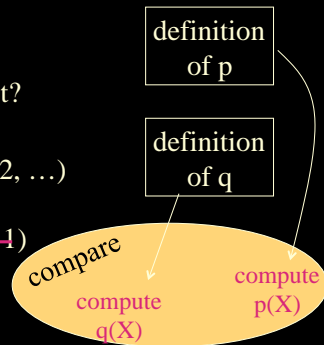


English vs. Polish



What is “X” in $p(X)$?

- Element of some implicit “event space”
 - e.g., race
 - e.g., sentence
- What if event is a whole text?
 - $p(\text{text})$
 - $= p(\text{sentence 1}, \text{sentence 2}, \dots)$
 - $= p(\text{sentence 1})$
 - * $p(\text{sentence 2} | \text{sentence 1})$
 - * ...



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:

$$p(\text{horses})$$
- But we rewrite $p(\text{horses})$ as

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

$$\approx p(x_1=\text{h}) * p(x_2=\text{o} | x_1=\text{h}) * \dots$$
- What does this variable=value notation mean?

Random Variables:

What is “variable” in “ $p(\text{variable}=\text{value})$ ”?

Answer: variable is really a function of Event

- $p(x_1=\text{h}) * p(x_2=\text{o} | x_1=\text{h}) * \dots$
 - Event is a sequence of letters
 - x_2 is the second letter in the sequence
- $p(\text{number of heads}=\text{2})$ or just $p(H=\text{2})$
 - Event is a sequence of 3 coin flips
 - H is the number of heads
- $p(\text{weather's clear}=\text{true})$ or just $p(\text{weather's clear})$
 - Event is a race
 - weather's clear is true or false

Random Variables:

What is “variable” in “ $p(\text{variable}=\text{value})$ ”?

Answer: variable is really a function of Event

- $p(x_1=\text{h}) * p(x_2=\text{o} | x_1=\text{h}) * \dots$
 - Event is a sequence of letters
 - $x_2(\text{Event})$ is the second letter in the sequence
- $p(\text{number of heads}=\text{2})$ or just $p(H=\text{2})$
 - Event is a sequence of 3 coin flips
 - $H(\text{Event})$ is the number of heads
- $p(\text{weather's clear}=\text{true})$ or just $p(\text{weather's clear})$
 - Event is a race
 - weather's clear (Event) is true or false

Random Variables:

What is “variable” in “ $p(\text{variable}=\text{value})$ ”?

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

TTT	TTH	HTT	HTH
THT	THH	HHT	HHH

All Events

31

600.465 – Intro to NLP – J. Eisner

Random Variables:

What is “variable” in “ $p(\text{variable}=\text{value})$ ”?

- $p(\text{weather's clear})$
 - Event is a race
 - weather's clear is true or false of the event
- So $p(\text{weather's clear})$
 - = $p(\text{weather's clear}(\text{Event})=\text{true})$
 - picks out the *set* of events with clear weather



$$p(\text{win} | \text{clear}) \equiv p(\text{win, clear}) / p(\text{clear})$$

All Events (races)

32

600.465 – Intro to NLP – J. Eisner

Random Variables:

What is “variable” in “ $p(\text{variable}=\text{value})$ ”?

- $p(x_1=h) * p(x_2=o | x_1=h) * \dots$
 - Event is a sequence of letters
 - x_2 is the second letter in the sequence
- So $p(x_2=o)$
 - = $p(x_2(\text{Event})=o)$ picks out the *set* of events with ...
 - = $\sum p(\text{Event})$ over all events whose second letter ...
 - = $p(\text{horses}) + p(\text{boffo}) + p(\text{xoyzkklp}) + \dots$

600.465 – Intro to NLP – J. Eisner

33

Back to trigram model of $p(\text{horses})$

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

$$\approx t_{\text{BOS, BOS, h}}$$

$$* t_{\text{BOS, h, o}}$$

$$* t_{\text{h, o, r}}$$

$$* t_{\text{o, r, s}}$$

$$* t_{\text{r, s, e}}$$

$$* t_{\text{s, e, s}}$$

* ... This notation emphasizes that they're just real variables whose value must be estimated

the parameters of our old trigram generator! Same assumptions about language.

values of those parameters, as naively estimated from Brown corpus.

4470/ 52108

395/ 4470

1417/ 14765

1573/ 26412

1610/ 12253

2044/ 21250

counts from Brown corpus

600.465 – Intro to NLP – J. Eisner

34

A Different Model

- Exploit fact that horses is a common word

$$p(W_1 = \text{horses})$$

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

$$= p(W_i = \text{horses} | i=1)$$

$$\approx p(W_i = \text{horses}) = 7.2e-5$$

independence assumption says that sentence-initial words w_1 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?

600.465 – Intro to NLP – J. Eisner

35

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.

$$p(W_1 = \text{horses}) = \sum_i p(W_1 = \text{horses}, T_1 = t)$$

$$= \sum_i p(W_1 = \text{horses} | T_1 = t) * p(T_1 = t)$$

$$= \sum_i p(W_i = \text{horses} | T_i = t, i=1) * p(T_1 = t)$$

$$\approx \sum_i p(W_i = \text{horses} | T_i = t) * p(T_1 = t)$$

$$= p(W_i = \text{horses} | T_i = \text{PlNoun}) * p(T_1 = \text{PlNoun})$$

$$+ p(W_i = \text{horses} | T_i = \text{Verb}) * p(T_1 = \text{Verb}) + \dots$$

$$= (72 / 55912) * (977 / 52108) + (0 / 15258) * (146 / 52108) + \dots$$

$$= 2.4e-5 + 0 + \dots + 0 = 2.4e-5$$

600.465 – Intro to NLP – J. Eisner

36

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:
$$\frac{p(h | \text{BOS}, \text{BOS})}{1/8} * \frac{p(o | \text{BOS}, h)}{1/8} * \frac{p(r | h, o)}{1/8} * \frac{p(s | o, r)}{1/16} \dots$$
- high prob → low xent by 3 cosmetic improvements:
 - Take logarithm (base 2) to prevent underflow:
$$\log(1/8 * 1/8 * 1/8 * 1/16 \dots)$$
$$= \log 1/8 + \log 1/8 + \log 1/8 + \log 1/16 \dots = (-3) + (-3) + (-3) + (-4) + \dots$$
 - Negate to get a positive value in *bits* $3+3+3+4+\dots$
 - 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)