# Lecture 2:
# Probability and Language Models

Intro to NLP, CS585, Fall 2014
Brendan O'Connor (http://brenocon.com)

1

# Admin

- Waitlist
- Moodle access: Email me if you don't have it
- Did you get an announcement email?
- Piazza vs Moodle?
- Office hours today

Thursday, September 4, 14

# Things today

- Homework: ambiguities
- Python demo
- Probability Review
- Language Models

3

4

# Python demo

- [TODO link ipython-notebook demo]

- For next week, make sure you can run
  - Python 2.7  (Built-in on Mac & Linux)
  - IPython Notebook  http://ipython.org/notebook.html
    - Please familiarize yourself with it.
    - Python 2.7, IPython 2.2.0
  - Nice to have: Matplotlib
- Python interactive interpreter
- Python scripts
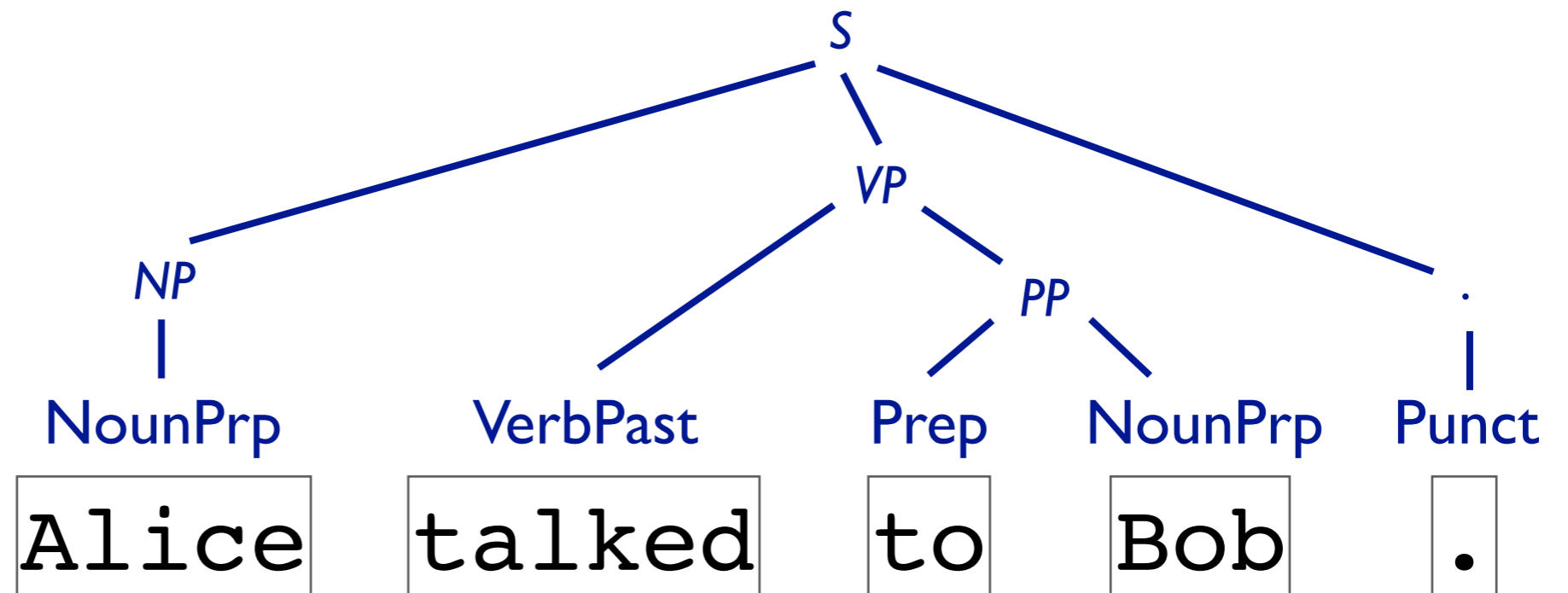
5

6

# Levels of linguistic structure

**Discourse**

**Semantics**

CommunicationEvent(e)  SpeakerContext(s)
Agent(e, Alice)  TemporalBefore(e, s)
Recipient(e, Bob)

**Syntax**

S
 NP
 VP
 PP
NounPrp   VerbPast   Prep   NounPrp   Punct

**Words**

Alice  talked  to  Bob  .

**Morphology**

talk -ed

**Characters**

Alice talked to Bob.

7

# Levels of linguistic structure

Words are fundamental units of meaning

and easily identifiable*

*in some languages

| Words | Alice | talked | to | Bob | . |

| Characters | Alice talked to Bob. |

# Probability theory

# Review: definitions/laws

$$\boxed{\phantom{XX}} = \sum_a P(A = a)$$

**Conditional Probability**

$$\boxed{\phantom{XXXXX}} = \frac{P(AB)}{P(B)}$$

**Chain Rule**

$$\boxed{\phantom{XXXXX}} = P(A|B)P(B)$$

**Law of Total Probability**

$$\boxed{\phantom{XXX}} = \sum_b P(A, B = b)$$

$$\boxed{\phantom{XXXX}} = \sum_b P(A|B = b)P(B = b)$$

**Disjunction (Union)** $\quad P(A \vee B) = \boxed{\phantom{XXXXXXXXXXX}}$

**Negation (Complement)** $\quad P(\neg A) = \boxed{\phantom{XXXXXX}}$

10

# Bayes Rule

Want P(H|D) but only have P(D|H)
e.g. H causes D, or P(D|H) is easy to measure...

H: who wrote this document?

Model: authors' word probs

Bayesian inference

D: words

Likelihood

Prior

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)}$$

Posterior

Normalizer

Rev. Thomas Bayes
c. 1701-1761

# Bayes Rule and its pesky denominator

Likelihood

Prior

$$P(h|d) = \frac{P(d|h)P(h)}{P(d)} = \frac{P(d|h)P(h)}{\sum_{h'} P(d|h')P(h')}$$

$$P(h|d) = \frac{1}{Z} P(d|h)P(h)$$

*Z*: whatever lets the posterior, when summed across *h*, to sum to 1
*Zustandssumme, "sum over states"*

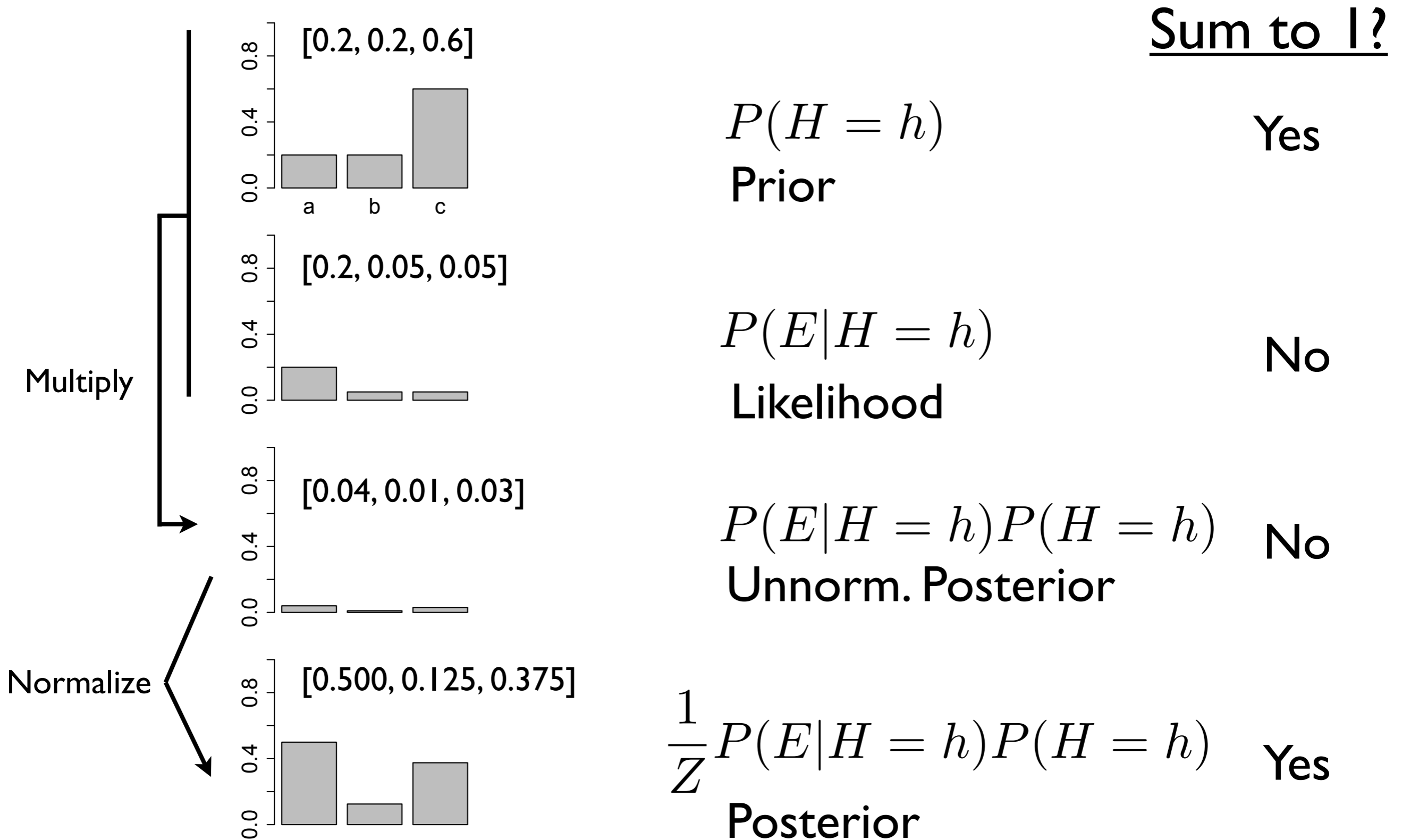$$P(h|d) \propto \underline{P(d|h)P(h)}$$

"Proportional to"
(implicitly for varying H.
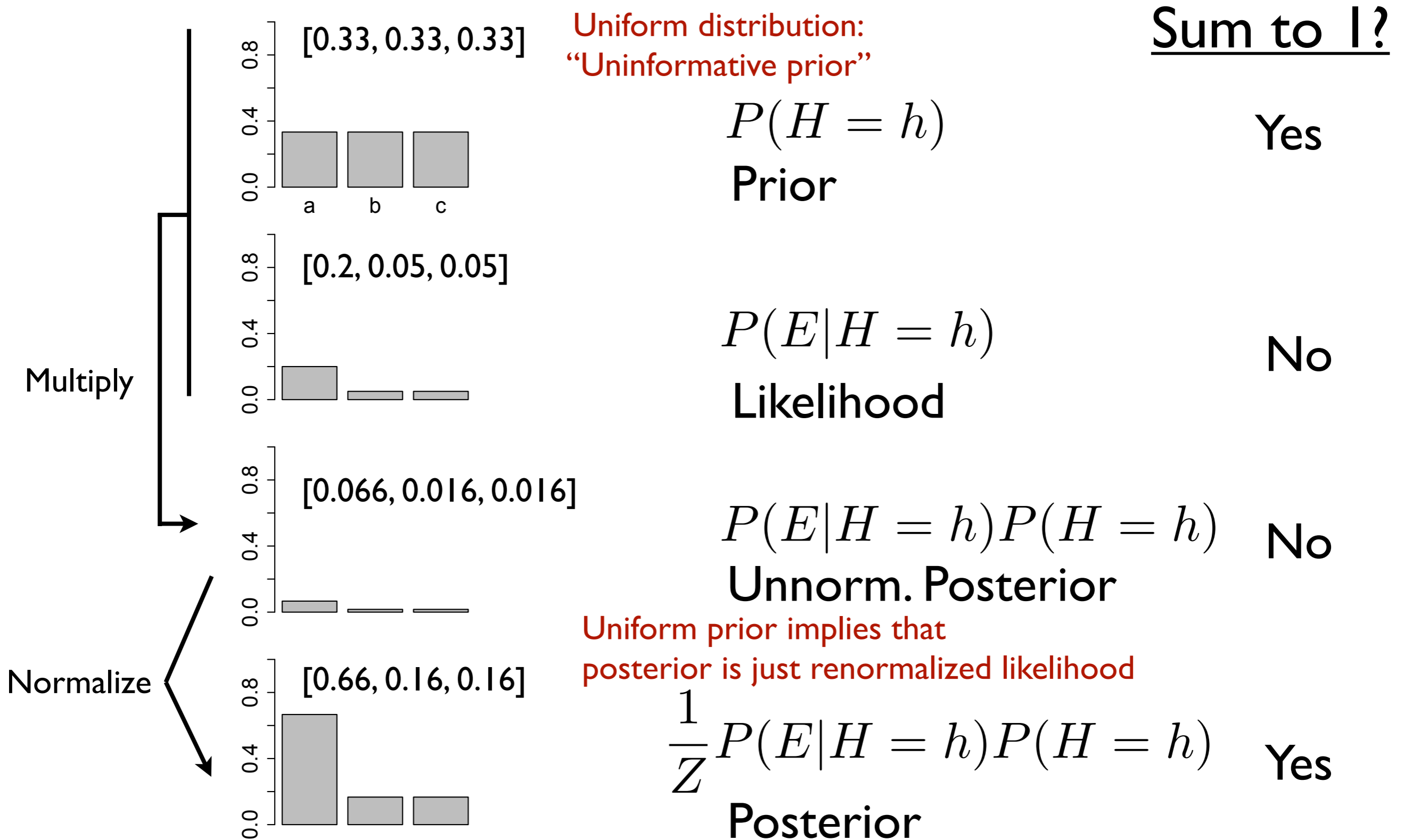This notation is very common, though slightly ambiguous.)
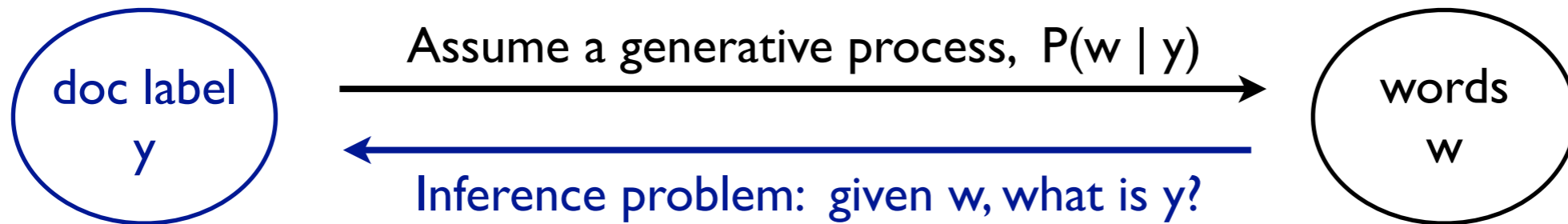
Unnormalized posterior
By itself does not sum to 1!

12

# Bayes Rule: Discrete

[0.2, 0.2, 0.6]

$$P(H = h)$$

Prior

Yes

[0.2, 0.05, 0.05]

$$P(E|H = h)$$

Likelihood

No

**Multiply**

[0.04, 0.01, 0.03]

$$P(E|H = h)P(H = h)$$

Unnorm. Posterior

No

**Normalize**

[0.500, 0.125, 0.375]

$$\frac{1}{Z}P(E|H = h)P(H = h)$$

Posterior

Yes

13

# Bayes Rule: Discrete, uniform prior

[0.33, 0.33, 0.33]

Uniform distribution:
"Uninformative prior"

$$P(H = h)$$

Prior

Yes

[0.2, 0.05, 0.05]

$$P(E|H = h)$$

Likelihood

No

[0.066, 0.016, 0.016]

$$P(E|H = h)P(H = h)$$

Unnorm. Posterior

No

Uniform prior implies that
posterior is just renormalized likelihood

[0.66, 0.16, 0.16]

Multiply

Normalize

$$\frac{1}{Z}P(E|H = h)P(H = h)$$

Posterior

Yes

14

# Bayes Rule for doc classification

(doc label y) → Assume a generative process, P(w | y) → (words w)

Inference problem: given w, what is y?

If we knew $P(w|y)$
We could estimate $P(y|w) \propto P(y)P(w|y)$

|  | *abracadabra* | *gesundheit* |
|---|---|---|
| Anna | 5 per 1000 words | 6 per 1000 words |
| Barry | 10 per 1000 words | 1 per 1000 words |

Look at random word.
It is *abracadabra*

Assume 50% prior prob
Prob author is Anna?

15

# Bayes Rule for doc classification

doc label
y

Assume a generative process, P(w | y)

Inference problem: given w, what is y?

words
w

If we knew
We could estimate

$$P(w|y)$$
$$P(y|w) \propto P(y)P(w|y)$$

|       | *abracadabra*       | *gesundheit*       |
|-------|---------------------|--------------------|
| Anna  | 5 per 1000 words    | 6 per 1000 words   |
| Barry | 10 per 1000 words   | 1 per 1000 words   |

Look at two random words.

$w_1$ = *abracadabra*

$w_2$ = *gesundheit*

Assume 50% prior prob
Prob author is Anna?

# Bayes Rule for doc classification

( doc label y )  →  Assume a generative process, P(w | y)  →  ( words w )

Inference problem: given w, what is y?

If we knew  $P(w|y)$

We could estimate  $P(y|w) \propto P(y)P(w|y)$

|        | *abracadabra*        | *gesundheit*        |
|--------|----------------------|---------------------|
| Anna   | 5 per 1000 words     | 6 per 1000 words    |
| Barry  | 10 per 1000 words    | 1 per 1000 words    |

Look at two random words.

$w_1$ = *abracadabra*

$w_2$ = *gesundheit*

Chain rule:

$P(w_1, w_2 \mid y) = P(w_1 \mid w_2\ y)\ P(w_2 \mid y)$

ASSUME conditional independence:

$P(w_1, w_2 \mid y) = P(w_1 \mid y)\ P(w_2 \mid y)$

Assume 50% prior prob
Prob author is Anna?

# Cond indep. assumption: "Naive Bayes"

doc label
y

Assume a generative process, $P(w \mid y)$

Inference problem: given w, what is y?

words
w

$$P(w_1 \ldots w_T \mid y) = \prod_{t=1}^{T} P(w_t \mid y)$$

$$\text{each } w_t \in 1..V \qquad V = \text{vocabulary size}$$
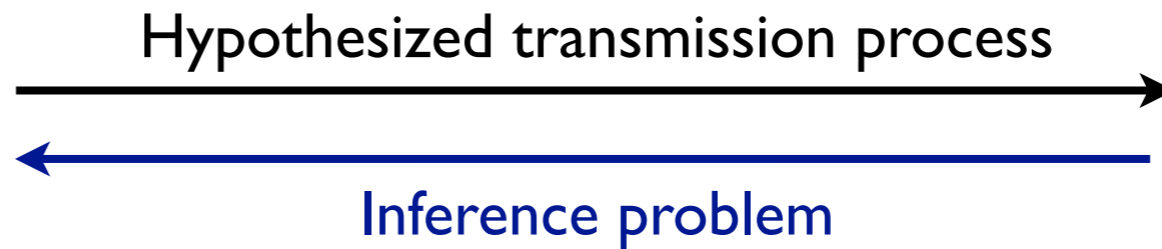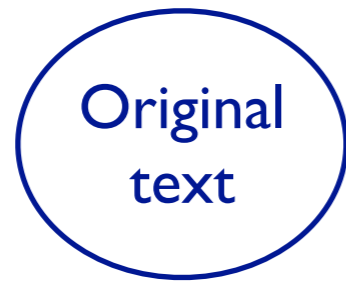
Generative story ("Multinom NB" [McCallum & Nigam 1998]):
  - For each token *t* in the document,
    - Author chooses a word
    by rolling the same weighted V-sided die
This model is wrong!
How can it possibly be useful for doc classification?

# Bayes Rule for *text* inference

Noisy channel model

Original text → Hypothesized transmission process → Observed data

Inference problem

## Codebreaking

$$P(\text{plaintext} \mid \text{encrypted text}) \propto P(\text{encrypted text} \mid \text{plaintext}) \, P(\text{plaintext})$$
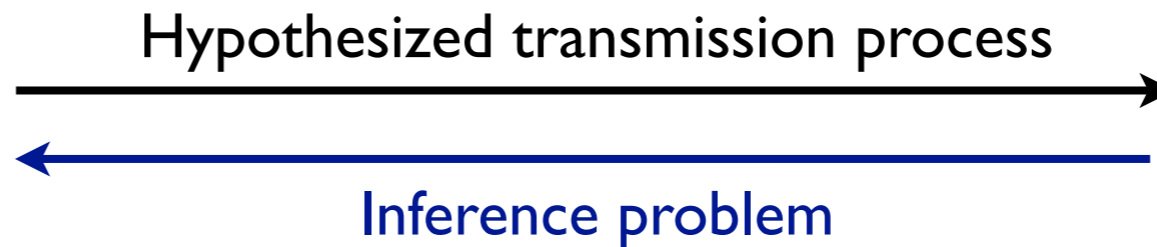


Bletchley Park (WWII)



Enigma machine

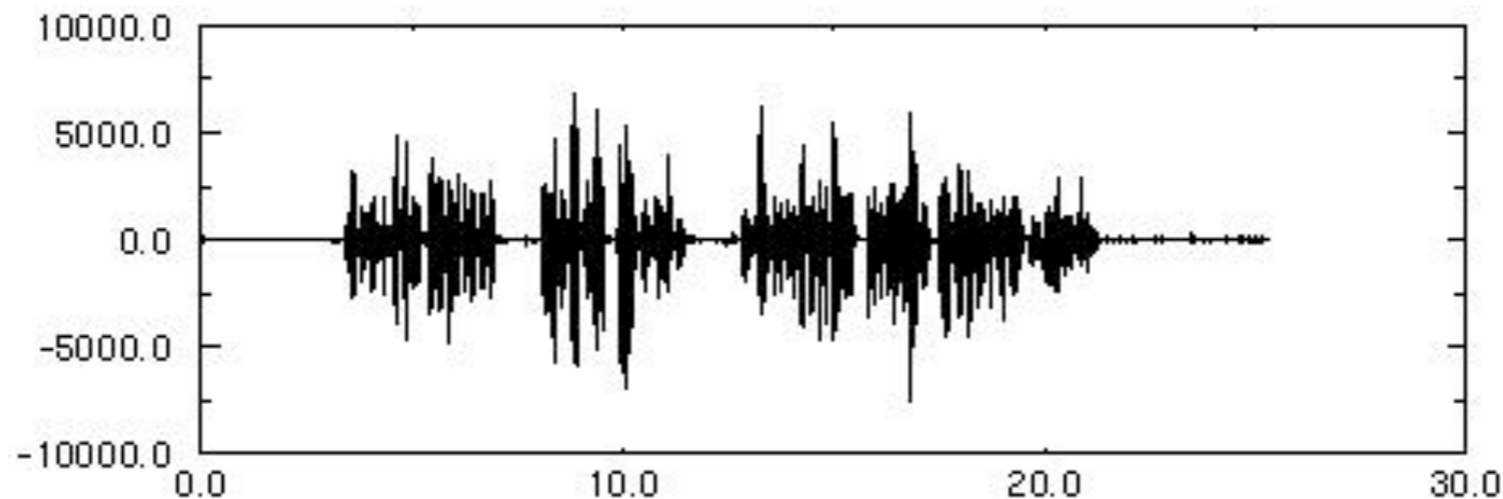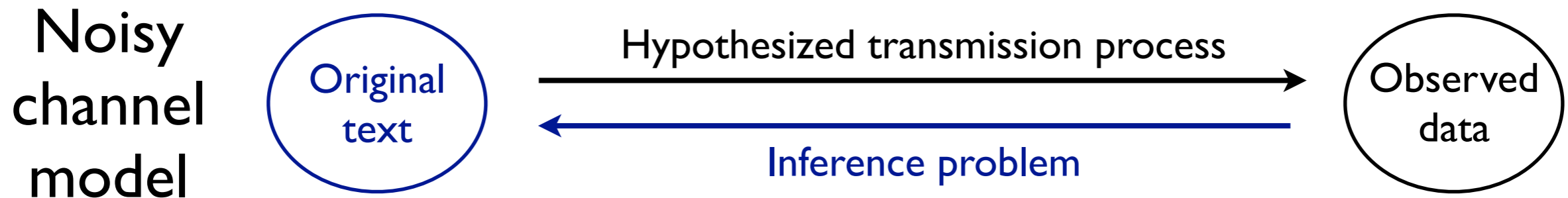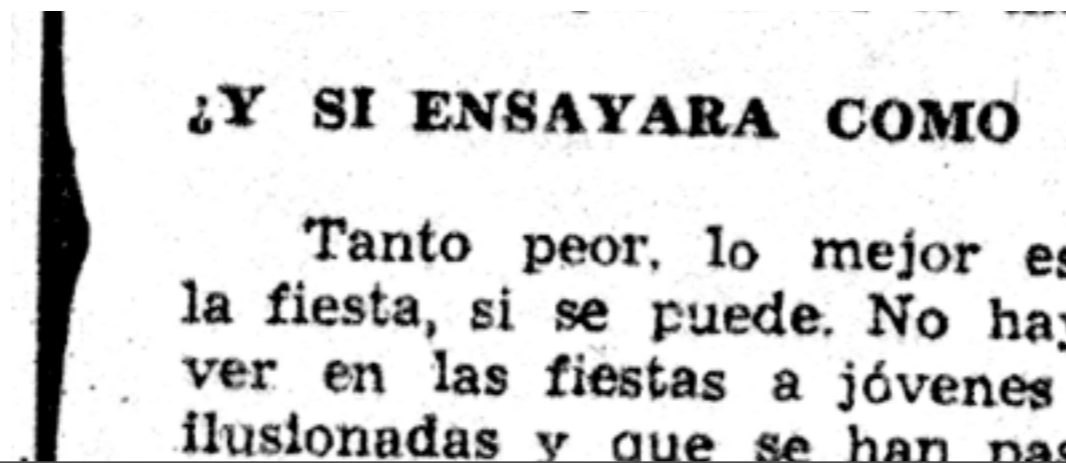# Bayes Rule for *text* inference

Noisy channel model

Original text  → Hypothesized transmission process → Observed data

← Inference problem

## Codebreaking

P(plaintext | encrypted text) $\propto$ P(encrypted text | plaintext) P(plaintext)

## Speech recognition

P(text | acoustic signal) $\propto$ P(acoustic signal | text) P(text)

# Bayes Rule for *text* inference

Noisy channel model

Original text → (Hypothesized transmission process) → Observed data

← Inference problem

## Codebreaking

P(plaintext | encrypted text) $\propto$ P(encrypted text | plaintext) P(plaintext)

## Speech recognition

P(text | acoustic signal) $\propto$ P(acoustic signal | text) P(text)

## Optical character recognition

P(text | image) $\propto$ P(image | text) P(text)

¿Y SI ENSAYARA COMO

Tanto peor, lo mejor es
la fiesta, si se puede. No hay
ver en las fiestas a jóvenes
ilusionadas y que se han pas

# Bayes Rule for *text* inference

Noisy channel model

Original text → (Hypothesized transmission process) → Observed data

Original text ← (Inference problem) ← Observed data
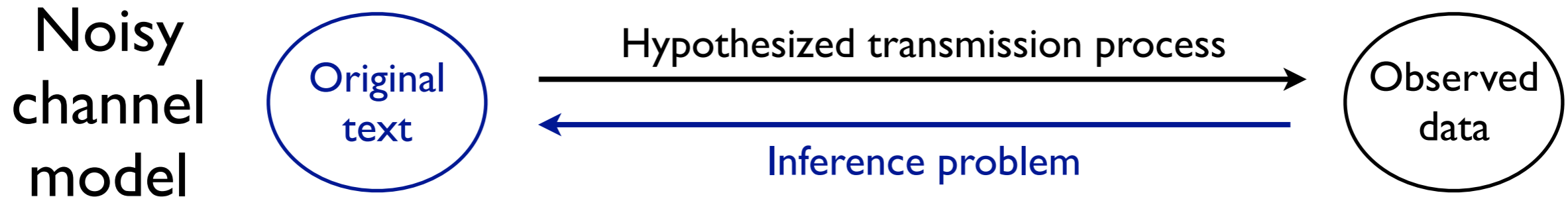
Codeb...  P(plaintex...  ...text)

Speech...  P(text | a...

Optical...  P(text | i...

One naturally wonders if the problem of translation could conceivably be treated as a problem in cryptography. When I look at an article in Russian, I say: 'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.'

-- Warren Weaver (1955)

# Bayes Rule for *text* inference

Noisy channel model



Original text → Hypothesized transmission process → Observed data

Inference problem ←

## Codebreaking

P(plaintext | encrypted text) $\propto$ P(encrypted text | plaintext) P(plaintext)

## Speech recognition

P(text | acoustic signal) $\propto$ P(acoustic signal | text) P(text)

## Optical character recognition

P(text | image) $\propto$ P(image | text) P(text)

## Machine translation?

P(target text | source text) $\propto$ P(source text | target text) P(target text)