

Information Theory



What is language about?

- A critical function of language is to communicate information.

Information theory is the study of how information is stored and exchanged (communicated).

Today we will explore some hypotheses about language as an *efficient information communication system*.

Communicative efficiency

Communicative efficiency hypothesis:

More predictable meanings are expressed with shorter / faster forms because this leads to efficient communication.

Communicative robustness

Communicative robustness hypothesis:

More predictable meanings are expressed with shorter / faster forms because it is important for infrequent meanings to be expressed in a way that is robust to error.

Probability review

- **Probability: $p(X)$**

$$p(X)$$

How likely an event is to occur.

- **Probability distribution:**

$$p(\text{even}) = 0.5 \quad p(\text{odd}) = 0.5$$

A description of a phenomenon in terms of the probabilities of all possible outcomes. Sums to 1.

- **Conditional probability: $p(X|Y)$**

$$p(X|Y)$$

The chance of event X occurring given that event Y occurs.

If events are truly independent, $p(X|Y) = p(X)$.

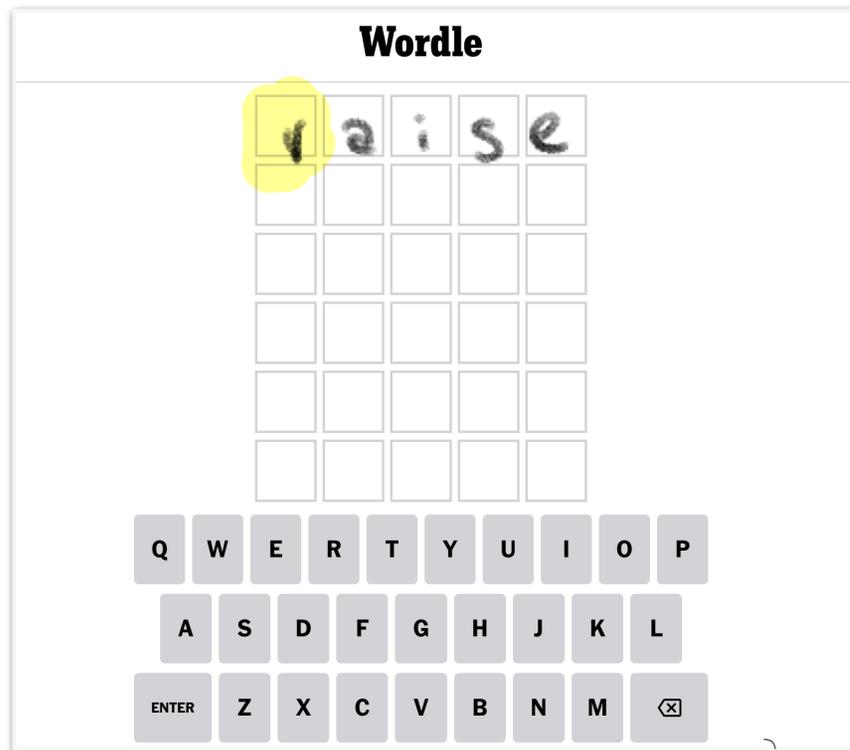
- **Joint probability: $p(x,y)$**

The chance of event X and event Y both occurring. If

events are truly independent, $p(X,Y) = p(X)p(Y)$.

Statistics in language

- You have implicit knowledge about the **probability** of letters in English.
- You also have implicit knowledge about the **conditional probability** of letters in English.



Estimating probability by frequency

Sample text:

"on wednesdays, we wear pink."

Total count of letters: 22

$$p(w) = 3/22$$

$$p(e) = 4/22$$

$$p(e | w) = 3/3$$

$$p(w, e) = p(w) p(e | w) \\ 3/22 \cdot 3/3 = 3/22$$

Zipfian hypotheses

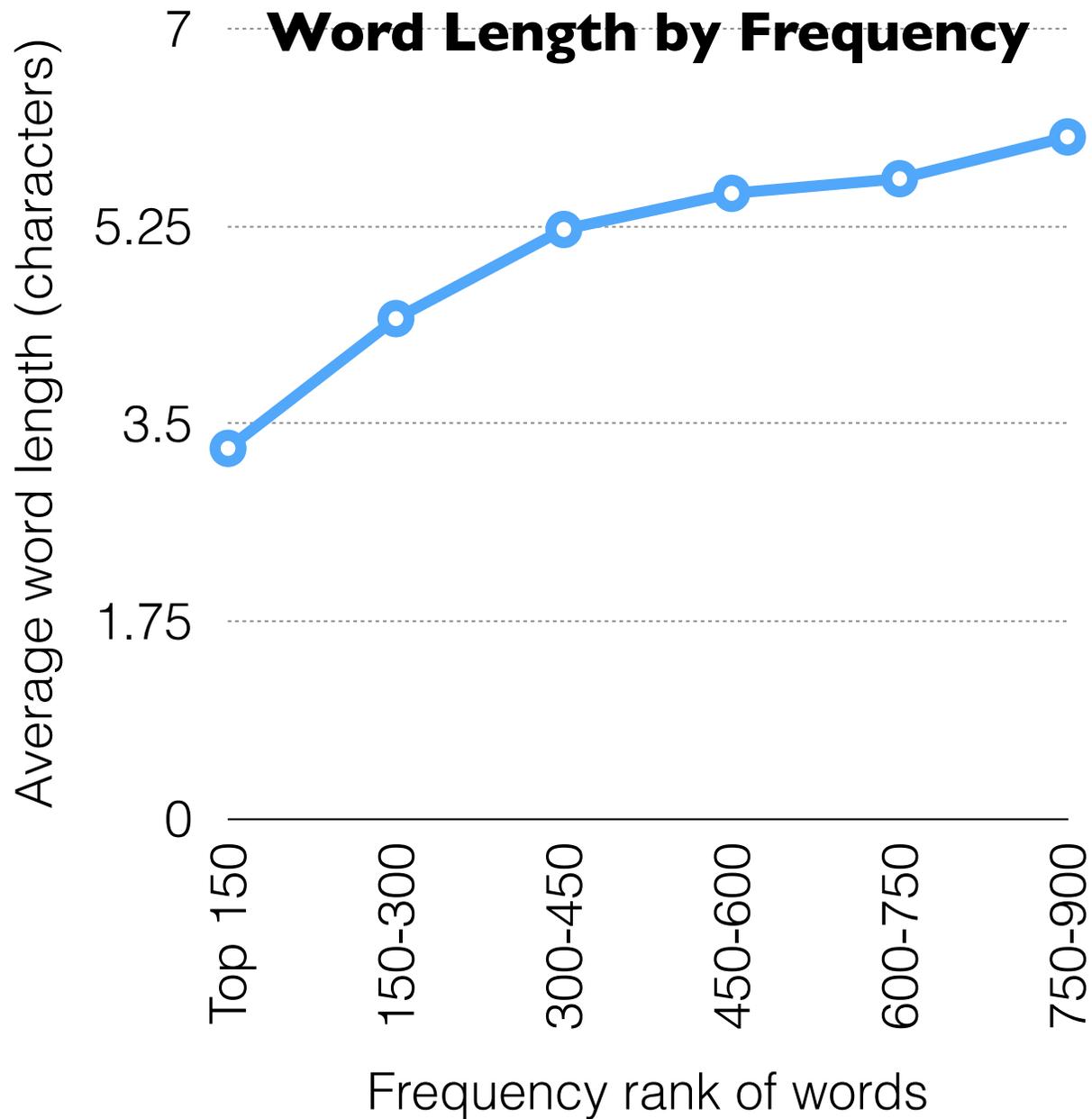
Zipf's law: The frequency of a word is inversely proportional to its frequency ranking.

We'll look into this next class!

Zipf's hypothesis:

Shorter words are more frequent because languages maximize efficiency: they assign common meanings to words that take less effort to produce.

Zipfian hypotheses



Zipf's hypothesis:
Shorter words are more frequent because languages maximize efficiency.

Statistics from the Brown corpus

Communicative efficiency

How can we code meanings efficiently?

Imagine I have a bag of marbles with three colors: **blue**, **red**, and **green**. There are twice as many red marbles as blue and twice as many blue as green.

I am going to reach into the bag, pick a marble, and tell you what color it is.

Here's the trick: the only words I'm allowed to say are **SNUFFLEUPAGUS**, **SHAMBLE**, and **SQUEAK**.

Communicative efficiency

How can we code meanings efficiently?

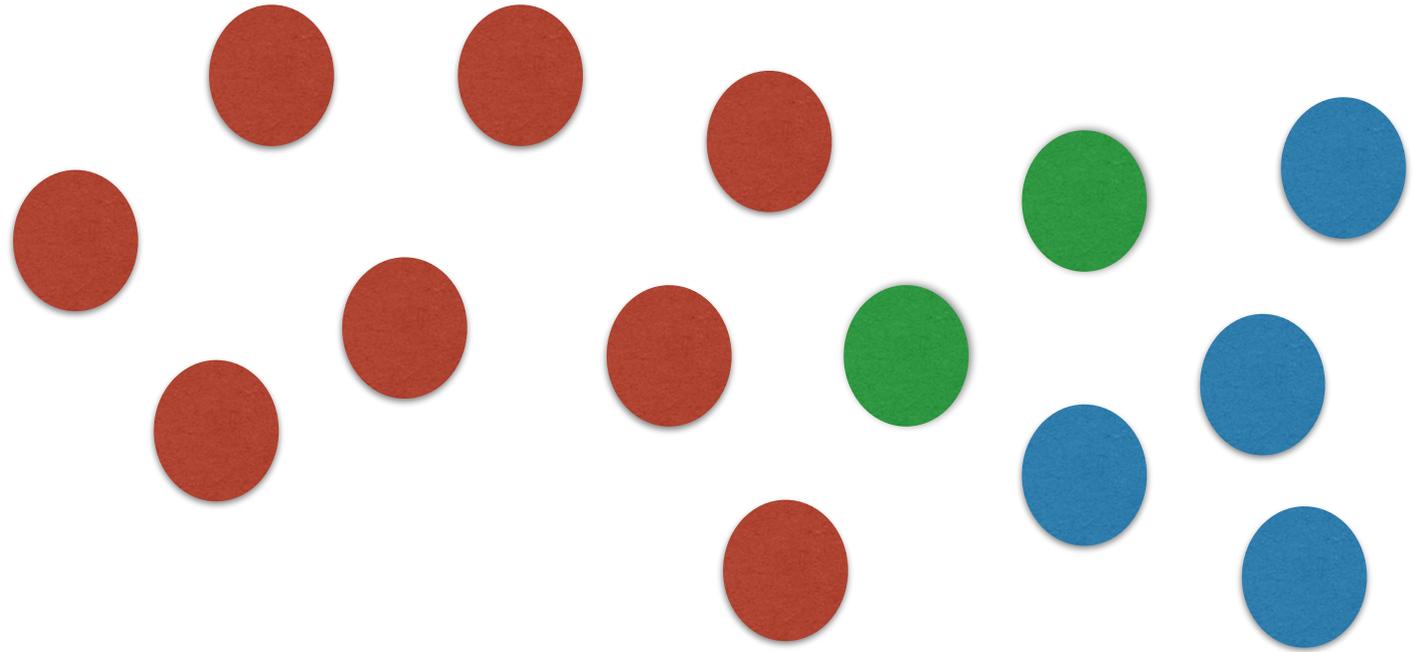
Imagine I have a bag of marbles with three colors: **blue**, **red**, and **green**. There are twice as many red marbles as blue and twice as many blue as green.

I am going to close my eyes, pick a marble out of the bag, and I want you to yell out what color it is.

Here's the trick: the only words you can yell are **SNUFFLEUPAGUS**, **SHAMBLE**, and **SQUEAK**. And we want to do this as **fast as possible**.

Communicative efficiency

READY?

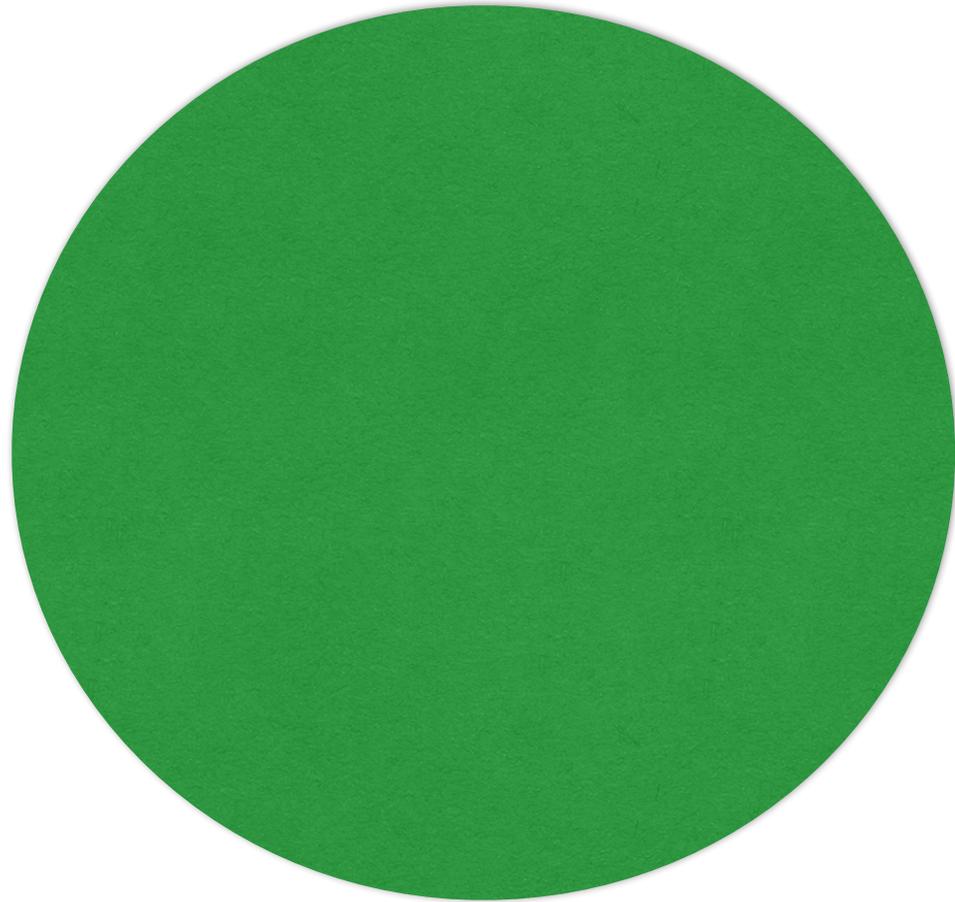


SNUFFLEUPAGUS

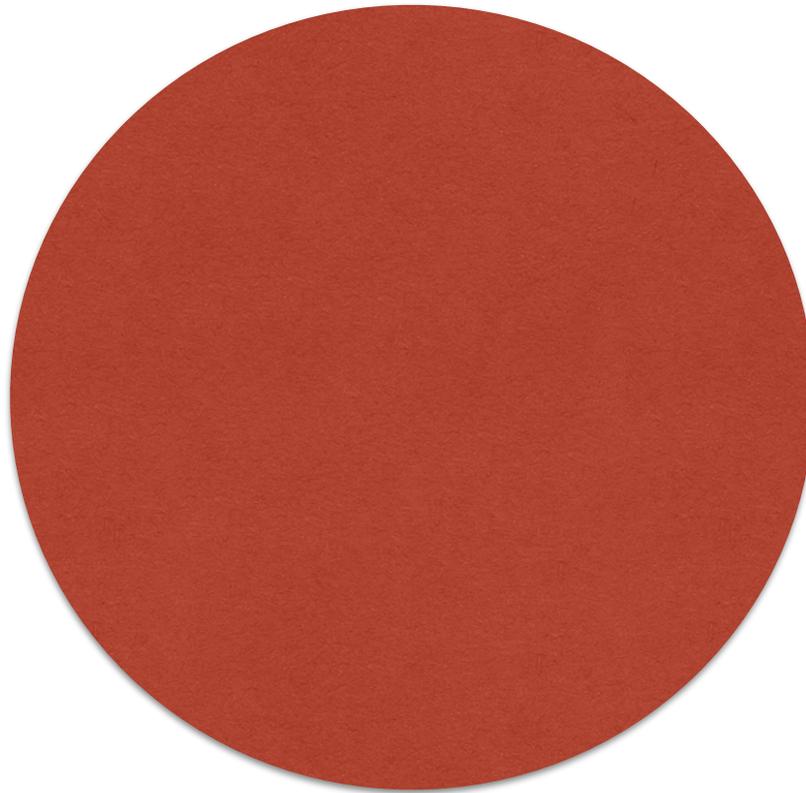
SHAMBLE

SQUEAK

Communicative efficiency

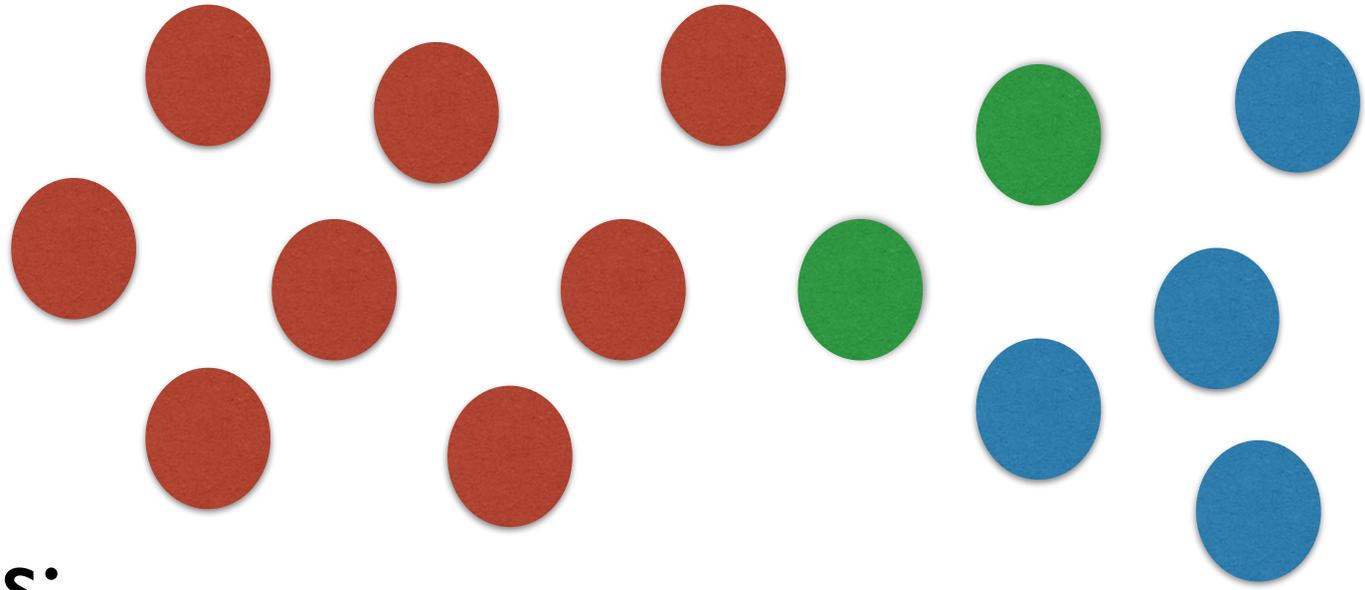


Communicative efficiency



Communicative efficiency

Distribution
of marbles:



Available words:

SNUFFLEUPAGUS

SHAMBLE

SQUEAK

What is the most **efficient** assignment?

Communicative efficiency

Distribution of marbles:

$p(\text{red})$	$= 8/14 = 57.1\%$
$p(\text{blue})$	$= 4/14 = 28.6\%$
$p(\text{green})$	$= 2/14 = 14.3\%$

Available words:

SNUFFLEUPAGUS SHAMBLE SQUEAK

What is the most **efficient** assignment?

Communicative efficiency

Assuming longer words are harder to say, the best arrangement is:

$p(\text{red})$	$= 8/14 = 57.1\%$	SQUEAK
$p(\text{blue})$	$= 4/14 = 28.6\%$	SHAMBLE
$p(\text{green})$	$= 2/14 = 14.3\%$	SNUFFLEUPAGUS

Communicative efficiency

Assuming longer words are harder to say, the best arrangement is:

$p(\text{red}) = 57.1\%$

$p(\text{blue}) = 28.6\%$

$p(\text{green}) = 14.3\%$

SQUEAK

SHAMBLE

SNUFFLEUPAGUS

Communicative efficiency

Communicative efficiency hypothesis:

More predictable meanings are expressed with shorter / faster forms because this leads to efficient communication.

Communicative robustness

Communicative robustness hypothesis:

More predictable meanings are expressed with shorter / faster forms because it is important for infrequent meanings to be expressed in a way that is robust to error.

Communicative robustness

Hypothesis: the more unlikely a word is, the worse it is to make a speech error.

Imagine we're playing the same weird marble game.

But this time, Jess is standing there with an airhorn, making earsplitting noises at random intervals.

Communicative robustness

I'd argue that the strategy of assigning longer words to rarer colors is still a good one, but for a different reason.

Why?

Communicative robustness

I'd argue that the strategy of assigning longer words to rarer colors is still a good one, but for a different reason.

Why?

If you hear nothing but airhorn on a particular turn, what color should you guess?

Communicative robustness

If you shout SNUFFLEUPAGUS, and the airhorn blocks out one syllable, I'll probably still hear enough to know what you said.

But if you shout SQUEAK, I might not.

If your message is **rare**, and the **channel is noisy**, then it makes sense to build some **redundancy** into your message.

Noisy channel model

The Noisy Channel model:

Listener's job

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

Model of the speaker

effort of signal

Bayes' rule!

Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

The girl put out a bowl of milk for her _____

Meaning

Signal

$$p(\text{🐱} \mid \text{hat}) =$$

$$\frac{p(\text{hat} \mid \text{👑}) p(\text{👑})}{p(\text{hat})}$$

which is bigger?

$$p(\text{👑} \mid \text{hat}) =$$

$$\frac{p(\text{hat} \mid \text{👑}) p(\text{👑})}{p(\text{hat})}$$

Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

The girl put out a bowl of milk for her _____

Probabilities of meanings:

$$p(\img alt="cat face" data-bbox="108 598 158 668") = 0.99$$

$$p(\img alt="top hat" data-bbox="108 700 158 760") = 0.01$$

Implicitly, these are **conditioned** on the context, but we'll ignore this for now.

Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

$p(\text{signal} \mid \text{meaning})$: probability of pronunciation given meaning (speech error rate)

5% speech error rate:

$$p(\text{"hat"} \mid \text{h}) = 0.05 \quad p(\text{"hat"} \mid \text{H}) = 0.95$$

Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

$p(\text{signal})$:

Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

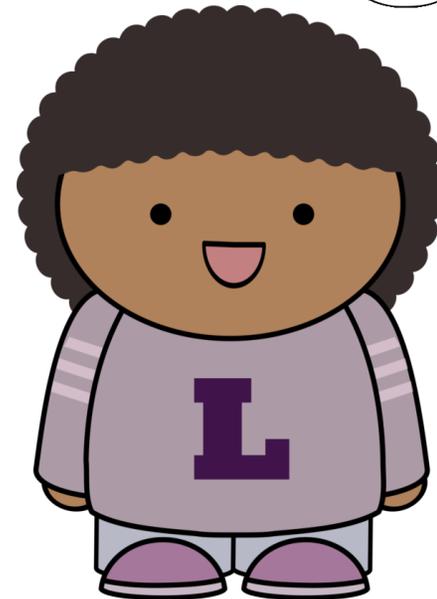
$p(\text{signal})$: we can ignore this, because we care about the relative probability of the meanings given the same signal.

Noisy channel model

The girl put out
a bowl of milk
for her hat.



hmm...



Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

$$p(\img alt="cat face emoji" data-bbox="100 540 155 610" \mid \text{"hat"}) = \frac{p(\text{hat} \mid \img alt="cat face emoji" data-bbox="440 450 495 520}) p(\img alt="cat face emoji" data-bbox="495 450 550 520})}{p(\text{hat})} = \frac{0.05 \cdot 0.99}{0.995} = 0.0495$$

$$p(\img alt="top hat emoji" data-bbox="100 730 155 800" \mid \text{"hat"}) = \frac{p(\text{hat} \mid \img alt="top hat emoji" data-bbox="460 650 515 720}) p(\img alt="top hat emoji" data-bbox="515 650 570 720})}{p(\text{hat})} = \frac{0.95 \cdot 0.01}{0.995} = 0.0095$$

Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

$$p(\img alt="cat face" data-bbox="95 468 145 535" \mid \text{"hat"}) \propto$$

$$p(\img alt="top hat" data-bbox="100 655 155 720" \mid \text{"hat"}) \propto$$

Noisy channel model

$$p(\text{meaning} \mid \text{signal}) = \frac{p(\text{signal} \mid \text{meaning})p(\text{meaning})}{p(\text{signal})}$$

$$p(\img alt="cat face emoji" data-bbox="98 465 148 535" \mid \text{"hat"}) \propto$$

$$p(\img alt="top hat emoji" data-bbox="98 688 148 758" \mid \text{"hat"}) \propto$$

In a noisy channel model, our prior belief can overcome the signal we receive.

Noisy channel model

The girl put out a bowl of milk for her _____

$$p(\text{🐱} | \text{"hat"}) \propto 0.0495$$

$$p(\text{👒} | \text{"hat"}) \propto 0.0095$$

If our intended message is 🐱, making a speech error isn't so bad— our listener will land on the correct message anyway.

Noisy channel model

The girl put out a bowl of milk for her _____

$$p(\text{🐱} | \text{"hat"}) \propto 0.0495$$

$$p(\text{🎩} | \text{"hat"}) \propto 0.0095$$

If our intended message is rare, making a speech error is bad — our listener's prior belief in the unlikeliness of the message makes it hard to communicate that message, even if we produce it perfectly.

Noisy channel model

Another reason that assigning longer words to rarer meanings makes sense is for *communicative robustness*:

a longer word is more robust to error on a single phoneme, because there are more phonemes.

Summary

- **Zipf's Law:** the frequency of a word is inversely proportional to its rank in the frequency table
- **Zipf's Hypothesis:** shorter words are used for more frequent meanings because they are more efficient.
- **Communicative efficiency:** languages evolve to express information efficiently
- **Communicative robustness:** languages evolve to express information in a noise-tolerant way