QUIZ 1: print notes non!

(if you want to ")

CS 333:

Natural Language Processing Fall 2025

Prof. Carolyn Anderson Wellesley College

Reminders

- HW I is due on Thursday
- I have help hours from 4-5 on Thursday
- Caroline has help hours today from 7:30-8:30
- Ashley has help hours tomorrow from 4-5
- All help hours are in W423

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 hypothesis:

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

Function words	
was	esopmagus
i S	gspar agus
in	mbrella
40	SIM X

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)

How likely an event is to occur.

- Probability distribution:

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

Conditional probability: p(X|Y)

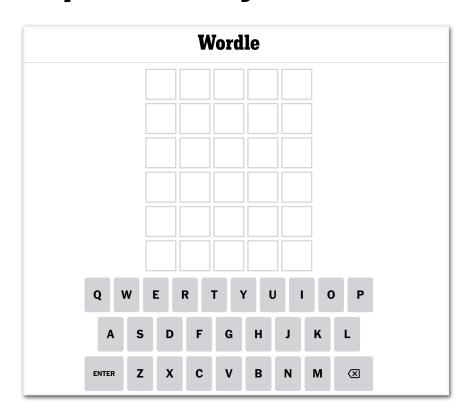
The chance of event X occurring given that event Y occurs. If events are truly independent, p(X|Y) = p(X).

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

Sample text:

"on wednesdays, we wear pink."

Total count of letters:

$$p(w) = \frac{3}{22}$$

$$p(e) = \frac{4}{22}$$

$$p(e|w) = \frac{\text{total count of e after w}}{\text{total count of w}} = \frac{3}{3}$$

$$p(w,e) = \rho(e | w) \rho(w)$$

 $3/3 \cdot 3/22 = 3/22$

Zipfian hypotheses

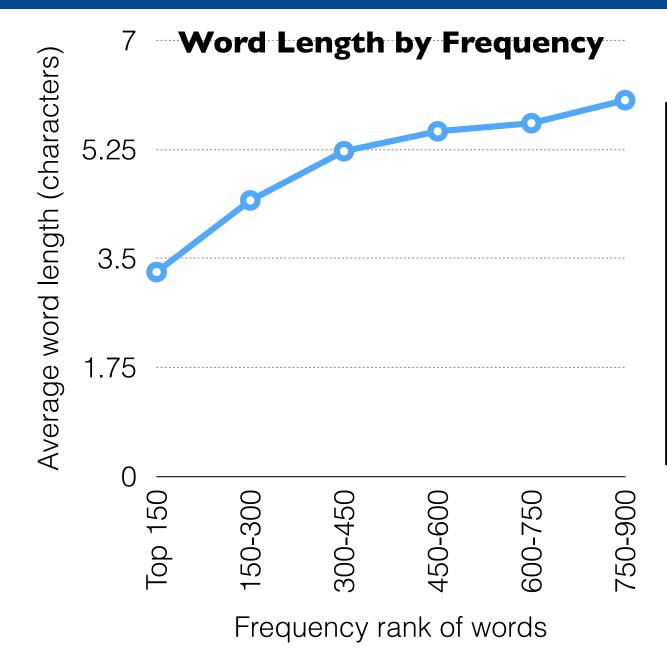
Zipf's law: The frequency of a word is inversely proportional to its frequency ranking.

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

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.

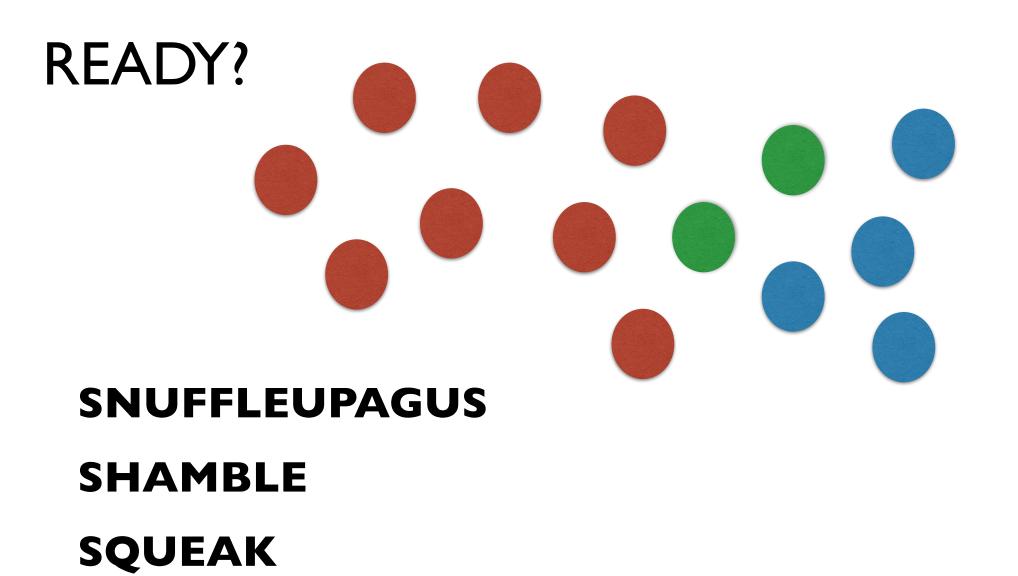
Here's the trick: the only words you can yell are **SNUFFLEUPAGUS**, **SHAMBLE**, and **SQUEAK**.

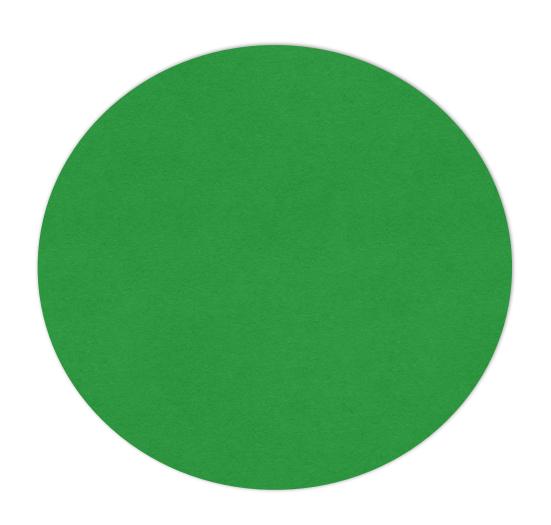
How can we code meanings efficiently?

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.

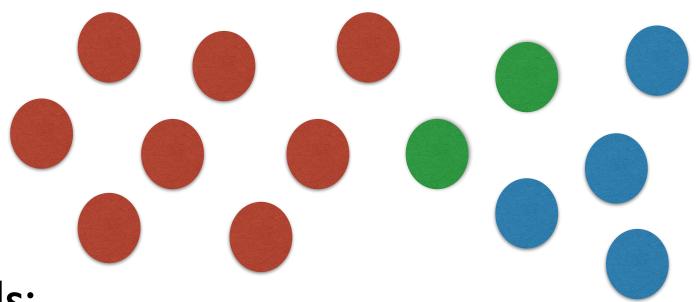
Here's the trick: the only words you can yell are **SNUFFLEUPAGUS**, **SHAMBLE**, and **SQUEAK**. And we want to do this as **fast as possible**.







Distribution of marbles:



Available words:

SNUFFLEUPAGUS SHAMBLE SQUEAK

What is the most **efficient** assignment?

```
Distribution p(red) = 8/14 = 57.1\% of marbles: p(blue) = 4/14 = 28.6\% p(green) = 2/14 = 14.3\%
```

Available words:

What is the most **efficient** assignment?

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

```
p(red) = 8/14 = 57.1\%

p(blue) = 4/14 = 28.6\%

p(green) = 2/14 = 14.3\%
```

Communicative efficiency hypothesis:

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

efficient communication.

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.

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, Seojean is standing there with an airhorn, making earsplitting noises at random intervals.

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

Why?

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?

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.

The Noisy Channel model: p(meaning | signal) = p(signal | meaning)p(meaning) p(signal)

Bayes' rule!

$$p(meaning | signal) = p(signal | meaning)p(meaning)$$

 $p(signal)$

The girl put out a bowl of milk for her _____

$$P(|| hat |) = \frac{p(hot | M) p(M)}{p(hot)}$$

$$P(|| hat |) = \frac{p(hot | M) p(M)}{p(hot)}$$

$$p(meaning | signal) = p(signal | meaning)p(meaning)$$

 $p(signal)$

The girl put out a bowl of milk for her _____

Probabilities of meanings:

$$p(3) = 0.99$$

$$p(-) = 0.01$$

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

p(meaning | signal) = p(signal | meaning)p(meaning)p(signal)

p(signal|meaning): probability of pronunciation given meaning (speech error rate)

Let's assume a 5% envol rete
$$p("nat!|\mathcal{D}) = 0.05$$

$$p("cat"|\mathcal{I}) = 0.05$$

$$p("cat"|\mathcal{D}) = 0.95$$

$$p("nat"|\mathcal{D}) = 0.95$$

```
p(meaning | signal) = p(signal | meaning)p(meaning)
p(signal)
```

p(signal|meaning): probability of pronunciation given meaning

Let's assume a 5% error rate:

$$p("hat" | \stackrel{\text{left}}{=}) = 0.05$$
 $p("hat" | \stackrel{\text{left}}{=}) = 0.95$ $p("cat" | \stackrel{\text{left}}{=}) = 0.05$

```
p(meaning | signal) = p(signal | meaning)p(meaning)
p(signal)
```

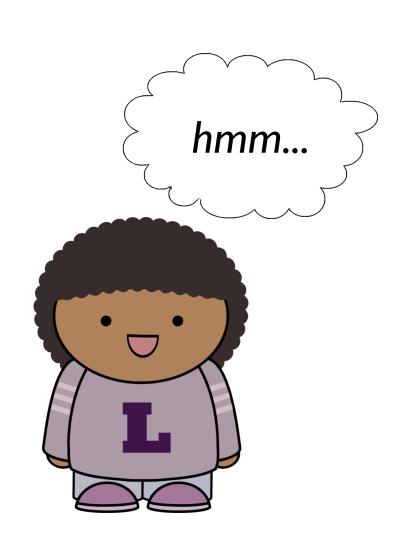
```
p(signal): ignore because me come about the relative probability
```

p(meaning | signal) = p(signal | meaning)p(meaning)p(signal)

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

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





p(meaning | signal) = p(signal | meaning)p(meaning)p(signal)

$$P(|| \text{hat}||) = P(|| \text{hat}|| || \text{mat}||) p(|| \text{mat}||)$$

$$P(|| \text{hat}||) = P(|| \text{hat}||) p(|| \text{mat}||)$$

$$P(|| \text{hat}||) = P(|| \text{hat}||) p(|| \text{mat}||)$$

```
p(meaning | signal) = p(signal | meaning)p(meaning)
p(signal)
```

$$p(\Psi)$$
 "hat") $\propto 0.05 \cdot 0.99 = 0.0495$

```
p(meaning | signal) = p(signal | meaning)p(meaning)
p(signal)
```

$$p(\mathbf{J} \mid \text{hat}') \propto (0.95)p(\mathbf{J})$$

$$p(meaning | signal) = p(signal | meaning)p(meaning)$$

 $p(signal)$

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

The girl put out a bowl of milk for her _____

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

The girl put out a bowl of milk for her _____

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.

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

Why Do We Care About This?

Understanding the information theory of natural language helps us design efficient techniques for text encoding and processing.

Next class: **TOKENIZATION**

How do we efficiently split text into useful chunks?