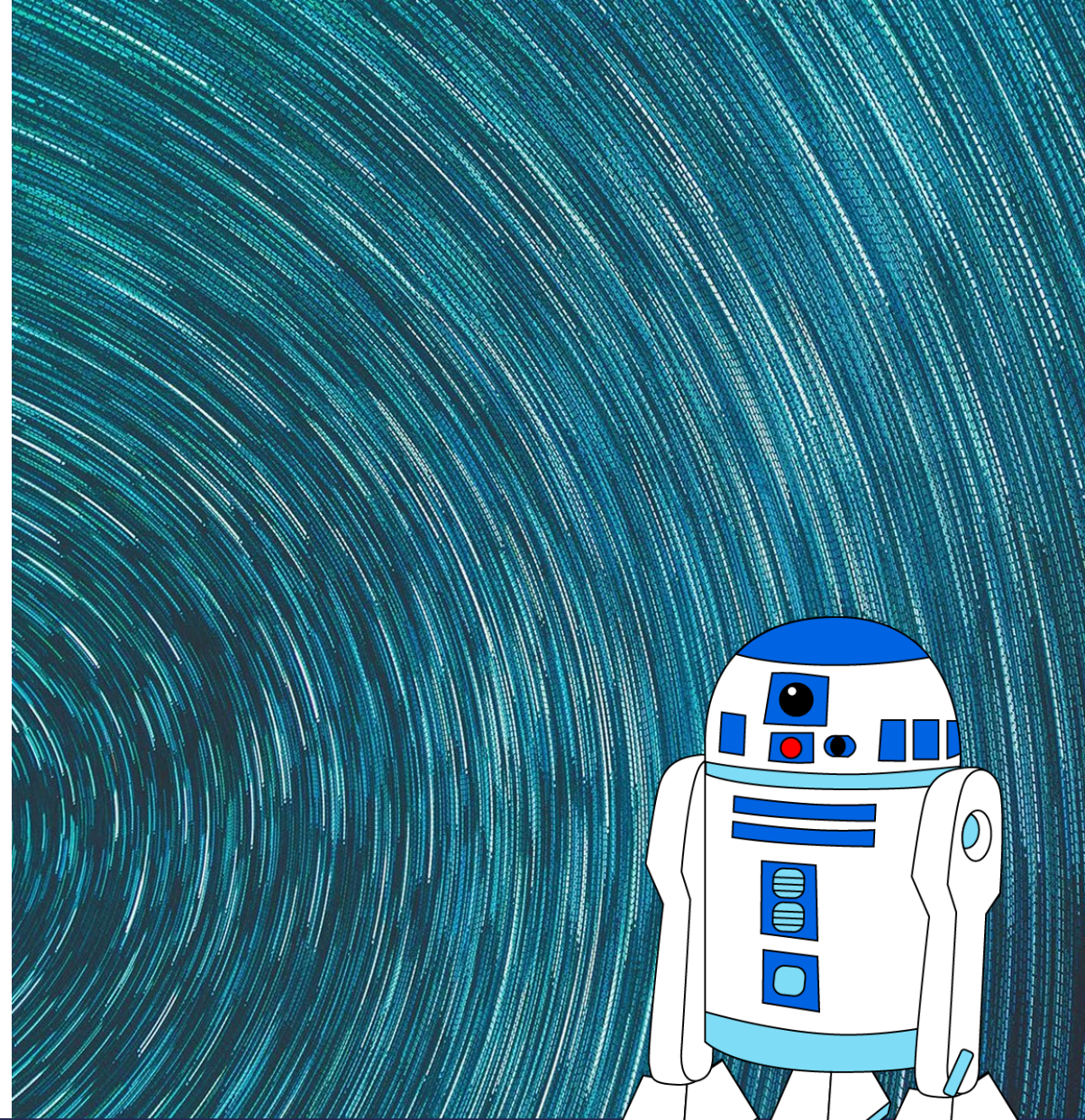


CIS 421/521:  
ARTIFICIAL INTELLIGENCE

# Probabilities and Language Models



# Uncertainty

- General situation:
  - **Observed variables (evidence):** Agent knows certain things about the state of the world (e.g., sensor readings or symptoms)
  - **Unobserved variables (states):** Agent needs to reason about other aspects (e.g. where an object is or what disease is present)
  - **Model:** Agent knows something about how the known variables relate to the unknown variables
- Probabilistic reasoning gives us a framework for managing our beliefs and knowledge



# What Probabilities Are About

- Like logical assertions, probabilities are about **possible worlds**. Instead of strictly ruling out possibilities (where a logical assertion is false), probabilities quantify **how likely** a particular possible world is.
- In probability theory, the possible worlds are called the **sample space**, and they **mutually exclusive** and **exhaustive**.
- A fully specified probability model associates a probability  **$P(w)$**  with each possible world  **$w$** .

# Random Variables

- A random variable is some aspect of the world about which we (may) have uncertainty
  - $R$  = Is it raining?
  - $U$  = Is the professor carrying an umbrella?
- We denote random variables with capital letters



# Axioms of Probability

- The probability of any possible world is between 0 and 1.
  - $0 \leq P(w) \leq 1$  for every  $w$
- The total probability of the set of all possible worlds is 1:
- $\sum_{w \in \Omega} P(w) = 1$

# Probability Distributions

- Unobserved random variables have distributions

T	P
hot	0.5
cold	0.5

W	P
sun	0.6
rain	0.1
fog	0.3

- A distribution is a TABLE of probabilities of values
- A probability (lower case value) is a single number

- Must have:  $P(W = rain) = 0.1$  and

$$\forall x \ P(X = x) \geq 0$$

Shorthand notation:

$$P(hot) = P(T = hot),$$
$$P(cold) = P(T = cold),$$
$$P(rain) = P(W = rain),$$

...

OK if all domain entries are unique

$$\sum_x P(X = x) = 1$$

# Joint Distributions

- A *joint distribution* over a set of random variables  $X_1, X_2, \dots, X_n$  specifies a real number for each assignment (or *outcome*):

$$P(X_1 = x_1, X_2 = x_2, \dots, X_n = x_n)$$

$$P(x_1, x_2, \dots, x_n)$$

- Must obey:

$$P(x_1, x_2, \dots, x_n) \geq 0$$

$$\sum_{(x_1, x_2, \dots, x_n)} P(x_1, x_2, \dots, x_n) = 1$$

- Size of distribution if n variables with domain sizes d?
  - For all but the smallest distributions, impractical to write out!

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

# Probabilistic Models

- A probabilistic model is a joint distribution over a set of random variables
- Probabilistic models:
  - (Random) variables with domains
  - Assignments are called *outcomes*
  - Joint distributions: say whether assignments (outcomes) are likely
  - *Normalized*: sum to 1.0
  - Ideally: only certain variables directly interact

Distribution over T,W

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3



# Events

- An *event* is a set  $E$  of outcomes

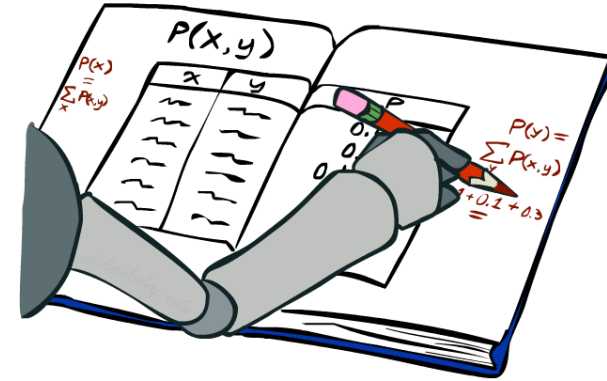
$$P(E) = \sum_{(x_1 \dots x_n) \in E} P(x_1 \dots x_n)$$

- From a joint distribution, we can calculate the probability of any event
  - Probability that it's hot AND sunny?
  - Probability that it's hot?
  - Probability that it's hot OR sunny?
- Typically, the events we care about are *partial assignments*, like  $P(T=\text{hot})$

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

# Marginal Distributions



- Marginal distributions are sub-tables which eliminate variables
- Marginalization (summing out): Combine collapsed rows by  $P(T)$  adding

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3



$$P(t) = \sum_s P(t, s)$$

T	P
hot	0.5
cold	0.5



$$P(s) = \sum_t P(t, s)$$

$P(W)$

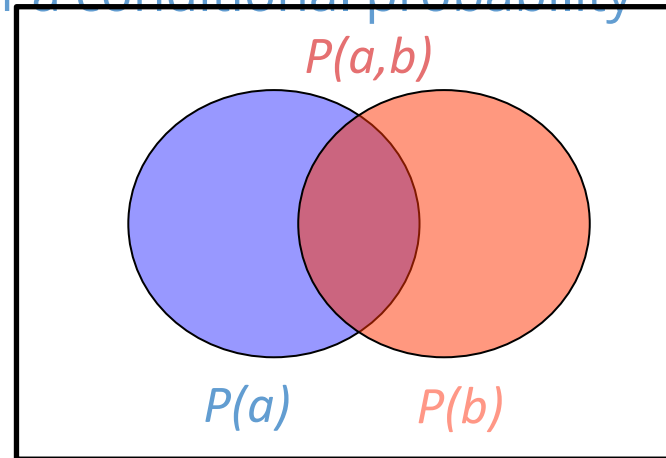
W	P
sun	0.6
rain	0.4

$$P(X_1 = x_1) = \sum_{x_2} P(X_1 = x_1, X_2 = x_2)$$

# Conditional Probabilities

- A simple relation between joint and conditional probabilities
  - In fact, this is taken as the *definition of a conditional probability*

$$P(a|b) = \frac{P(a, b)}{P(b)}$$



$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

$$P(W = s|T = c) = \frac{P(W = s, T = c)}{P(T = c)} = \frac{0.2}{0.5} = 0.4$$

$$\begin{aligned} &= P(W = s, T = c) + P(W = r, T = c) \\ &= 0.2 + 0.3 = 0.5 \end{aligned}$$

# Conditional Distributions

- Conditional distributions are probability distributions over some variables given fixed values of others

Conditional Distributions

$P(W|T)$

$P(W T = hot)$	
W	P
sun	0.8
rain	0.2

$P(W T = cold)$	
W	P
sun	0.4
rain	0.6

Joint Distribution

$P(T, W)$

T	W	P
hot	sun	0.4
hot	rain	0.1
cold	sun	0.2
cold	rain	0.3

# Probabilistic Inference

- Probabilistic inference: compute a desired probability from other known probabilities (e.g. conditional from joint)
- We generally compute conditional probabilities
  - $P(\text{on time} \mid \text{no reported accidents}) = 0.90$
  - These represent the agent's *beliefs* given the evidence
- Probabilities change with new evidence:
  - $P(\text{on time} \mid \text{no accidents, 5 a.m.}) = 0.95$
  - $P(\text{on time} \mid \text{no accidents, 5 a.m., raining}) = 0.80$
  - Observing new evidence causes *beliefs to be updated*



# Inference by Enumeration

- We want:

$$P(Q|e_1 \dots e_k)$$

- Step 3: Normalize

$$\times \frac{1}{Z}$$

$$Z = \sum_q P(Q, e_1 \dots e_k)$$

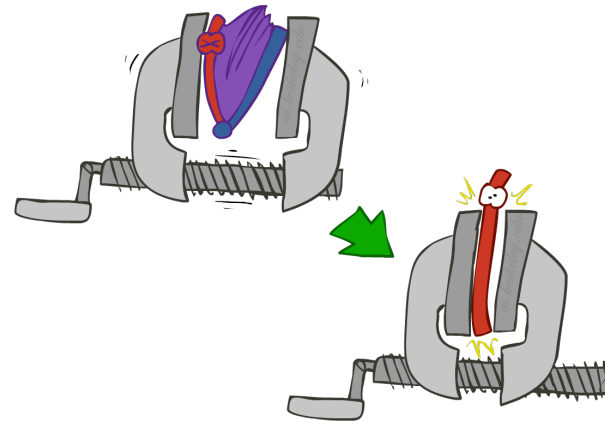
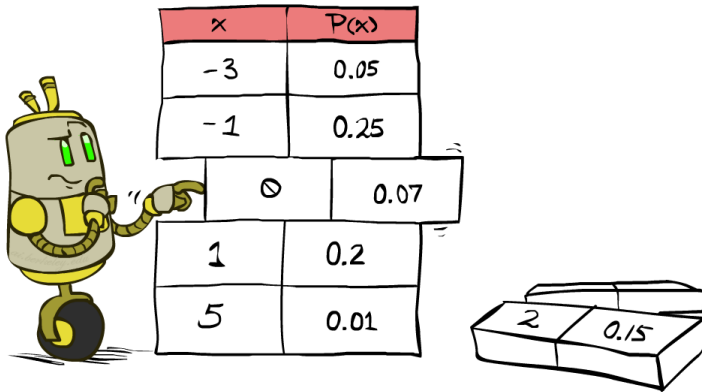
$$P(Q|e_1 \dots e_k) = \frac{1}{Z} P(Q, e_1 \dots e_k)$$

- General case:
  - Evidence variables:  $E_1 \dots E_k = e_1 \dots e_k$
  - Query\* variable:  $Q$
  - Hidden variables:  $H_1 \dots H_r$

$X_1, X_2, \dots, X_n$   
All variables

- Step 1: Select the entries consistent with the evidence

- Step 2: Sum out H to get joint of Query and evidence



$$P(Q, e_1 \dots e_k) = \sum_{h_1 \dots h_r} P(Q, \underbrace{h_1 \dots h_r}_{X_1, X_2, \dots, X_n}, e_1 \dots e_k)$$

# The Product Rule

$$P(y)P(x|y) = P(x, y)$$

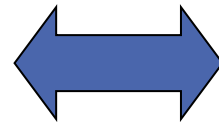
- Example:

$P(W)$

R	P
sun	0.8
rain	0.2

$P(D|W)$

D	W	P
wet	sun	0.1
dry	sun	0.9
wet	rain	0.7
dry	rain	0.3



$P(D, W)$

D	W	P
wet	sun	
dry	sun	
wet	rain	
dry	rain	

# The Chain Rule

- More generally, can always write any joint distribution as an incremental product of conditional distributions

$$P(x_1, x_2, x_3) = P(x_1)P(x_2|x_1)P(x_3|x_1, x_2)$$

$$P(x_1, x_2, \dots, x_n) = \prod_i P(x_i|x_1 \dots x_{i-1})$$

- Why is this always true?



# Bayes' Rule

- Two ways to factor a joint distribution over two variables:

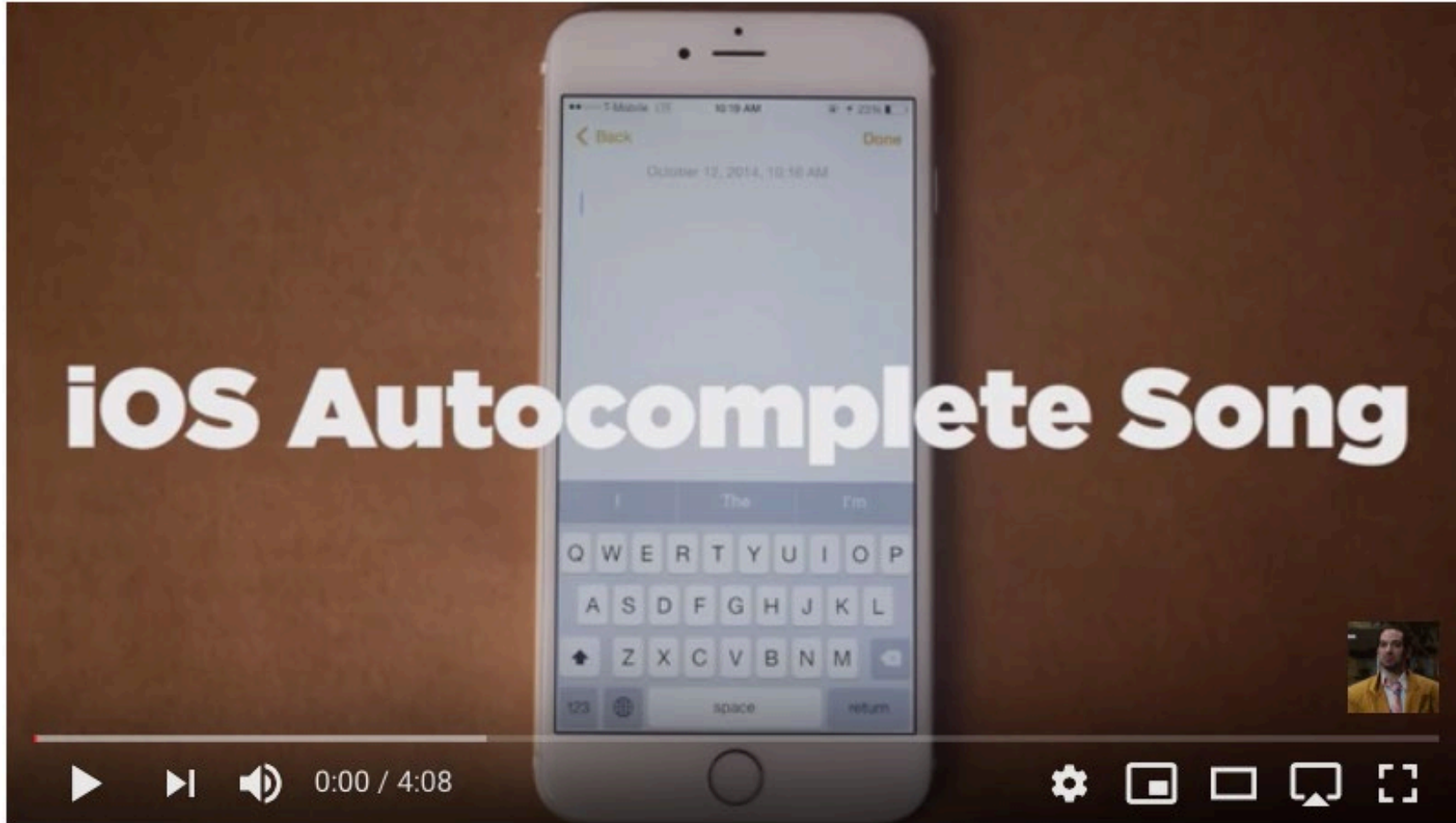
$$P(x, y) = P(x|y)P(y) = P(y|x)P(x)$$

- Dividing, we get:

$$P(x|y) = \frac{P(y|x)P(x)}{P(y)}$$

- Why is this at all helpful?
  - Lets us build one conditional from its reverse
  - Often one conditional is tricky but the other one is simple
  - Foundation of many systems we'll see later (e.g. ASR, MT)
- In the running for most important AI equation!





🎵 iOS Autocomplete Song | Song A Day #2110

<https://www.youtube.com/watch?v=M8MJFrdfGe0>

# Probabilistic Language Models

# Probabilistic Language Models

- One goal: assign a probability to a sentence
  - Autocomplete for texting
  - Machine Translation
  - Spelling Correction
  - Speech Recognition
- Other Natural Language Generation tasks: summarization, question-answering, dialog systems

# Probabilistic Language Modeling

- Goal: compute the probability of a sentence or sequence of words
- Related task: probability of an upcoming word
- A model that computes either of these is called a **language model** or **LM**

# Probabilistic Language Modeling

- Goal: compute the probability of a sentence or sequence of words

$$P(W) = P(w_1, w_2, w_3, w_4, w_5 \dots w_n)$$

- Related task: probability of an upcoming word

$$P(w_5 | w_1, w_2, w_3, w_4)$$

- A model that computes either of these

$P(W)$  or  $P(w_n | w_1, w_2 \dots w_{n-1})$  is called a **language model**.

- Better: **the grammar** But **language model** or **LM** is standard

# How to compute $P(W)$

- How to compute this joint probability:

$P(\text{the, underdog, Philadelphia, Eagles, won})$

- Intuition: let's rely on the Chain Rule of Probability

# The Chain Rule





# The Chain Rule

- Recall the definition of conditional probabilities

$$p(\mathbf{B} | \mathbf{A}) = \mathbf{P}(\mathbf{A}, \mathbf{B}) / \mathbf{P}(\mathbf{A}) \quad \text{Rewriting: } \mathbf{P}(\mathbf{A}, \mathbf{B}) = \mathbf{P}(\mathbf{A})\mathbf{P}(\mathbf{B} | \mathbf{A})$$

- More variables:

$$\mathbf{P}(\mathbf{A}, \mathbf{B}, \mathbf{C}, \mathbf{D}) = \mathbf{P}(\mathbf{A})\mathbf{P}(\mathbf{B} | \mathbf{A})\mathbf{P}(\mathbf{C} | \mathbf{A}, \mathbf{B})\mathbf{P}(\mathbf{D} | \mathbf{A}, \mathbf{B}, \mathbf{C})$$

- The Chain Rule in General

$$\mathbf{P}(x_1, x_2, x_3, \dots, x_n) = \mathbf{P}(x_1)\mathbf{P}(x_2 | x_1)\mathbf{P}(x_3 | x_1, x_2) \dots \mathbf{P}(x_n | x_1, \dots, x_{n-1})$$

# Joint probability of words in sentence



# Joint probability of words in sentence

$$P(w_1 w_2 \square \dots w_n) = \prod_i P(w_i | w_1 w_2 \square \dots w_{i-1})$$

P(“the underdog Philadelphia Eagles won”) =

P(the) ×  
P(underdog | the) ×  
P(Philadelphia | the underdog) ×  
P(Eagles | the underdog Philadelphia) ×  
P(won | the underdog Philadelphia Eagles)

# How to estimate these probabilities

- Could we just count and divide?

# How to estimate these probabilities

- Could we just count and divide? Maximum likelihood estimation (MLE)

$$P(\text{won} | \text{the underdog team}) = \frac{\text{Count}(\text{the underdog team won})}{\text{Count}(\text{the underdog team})}$$

- Why doesn't this work?

# Simplifying Assumption = Markov Assumption

# Simplifying Assumption = Markov Assumption

- $P(\text{won} | \text{the underdog team}) \approx P(\text{won} | \text{team})$
- Only depends on the previous  $k$  words, not the whole context
- $\approx P(\text{won} | \text{underdog team})$
- $\approx P(w_i | w_{i-2} w_{i-1})$
- $P(w_1 w_2 w_3 w_4 \dots w_n) \approx \prod_i^n P(w_i | w_{i-k} \dots w_{i-1})$
- $k$  is the number of context words that we take into account

# How much history should we use?

unigram	<b>no history</b>	$\prod_i^n p(w_i)$	$p(w_i) = \frac{\text{count}(w_i)}{\text{all words}}$
bigram	1 word as history	$\prod_i^n p(w_i w_{i-1})$	$p(w_i w_{i-1}) = \frac{\text{count}(w_{i-1}w_i)}{\text{count}(w_{i-1})}$
trigram	2 words as history	$\prod_i^n p(w_i w_{i-2}w_{i-1})$	$\begin{aligned} p(w_i w_{i-2}w_{i-1}) \\ &= \frac{\text{count}(w_{i-2}w_{i-1}w_i)}{\text{count}(w_{i-2}w_{i-1})} \end{aligned}$
4-gram	3 words as history	$\prod_i^n p(w_i w_{i-3}w_{i-2}w_{i-1})$	$\begin{aligned} p(w_i w_{i-3}w_{i-2}w_{i-1}) \\ &= \frac{\text{count}(w_{i-3}w_{i-2}w_{i-1}w_i)}{\text{count}(w_{i-3}w_{i-2}w_{i-1})} \end{aligned}$



# Simplest case: Unigram model

$$P(w_1 w_2 \square \dots w_n) \approx \prod_i P(w_i)$$

Some automatically generated sentences from a unigram model

fifth an of futures the an incorporated a a the inflation most  
dollars quarter in is mass

thrift did eighty said hard 'm july bullish

that or limited the

# Bigram model

- Condition on the previous word:

$$P(w_i | w_1 w_2 \dots w_{i-1}) \approx P(w_i | w_{i-1})$$

texaco rose one in this issue is pursuing growth in a boiler house said  
mr. gurria mexico 's motion control proposal

without permission from five hundred fifty five yen

outside new car parking lot of the agreement reached

this would be a record november

# N-gram models

- We can extend to trigrams, 4-grams, 5-grams
- In general this is an insufficient model of language
  - because language has **long-distance dependencies**:

“The computer(s) which I had just put into the machine room on the fifth floor is (are) crashing.”

- But we can often get away with N-gram models

# Language Modeling

Estimating N-gram Probabilities

# Estimating bigram probabilities

- The Maximum Likelihood Estimate

$$P(w_i | w_{i-1}) = \frac{\textit{count}(w_{i-1}, w_i)}{\textit{count}(w_{i-1})}$$

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

# An example

$$P(w_i | w_{i-1}) = \frac{c(w_{i-1}, w_i)}{c(w_{i-1})}$$

<s> I am Sam </s>

<s> Sam I am </s>

<s> I do not like green eggs and ham </s>

$$P(\text{I} | \text{<s>}) = \frac{2}{3} = .67$$

$$P(\text{Sam} | \text{<s>}) = \frac{1}{3} = .33$$

$$P(\text{am} | \text{I}) = \frac{2}{3} = .67$$

$$P(\text{</s>} | \text{Sam}) = \frac{1}{2} = 0.5$$

$$P(\text{Sam} | \text{am}) = \frac{1}{2} = .5$$

$$P(\text{do} | \text{I}) = \frac{1}{3} = .33$$

# Problems for MLE

Train	Test
denied the allegations	denied the memo
denied the reports	
denied the claims	
denied the requests	

- $P(\text{memo} | \text{denied the}) = 0$
- And we also assign 0 probability to all sentences containing it!

# Problems for MLE

- Out of vocabulary items (OOV)
- `<unk>` to deal with OOVs
- Fixed lexicon  $L$  of size  $V$
- Normalize training data by replacing any word not in  $L$  with `<unk>`
  
- Avoid zeros with smoothing



# Practical Issues

- We do everything in log space
  - Avoid underflow
  - (also adding is faster than multiplying)

$$\log(p_1 \times p_2 \times p_3 \times p_4) = \log p_1 + \log p_2 + \log p_3 + \log p_4$$

# Google N-Gram Release, August 2006

AUG

3

## All Our N-gram are Belong to You

Posted by Alex Franz and Thorsten Brants, Google Machine Translation Team

Here at Google Research we have been using word [n-gram models](#) for a variety of R&D projects,

...

That's why we decided to share this enormous dataset with everyone. We processed 1,024,908,267,229 words of running text and are publishing the counts for all 1,176,470,663 five-word sequences that appear at least 40 times. There are 13,588,391 unique words, after discarding words that appear less than 200 times.

# Google N-Gram Release

- o serve as the incoming 92
- o serve as the incubator 99
- o serve as the independent 794
- o serve as the index 223
- o serve as the indication 72
- o serve as the indicator 120
- o serve as the indicators 45
- o serve as the indispensable 111
- o serve as the indispensable 40
- o serve as the individual 234

<http://googleresearch.blogspot.com/2006/08/all-our-n-gram-are-belong-to-you.html>