

# Natural Language Processing

## CSE 325/425



Sihong Xie

### Lecture 2:

- Language models
- MLE and Laplacian

# Sample space, events, random variable

- NLP is about using probability to model languages, and using optimization to find the right model.
- The universal set  $\Omega$  contains all objects that can occur in a specific context.
  - Discrete:  $\Omega$  = All words in a corpus (vocabulary)
  - Continuous:  $\Omega$  = Sentimental score in  $[-1, 1]$
- Event: subsets of  $\Omega$  ( $A \subseteq \Omega$ )
  - A document: bag of words (subset of the vocabulary)
  - Positive sentimental score  $[0, 1]$
- Random variable  $X : \Omega \rightarrow \text{set of numbers}$ 
  - e.g., map a word to an integer index.
  - The sentimental score is a number already.
- Probability distribution of a random variable  $P(X = 1) = P(\omega \in \Omega : X(\omega) = 1)$
- Probability of an event  $A \subseteq \Omega$  :

$$P(A) \rightarrow [0, 1]$$

$|V|$  words in the vocabulary  $V$   
 $m$  words

Random variable  
 $\downarrow$   
 $P(\text{a doc}) = P(w_1, w_2, \dots, w_m)$   
 $m$  words in doc

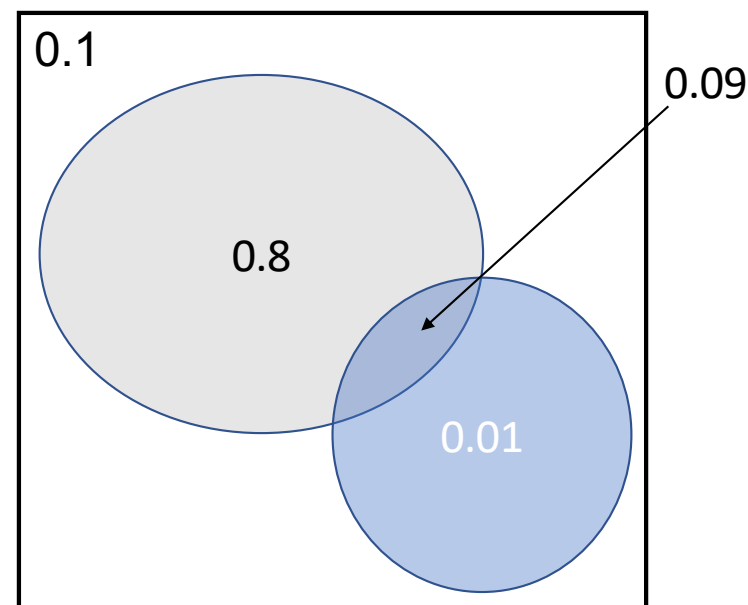
## Joint Distribution

$|V|^m$  Comb

Recipe for making a joint distribution of  $m$  Rvs

1. List **all** possible combinations of values of the RVs.
2. Assign **valid** probability to each combination.

A	B	P(A, B)	
T	T	0.09	$=P(A \cap B)$
T	F	0.8	$=P(A \cap B^c)$
F	T	0.01	$=P(B \cap A^c)$
F	F	0.1	$=P(B^c \cap A^c)$



$$P(A) = P(A \cap B) + P(A \cap B^c)$$

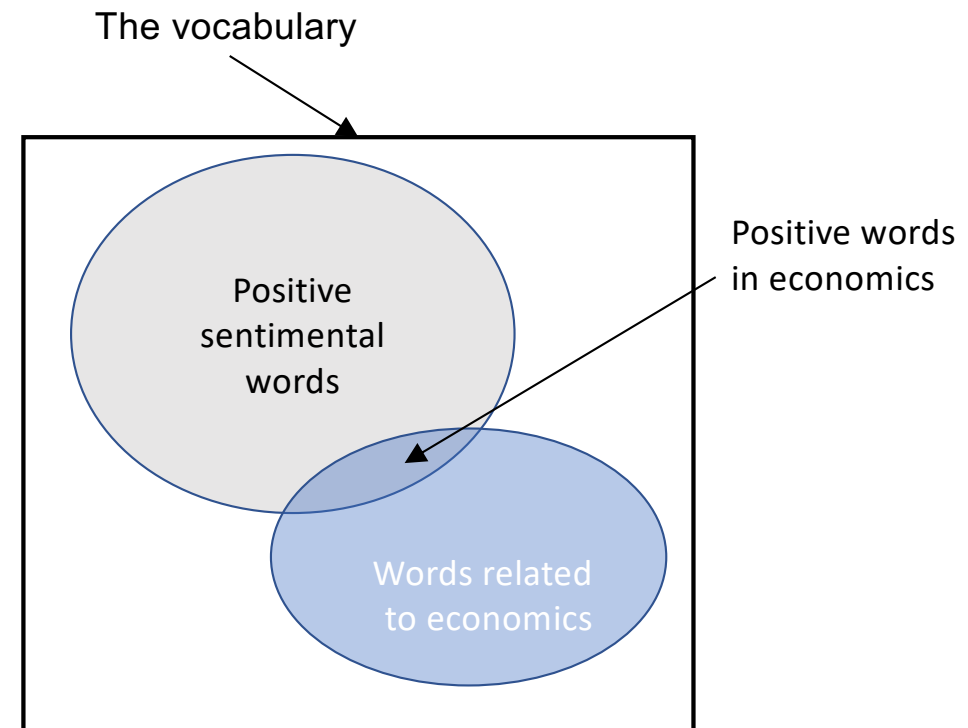
$$P(B|A) = \frac{P(A \cap B)}{P(A)}$$

Venn Diagram

# Axioms

## The Axioms of Probability

- $P(\Omega) = 1$
- $0 \leq P(A)$   $A \subseteq \Omega$
- $P(A \cup B) = P(A) + P(B) - P(A \cap B)$



# Useful theorems

Theorems:

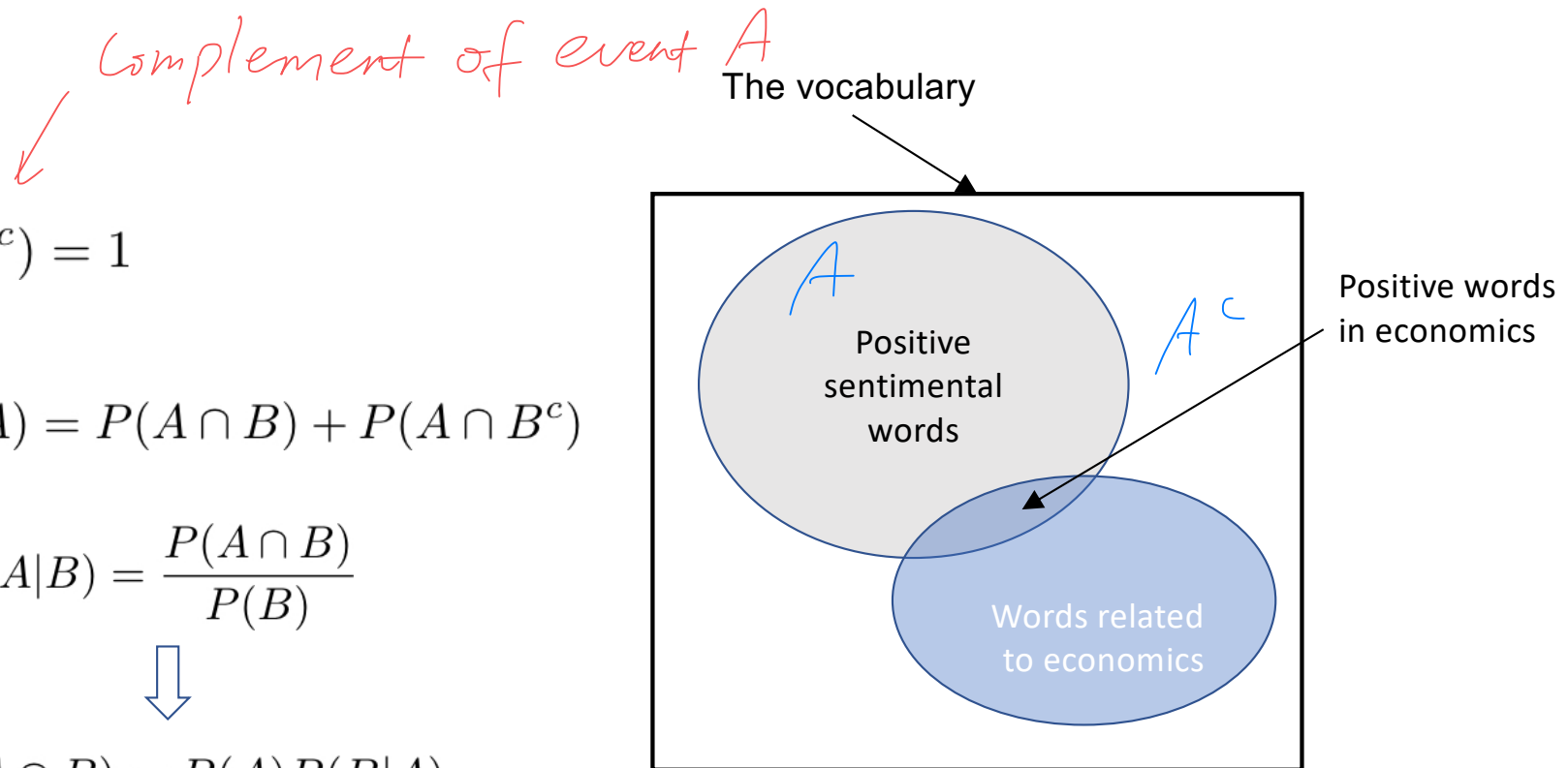
- $P(A) + P(A^c) = 1$
- Total Prob  $P(A) = P(A \cap B) + P(A \cap B^c)$
- Conditional:  $P(A|B) = \frac{P(A \cap B)}{P(B)}$

Definition



$$P(A \cap B) = P(A)P(B|A)$$

$$= P(B)P(A|B)$$



# Useful theorems

Bayes Rule:

- $P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{P(A|B)P(B)}{P(A)}$

- Prior:  $P(B)$

- Likelihood:  $P(A|B)$

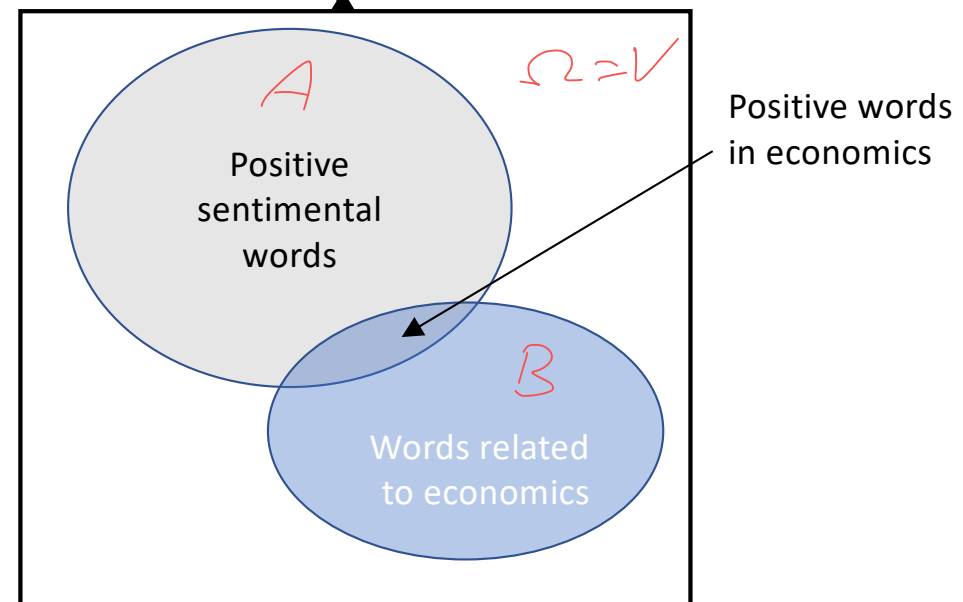
- Posterior:  $P(B|A)$

- Marginal:  $P(A)$

By definition of conditional Prob

By definition of Joint Prob.

The vocabulary



# Useful theorems

Bayes Rule:

- $$P(B|A) = \frac{P(A \cap B)}{P(A)} = \frac{P(A|B)P(B)}{P(A)}$$



$$P(B|A) = \frac{P(A|B)P(B)}{P(A|B)P(B) + P(A|B^c)P(B^c)}$$

$$B \cap B^c = \emptyset$$

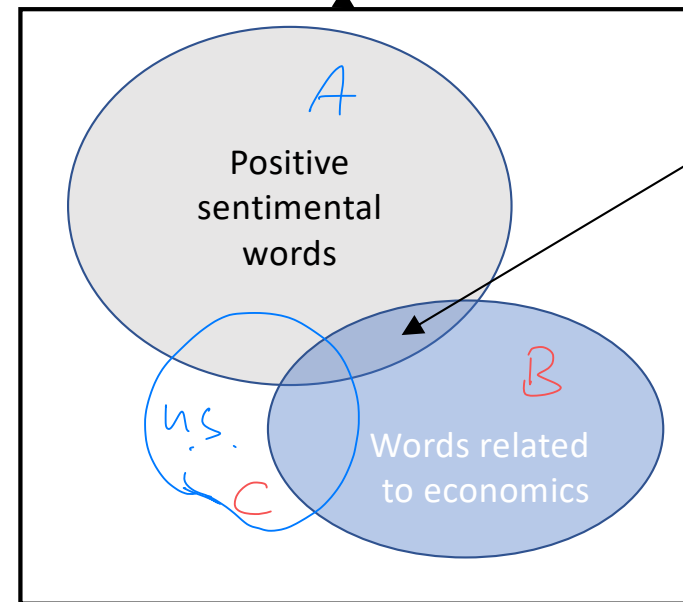
$$P(B|A \cap C) = \frac{P(A|B \cap C)P(B \cap C)}{P(A \cap C)}$$

↑

$$= \frac{P(C|B \cap A)P(A \cap B)}{P(A \cap C)}$$

$$\neq P(A \cap C | B)$$

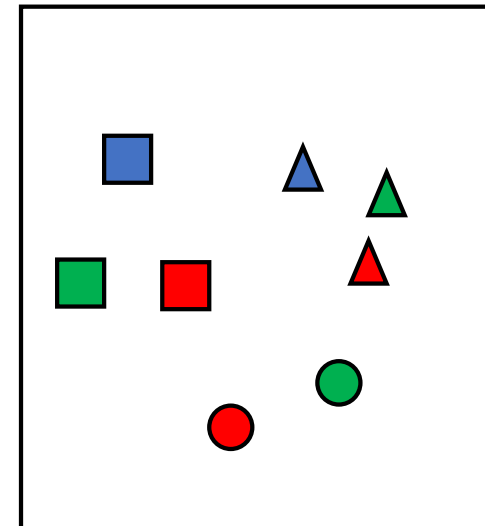
The vocabulary



Positive words  
in economics

# More exercises

Sample space  $\Omega$



- Joint probability

- $P(\text{red } \square) = 1/8$

- Conditional probability

- $P(\text{red} | \square) = 1/3$

- $P(\square | \text{red}) = 1/3$  1/4

- Total probability

- $P(\text{red}) = P(\text{red} | \square) \times P(\square) + P(\text{red} | \triangle) \times P(\triangle) + P(\text{red} | \circ) \times P(\circ)$

- Bayes rule

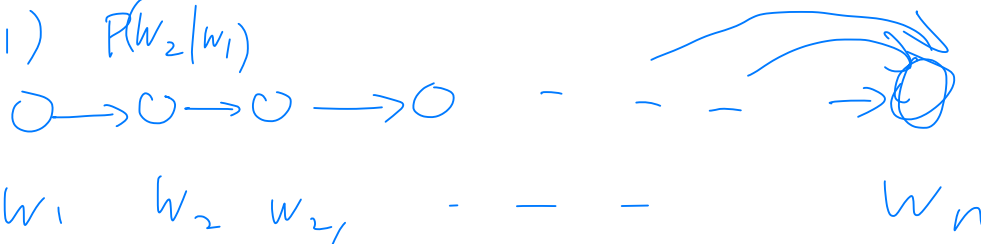
$$\frac{1/3}{3/8} \quad \frac{1/3}{3/8} \quad \frac{1/2}{2/8}$$

- $P(\text{red} | \square) = P(\text{red} \square) / P(\square)$

$$\frac{1/8}{3/8} = 1/3$$



$$P(w_1) \quad P(w_2|w_1)$$

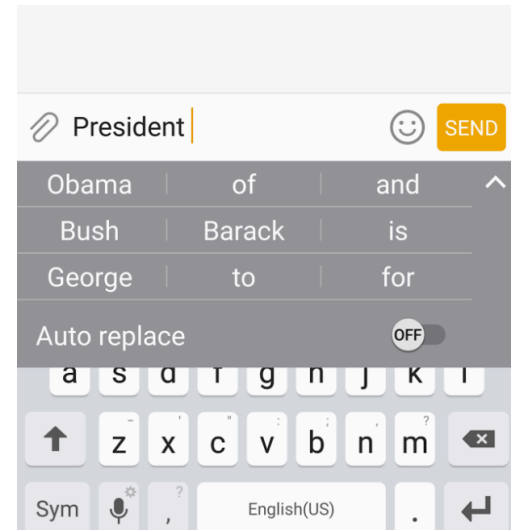


Markov chain

# Language models

## Language models

- Define the probability of a sequence of words (e.g., a sentence).
- $P(w_1, \dots, w_n)$  ?
- Can be used for spell-checking:  $P(\text{"the book"}) \gg P(\text{"book the"})$
- Can be used to predict the next word:  $P(\text{"Obama"}|\text{"President"}) \gg P(\text{"book"}|\text{"President"})$  ~~Bi-gram~~
- Many different ways to interpret this probability:



- Uni-gram:  $P(w_1)P(w_2) \dots P(w_n)$
- Bi-gram:  $P(w_1)P(w_2|w_1) \dots P(w_n|w_{n-1})$
- Tri-gram:  $P(w_1)P(w_2|w_1)P(w_3|w_1w_2) \dots P(w_n|w_{n-1}w_{n-2}) = \frac{P(w_nw_{n-1}w_{n-2})}{P(w_{n-1}w_{n-2})}$
- Where do these probabilities come from?

won't be able to  
find long-range dependency

# Language models

## Language models training

- Training corpus.

- Uni-gram: relative frequency of each word  $P(w) = \frac{\text{Count of } w}{\text{Count of the total tokens}}$

- Bi-gram: conditional probability of a word given another word  $P(w|v) = \frac{\text{Count of } (vw)}{\text{Count of } v}$

- Example

- Estimate  $P(\text{Sam})$  and  $P(\text{I}|\text{do})$

The corpus = {

$D1 = [<s> \text{ I am Sam } </s>]$  = length 5

$D2 = [<s> \text{ Sam I am } </s>]$  = length 5

$D3 = [<s> \text{ I do not like green eggs and ham } </s>]$

}

= length 10

$$P(\text{Sam}) = \frac{2}{20} = 1/10$$

$$P(\text{I}|\text{do}) = \frac{0}{10} = 0$$

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

# Language models

## Common issues in language models

- Time/space complexity: n-gram has  $O(|V|^n)$  complexity.
- Sample complexity: how many data are needed to estimate a quantity.
  - Data sparsity: corpus is not enough to reliably estimate the probabilities.
  - Recall the large-number-theorem.
  - Flipping a coin once and estimate  $P(head)$

$|V|$  : size of the vocab  $V$

$P(Sam | I)$

# Laplacian smoothing

- Missing words and phrases (in training corpus)

- Laplacian smoothing for unigram

$$0 \leq P_{\text{Laplacian}}(w) = \frac{(\text{Count of } w) + 1}{(\text{Count of the total tokens}) + |V|}$$

$$HW: \sum_{w \in V} P_{\text{Lap}}(w)$$

$$= 1$$

$$\leq 1$$

Size of the vocab  $V$

- Verify that the smoothed estimation is a distribution over the vocabulary.
  - Laplacian smoothing for bi-gram

$$P_{\text{Laplacian}}(w_2|w_1) = \frac{(\text{Count of } w_1 w_2) + 1}{(\text{Count of } w_1) + |V|}$$

- Use the following corpus to make Laplacian estimation.

The corpus = {

D1=[<s> I am Sam </s>]

D2=[<s> Sam I am </s>]

D3=[<s> I do not like green eggs and ham </s>]

{ you  $\in V$  }

$$P_{\text{Lap}}("you") = \frac{0+1}{20+|V|}$$