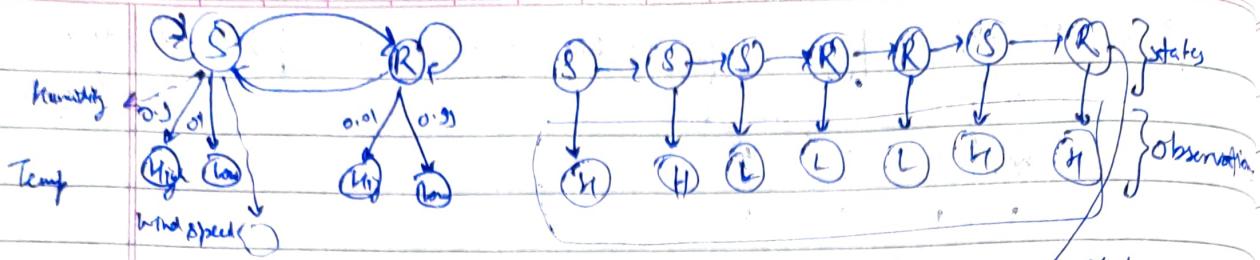


## Hidden Markov Model

Date \_\_\_\_\_  
DELTA Pg No. \_\_\_\_\_



- Transition probability table

- Prior values

- Observation (on sunny day  $P(H|S) = 0.9$ ,  $P(L|S) = 0.1$ )  
 (on rainy day  $P(H|R) = 0.01$ ,  $P(L|R) = 0.99$ )

(states)	S	H	L	(obs)
S	0.9	0.1		
R	0.01	0.99		

Predict next state  
 $\rightarrow p_{predict}$   
 $\frac{1}{2}$

$$P(S|SSSRERHHLLMH) = ?$$

$$P(R|SSSERRSHHLLLM) = ?$$

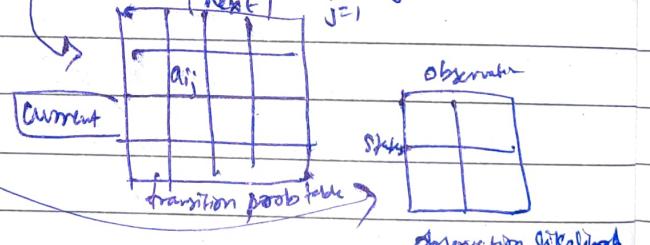
Ans: An HMM is specified by the following components

(Set of  $N$  states)  $Q = q_1, q_2, \dots, q_N$        $A = a_{11} a_{12} \dots a_{nn}$

(Sequence of  $T$  observations)  $O = O_1, O_2, \dots, O_T$

transition prob matrix  $A$  each  $a_{ij}$

represents of moving from state  $i$  to state  $j$  such that  $\sum_{j=1}^n a_{ij} = 1 \quad \forall i$



priors (initial probability distribution) over states  $\pi = \pi_1, \pi_2, \dots, \pi_N$

a set  $Q_A \subseteq Q$  of legal accepting states  $Q_A = \{q_1, q_2, \dots\}$

markov assumption:-  $P(q_i | q_1, q_2, \dots, q_{i-1}) = P(q_i | q_{i-1})$

output Assumption:-  $P(O_i | q_1, q_2, \dots, q_{i-1}, O_1, \dots, O_{i-1}) = P(O_i | q_i)$

problem 1 (likelihood): Given an HMM  $\lambda = (A, B)$  and an observation seq  $O$ , determine the likelihood of  $P(O | \lambda)$

Problem 2 (decoding): Given an observation sequence  $O$  and HMM  $\lambda = (A, B)$  discover the best state sequence  $Q$ .

Problem 3 (learning): Given an observation sequence  $O$  and the set of states  $Q$ , learn the HMM parameter  $A$  and  $B$

Problem 1 (likelihood) Given: HMM =  $\lambda = (A, B)$   
 transition observation  
 $[A] [B]$  find likelihood  $P(O|\lambda)$  **DELTA** Pg No.

$$P(3, 1, 3 | \lambda) = ?$$

no of ice cream <sup>odd</sup> = 0 (observation sequence)  
 3, 1, 3

Problem 2 (Decoding) Given: HMM =  $\lambda = (A, B)$   
 Observation seq  $O$ ,  
 discover best hidden state seq  $Q$

$$O \rightarrow 3, 1, 3$$

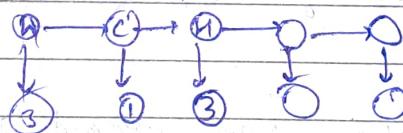
best  $P(\text{States } | 3, 1, 3)$

$$\left. \begin{array}{l} P(\text{hot hot hot } | 3, 1, 3) \\ P(\text{hot hot cold } | 3, 1, 3) \\ P(\text{cold cold cold } | 3, 1, 3) \end{array} \right\} \text{2 patterns}$$

Problem 3 (learning of HMM) given Observation seq  $O$  and set of states.

learn HMM parameters A and B  
 $S \rightarrow \text{hot cold hot} \quad \text{--- cold, hot cold hot hot}$

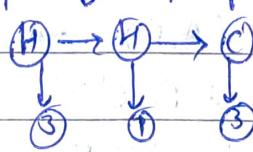
$$O \rightarrow 3, 1, 3 \quad \text{--- } 1, 6, 1, 2, 8, 9$$



$$[A?] [B?]$$

problem 2 given  $\lambda = (A, B)$ ,  $O = 3, 1, 3$   
 $P(O|A) = ?$  take state  $O$

Assumption suppose corresponding seq of state is hot hot cold



	H	C	G
H	1	0	0
C	0	2	0
G	0	0	3

Prob ~~of~~ ~~prob~~ (3, 1, 3 | HMC),

$$P(3, 1, 3 | \text{HMC}) = P(3|H) \cdot P(1|H) \cdot P(3|C)$$

Since we don't know the hidden state seq, we need to compute the probability of  $3,1,3$  by summing over all possible ~~other~~<sup>state</sup> sequences weighted by their probability.

$$P(O, Q) = P(O|Q) \cdot P(Q) = \prod_{i=1}^N P(O_i|q_i) \cdot \prod_{i=1}^N P(q_i|q_{i-1})$$

$$P(3,1,3 | \text{hot hot cold}) = P(H|H) \cdot P(H|H) \cdot P(C|H) \leftarrow \\ P(S|H) \cdot P(C|H) \cdot P(S|C) \leftarrow$$

$$P(O) = \sum_Q P(O|Q)$$

	1	→ $P(O=1)$
0,	2	→ $P(O=2)$
	3	→ $P(O=3)$

$$P(O) = \sum_Q P(O|Q) \cdot P(Q)$$

$$P(3,1,3) = P(3,1,3 | \text{hot hot hot}) \cdot P(\text{hot hot hot}) + \\ P(3,1,3 | \text{x x x}) \cdot P(\text{x x x}) + \\ P(3,1,3 | \text{y y y}) \cdot P(\text{y y y}) + \dots \quad (\text{8 sequences})$$

This is not practically feasible  
 (Computationally)  
 direct method

So it's dynamic programming based forward algorithm.

Cee16

Likelihood computation using forward algorithm computing recursively the probabilities of observations, one-by-one it proceeds to the states and calculates the obs probabilities with the help of ~~prev state~~ observation probabilities

Given an HMM  $Y = (A, B)$ , and observation seq  $O$

determine the likelihood  $P(O|Y)$

Eg: given an ice cream eating HMM, what is probability of 3,1,3  
 icecreams are sold.

The problem with hidden states is that we don't know the underlying seq of states.

day 1 = 3

observation seq.

Given  $Q = q_0, q_1, \dots, q_T$  and  $O = o_1, o_2, \dots, o_T$

PL the likelihood of observations gives states

$$P(O|Q) = \prod_{i=1}^T P(o_i|q_i)$$

$$\text{eg } P(3|3, \text{hot hot cold}) \rightarrow P(3|\text{hot}) \cdot P(1|\text{hot}) \cdot P(3|\text{cold})$$

$$P(O|Q) = P(O|Q) \cdot P(Q) = \prod_{i=1}^T P(o_i|q_i) \times \prod_{i=1}^T P(q_i|q_{i-1})$$

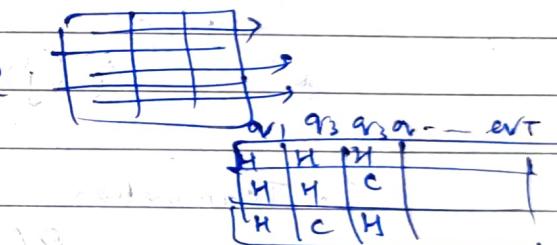
$$\text{eg } P(3|3, \text{hot hot cold}) = P(\text{hot}|\text{start}) \cdot P(\text{hot}|\text{hot}) \cdot P(\text{cold}|\text{hot}) \\ \cdot P(3|\text{hot}) \cdot P(1|\text{hot}) \cdot P(3|\text{cold})$$

W.K.t.hat  $P(A) = \sum_B P(A|B)$

$$P(O) = \sum_Q P(O|Q)$$

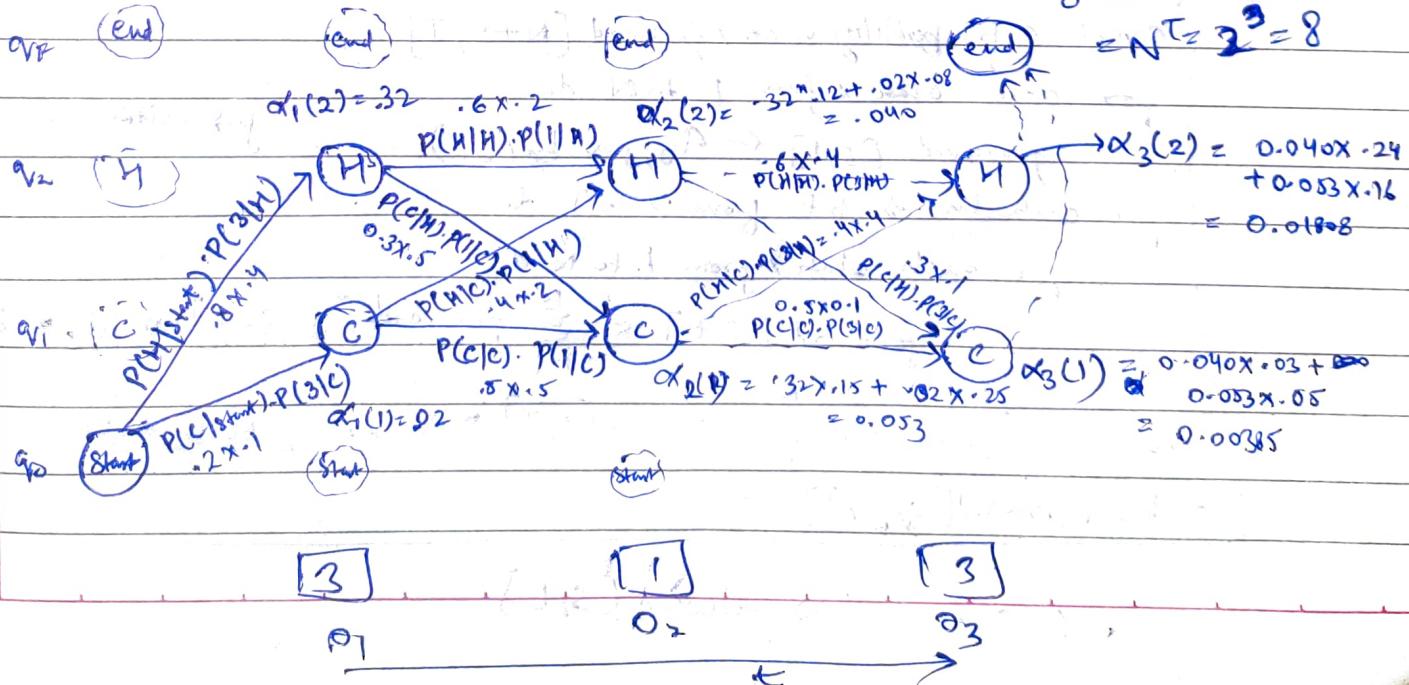
$$= \sum_Q P(o_i|Q) \cdot P(Q)$$

$$= P(O_1|Q_1) \cdot P(Q_1) \cdot \dots \cdot P(O_T|Q_T) \cdot P(Q_T)$$



$$P(3, 1, 3) = P(3|3, \text{hot hot hot}) + P(3|3, \text{cold, cold, cold})$$

8 combination



$\alpha_t(j)$  = represents the probability of being in state  $j$  after seeing  $t$  observations, given  $Y$  (probability table for markov model)

$$\alpha_t(j) = P(o_1, o_2, \dots, o_t | v_{j=1}^t | \lambda)$$

↳ the  $t^{th}$  state in the seq of states is state  $j$

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) [a_{ij} b_j(o_t)]$$

↓                    ↓  
 transition prob. observation probability

### 1. Initialization

$$\alpha_1(j) = a_{0j} b_j(o_1) \quad i \leq j \leq N$$

### 2. Recursion

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i) a_{ij} b_j(o_t) \quad 1 \leq t \leq T \quad 1 \leq j \leq N$$

### 3. Termination

$$P(O|Y) = \alpha_T(q_F) = \sum_{i=1}^N \alpha_T(i) q_i F$$

function Forward (observations lenT, state-graph lenN) returns forward prob

Create a probability matrix forward [N+2, T]

for each state  $s$  from 1 to N do

$$\text{forward}[s, 1] \leftarrow a_{0,s} * b_s(o_1)$$

    ] initialize

for each time step  $t$  from 2 to T do

    for each state  $s$  from 1 to N do

$$\text{forward}[s, t] \leftarrow \sum_{i=1}^N \text{forward}[s', t-1] * a_{s,i} * b_s(o_t)$$

$$\text{forward}[q_F, T] \leftarrow \sum_{i=1}^N \text{forward}[s, T] * q_i F \quad ] \text{ termination}$$

return  $\text{forward}[q_F, T]$

$$\alpha_t(s) \longrightarrow \text{forward}[s, t]$$

# Lecture 18

## Vector Semantics

Date: \_\_\_\_\_ DELTA Pg No. \_\_\_\_\_

We define a word as a vector called an "embedding" because it's embedded into a space.

The standard way to represent meaning in NLP

Fine-grained model of meaning for similarity.

- NLP tasks like sentiment analysis

- with words, requires same word to be in training and test.

- with embedding is ok if similar words occurred.

~~words~~

Term-document matrix (Each document is represented by a

Documents or context	→	D <sub>1</sub> As you like it	D <sub>2</sub> Twelfth Night	D <sub>3</sub> Julius Caesar	D <sub>4</sub> Henry V
battle	←	1	0	7	13
good	←	14	80	62	89
fool	←	36	58	1	4
wit	←	20	15	2	3

when we represent doc as count of words we call them Document Vectors

$$\text{battle} = [1, 0, 7, 13]$$

word  
vector

All words present in D<sub>1</sub> are known as content of ~~battle~~ battle

N content  
N words

word cooccurrence matrix
--------------------------

or  
word word matrix

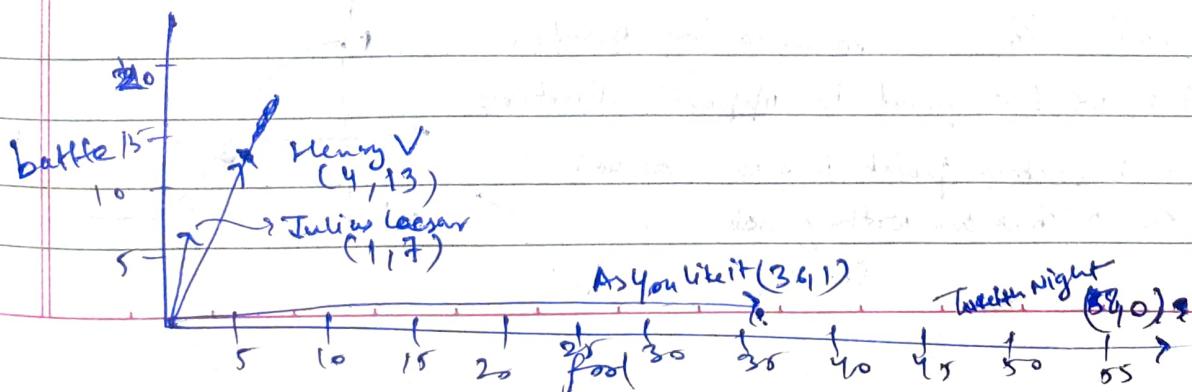
or  
word context matrix



good is the context  
battle is word

battle is ~~content~~ in  
context with good  
2 times

rows are known as vectors or word matrix



## Word-word matrix (or term-context matrix)

Two words are similar in meaning if their context vectors are similar.

content of apricot

	computer	data	pinch	result	sugar	more words
apricot	0	0	0	1	0	
pineapple	0	0	0	1	0	
digital	0	2	1	0	1	0
information	0	1	6	0	4	0

in context of information, data is present 6 times

Apricot and pineapple have similar values in their vectors  
then we say they are similar

→ Apricot is similar to pineapple

## Cosine for computing word similarity

$$\text{cosine}(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{\|\vec{v}\| \|\vec{w}\|} = \frac{\sum_{i=1}^N v_i w_i}{\sqrt{\sum_{i=1}^N v_i^2} \sqrt{\sum_{i=1}^N w_i^2}}$$

$v_i$  = Count for word  $v$  in context  $i$

$w_i$  = Count for word  $w$  in context  $i$

$\text{Cos}(v, w)$  is the cosine similarity of  $\vec{v}$  and  $\vec{w}$

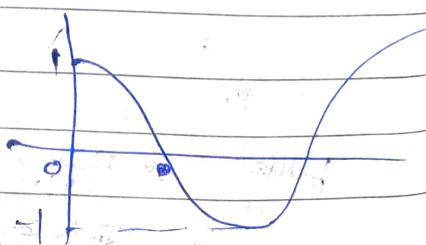
$$\begin{aligned} \vec{a} \cdot \vec{b} &= |\vec{a}| |\vec{b}| \cos \theta \\ \frac{\vec{a} \cdot \vec{b}}{|\vec{a}| |\vec{b}|} &= \cos \theta \end{aligned}$$

## Cosine as a similarity metric

-1: vectors point in opposite directions

0+1: vectors point in same direction

0: vectors are orthogonal



	large	data	computer
apricot	1	0	0
digital	0	1	2
information	1	6	1

1.1 0.6 0.1

$$\text{cosine}(\text{apricot}, \text{information}) = \frac{1 + 0 + 0}{\sqrt{1+0+0} \cdot \sqrt{1+36+1}} = \frac{1}{1 \cdot \sqrt{38}} = 0.16$$

$$\text{cosine}(\text{digital}, \text{information}) = \frac{0.1 \quad 1.6 \quad 2.1}{\sqrt{0+1+4} \cdot \sqrt{0+1+36+1}} = \frac{8}{\sqrt{38} \cdot \sqrt{38}} = 0.58$$

$$\text{cosine}(\text{apricot}, \text{digital}) = \frac{1.0 \quad 0.1 \quad 0.2}{\sqrt{1+0+0} \cdot \sqrt{0+0+4}} = 0$$

raw freq is bad representation

tf-idf : value for the word  $t$  in document  $d$ :  $w_{t,d} = t_f \times idf_t$   
 from freq-inverse document freq words like "the" or "good" have very low idf

(pointwise mutual information)  $PMI(w_1, w_2) = \log \frac{p(w_1, w_2)}{p(w_1) \cdot p(w_2)} = \frac{\text{probability of } w_1 \text{ and } w_2 \text{ occurring together}}{\text{probability of } w_1 \text{ and } w_2 \text{ occurring independently}}$   
 See if words like "good" appear more often with "great" than we would expect by chance

PMI bw a target word  $w$  and context word  $c$  may be defined as

$$PMI(w, c) = \log \frac{p(w, c)}{p(w) \cdot p(c)}$$

Co-occuring things less often leads to low PMI values. However, we are concerned with co-occuring things very often.

$$\Rightarrow \text{use Pointwise PMI} , PPMI(w, c) = \max\left(\log \frac{p(w, c)}{p(w) \cdot p(c)}, 0\right)$$

Given a co-occurrence matrix with  $W$  rows of words and  $C$  columns of contexts and  $f_{ij}$  is the no. of times word  $w_i$  occurs in context  $c_j$ . Then PPMI matrix :-

$$PPMI_{ij} = \max\left(\log \frac{p_{ij}}{p_{iw} \cdot p_{jc}}, 0\right)$$

$$P_{ij} = \frac{f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{iw} = \frac{\sum_{j=1}^C f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}, \quad p_{jc} = \frac{\sum_{i=1}^W f_{ij}}{\sum_{i=1}^W \sum_{j=1}^C f_{ij}}$$

	Computer	data	result	pie	sugar	Count(w)
Cherry	2	8	9	442	25	486
Strawberry	0	0	1	60	19	80
Digital Information	1670	1683	85	5	4	3447
	3325	3982	378	5	13	7703
Count(context)	4997	5673	473	512	61	11716

$$P(w = \text{information}, c = \text{data}) = \frac{3982}{11716} = 0.3339$$

$$P(w = \text{information}) = \frac{7703}{11716}, P(c = \text{data}) = \frac{5673}{11716} = 0.4842$$

$$= 0.6575$$

$$\phi_{PMI}(\text{information}, \text{data}) = \max\left(\log_2 \frac{0.3339}{0.6575 \times 0.4842}, 0\right) = 0.0944$$

$$\geq \max(0.0944, 0)$$

$$= 0.0944$$

## Vectors

— Vectors based on term-freq (counts of words in a corpus wrt to each document)

~~— Vectors based on Positive PMI.~~

- Characteristics :-
- Sparse vectors (large no. of zero entries)
  - Large size vectors (length of vector is equal to vocabulary size  $|V|$ )
  - However, we can reduce the size using dimensionality reduction techniques such as PCA, SVD etc.

Word2Vec — uses a technique called skip gram

— negative sampling

Aim :- Instead of counting the values or calculating some values, train a

Task - Prediction - If a ~~target word~~ 't' is likely to show up near the context 'c'. <sup>classifier</sup>

However, we skip the prediction and the weights of the learned classifier will be used as embeddings.

For a classifier, we require a supervised or labelled training data.

The running text from the corpus can be considered as supervised training data.

Skip grams follows the steps as

- 1) Treat the target word  $t$  and its neighbouring context  $c$  as the sample
- 2) Randomly use the other words in the vocab to obtain the negative
- 3) Train the classifier such as logistic regression
- 4) Use the regression weights as embeddings

### The classifier.

Train a classifier, such that the tuple  $(t, c)$  generate the probabilities as

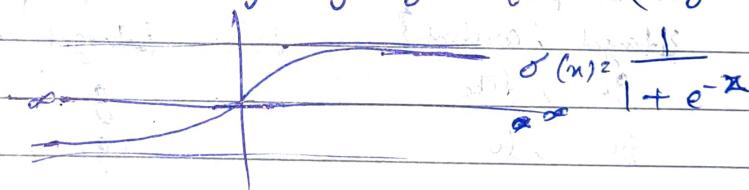
$$P(+ | t, c)$$

$$P(- | t, c) = 1 - P(+ | t, c)$$

A target word is likely to appear near the context, if its embedding is similar to the context embedding.

$$\text{similarity}(t, c) = \frac{t \cdot c}{\|t\| \|c\|} \quad \begin{matrix} \text{dot product} \\ \text{is used in place of} \end{matrix}$$

We bound this similarity using logistic function (sigmoid function)



$$P(+ | t, c) = \frac{1}{1 + e^{-t \cdot c}} \quad , \quad P(- | t, c) = 1 - P(+ | t, c) = 1 - \frac{1}{1 + e^{-t \cdot c}}$$

positive examples +	$t$	$c$
apricot	tablespoon	
apricot	of	
apricot	jam	
apricot	a	

negative examples -	$t$	$c$
apricot	apple	juice
apricot	my	forever
apricot	where	clear
apricot	cooking	if

$$P(+ | t, c_{1:k}) = \prod_{i=1}^k \frac{1}{1 + e^{-t \cdot c_i}}$$

$$\log P(+ | t_i, c_{1:k}) = b_0 + \sum_{j=1}^k \log \frac{1}{1 + e^{-t_i \cdot a_j}}$$

Since all the context words of target word are independent of each other we can multiply the probabilities together

$$L(\theta) = \sum_{(t,c) \in +} \log P(+|t,c) + \sum_{(t,c) \in -} \log P(-|t,c)$$

Maximize      Minimize

Optimizer like gradient descent may be used to optimize the weights or specifically, the randomly generated embeddings.

## Lecture 21

I would like to go someday nearby for lunch.

## True Sample

{target, context}

go , like

go, to

go, someplace

go to nearby

## The Sample

(target, context)

go, I

go - , lunch

go , twelve

go, therefore

$$S_m(t, c) = t \cdot c =$$