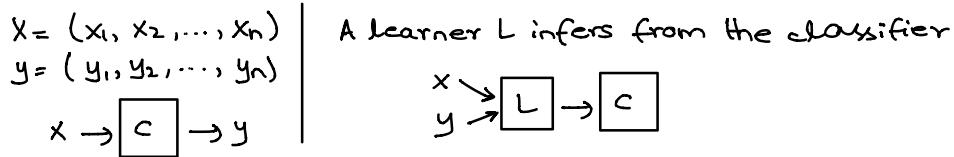


## Lecture # 8: Text Classification

The goal  $\Rightarrow$  take single observation, extract some useful features and thereby classify the observation into one set of discrete classes.



$\Rightarrow$  N-gram models can be used for classification, but for many tasks sequential relationships b/w words are largely irrelevant.

Naive Bayes  $\rightarrow$  family of probabilistic algorithms

For a document  $d$ , out of all classes  $c \in C$  the classifier returns the class  $\hat{c}$  which has max. posterior prob. given the document.  

$$\hat{c} = \operatorname{argmax}_{c \in C} P(c|d)$$

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

so,

$$\hat{c} = \operatorname{argmax}_{c \in C} \frac{P(d|c)P(c)}{P(d)}$$

since the  $P(d)$  or prob. of document is always same, we can drop it.

$$\hat{c} = \operatorname{argmax}_{c \in C} P(d|c)P(c)$$

### Modelling $P(d|c)$

Define a set of features that help classify the doc.

$$d = (f_1, f_2, \dots, f_n)$$

$$P(d|c) = P(f_1, f_2, \dots, f_n|c)$$

Since data is sparse, make the naive assumption: features are conditionally independent given the class  $c$ .

$$P(f_1, f_2, \dots, f_n|c) = P(f_1|c)P(f_2|c), \dots, P(f_n|c)$$

$\Rightarrow$  Prob. of a word occurring depends only on a class.

So, we can say:

$$\hat{c} = \operatorname{argmax}_{c \in C} P(c) \prod_{i=1}^n P(f_i|c)$$

$P(c) \rightarrow$  prob. of class before observing any data

$P(f_i|c) \rightarrow$  prob. of  $f_i$  in class  $c$

### Estimating Class Priors

$$\hat{P}(c) = \frac{N_c}{N}$$

$\Rightarrow P(c)$  is estimated with MLE

$N_c \rightarrow$  total number of training documents in class c

$N \rightarrow$  total training documents

In other words,  $\hat{P}(c)$  is simply proportion of training documents belonging to class c.

Consider the example:

	the	your	model	cash	Money	class	account	orderz	spam?
doc 1	12	3	1	0	0	2	0	0	-
doc 2	10	4	0	4	0	0	2	0	+
doc 3	25	4	0	0	0	1	1	0	-
doc 4	14	2	0	1	3	0	1	1	+
doc 5	17	5	0	2	0	0	1	1	+

$$\hat{P}(\text{spam}+) = \frac{3}{5}$$

$$\hat{P}(\text{spam}-) = \frac{2}{5}$$

$\Rightarrow P(f_i|c)$  can be estimated with simple smoothing.

$$\hat{P}(f_i|c) = \frac{\text{count}(f_i, c) + \alpha}{\sum_{f \in F} (\text{count}(f, c) + \alpha)}$$

$\text{count}(f_i, c) \rightarrow$  # of times  $f_i$  occurs in class c.

$F \rightarrow$  set of possible features

$\alpha \rightarrow$  smoothing parameter, optimized on held out data.

$$\hat{P}(\text{your}|+) = \frac{(4+2+5+\alpha)}{(\text{all words in } + \text{ class}) + \alpha |F|} = \frac{11+\alpha}{(68 + \alpha |F|)}$$

$$\hat{P}(\text{your}|-) = \frac{7+\alpha}{49 + \alpha |F|} \quad \hat{P}(\text{orderz}|+) = \frac{2+\alpha}{(68 + \alpha |F|)}$$

Let a test document be: get your cash and your orderz  $\rightarrow d$

$$\begin{aligned}
 p(+|d) &\propto p(+) \prod_{i=1}^n p(f_i|+) \\
 &\text{Let } \hat{p}(f_i|+) = \frac{\alpha}{68 + \alpha |F|} \text{ for all words (features) not shown} \\
 &= p(+). \frac{\alpha}{68 + \alpha |F|} \cdot \frac{11+\alpha}{(68 + \alpha |F|)} \cdot \frac{7+\alpha}{(68 + \alpha |F|)} \cdot \frac{2+\alpha}{(68 + \alpha |F|)} \cdot \frac{11+\alpha}{(68 + \alpha |F|)} \cdot \frac{2+\alpha}{(68 + \alpha |F|)}
 \end{aligned}$$

- Repeat the same for  $P(-|d)$
- Choose the class (+ or -) with high score

Since multiplying large numbers of small prob is problematic.

Actual implementation of Naive Bayes uses costs.

- Costs are -ve log prob.
- we can sum them and avoid underflow
- look for lowest cost overall

So, the equation becomes:

$$\hat{c} = \underset{c \in C}{\operatorname{argmin}} \left( -\log(P(c)) + \sum_{i=1}^n -\log(P(f_i|c)) \right)$$

⇒ Choosing features can be tricky.

### Advantages of Naive Bayes

- Very easy to implement, fast to train & test.
- Doesn't require much data.

### Problems

- Naive Bayes is naive.
- Features are not usually independent

→ Adding multiple feature types leads to stronger correlations.

### Maximum Entropy Model (MaxEnt)

→ Commonly known as Multinomial Logistic Regression.



more than 2 classes.

Also called log-linear model, one-layer neural network, single neuron classifier.

Can be defined as:  $\hat{c} = \underset{c \in C}{\operatorname{argmax}} P(c|d)$

Unlike NB, we do not apply Bayes Rule, we simply model  $P(c|d)$  directly.

Example: Given a webpage document, what topic does it belong to?

$\vec{x}$  → all the words in document  
 $c$  → latent class

In MaxEnt,  $\vec{x}$  is the observed data.

Features are functions on both observations  $\vec{x}$  and class  $c$ .

$c$	class
1	Travel
2	Sport
3	Finance

MaxEnt is used when dependent variable in question is nominal, and for which there are more than two categories.

If there are three classes, features come in groups of three.  
For example,

$$\begin{aligned}f_1 &: \text{contains('ski')} \& c = 1 \\f_2 &: \text{contains('ski')} \& c = 2 \\f_3 &: \text{contains('ski')} \& c = 3\end{aligned}$$

- training docs from  $c_1$  with 'ski' in it will have  $f_1$  active.

Each feature  $f_i$  has a real-valued weight  $w_i$  (learned in training).

### Classification with MaxEnt

Choose the class that has highest probability.

$$P(c|x) = \frac{1}{Z} \exp\left(\sum_i w_i f_i(\vec{x}, c)\right)$$

where  $Z = \sum_c \exp\left(\sum_i w_i f_i(\vec{x}, c)\right)$   $\rightarrow$  normalization constant

$w \cdot f \rightarrow$  is just a dot product

$\Rightarrow P(c|x)$  is just a monotonic fit of this dot product

so, we will end up choosing the class for which  $\vec{w} \cdot \vec{f}$  is highest.

Consider the following features and weights.

$f_1$ :	$\text{contains('ski')}$ & $c = 1$	$w_1 = 1.2$
$f_2$ :	$\text{contains('ski')}$ & $c = 2$	$w_2 = 2.3$
$f_3$ :	$\text{contains('ski')}$ & $c = 3$	$w_3 = -0.5$
$f_4$ :	$\text{link\_to('expedia.com')}$ & $c = 1$	$w_4 = 4.6$
$f_5$ :	$\text{link\_to('expedia.com')}$ & $c = 2$	$w_5 = -0.2$
$f_6$ :	$\text{link\_to('expedia.com')}$ & $c = 3$	$w_6 = 0.5$
$f_7$ :	$\text{num\_links}$ & $c = 1$	$w_7 = 0.0$
$f_8$ :	$\text{num\_links}$ & $c = 2$	$w_8 = 0.2$
$f_9$ :	$\text{num\_links}$ & $c = 3$	$w_9 = -0.1$

Suppose test document contains 'ski' and 6 outgoing links.

$$\text{Travel: } \sum_i w_i f_i(\vec{x}, c) = (1.2) + (0.0)(6) = 1.2$$

$$\text{Sports: } \sum_i w_i f_i(\vec{x}, c=2) = (2.3) + (0.2)(6) = 3.5$$

$$\text{Finance: } \sum_i w_i f_i(\vec{x}, c=3) = -0.5 + (-0.1)(6) = -1.1$$

We can conclude that Sports will be most probable.

### Training the model

Given annotated data, choose weights that make the labels most probable under the model.

$$\text{Given } x^1, \dots, x^n \text{ choose } c^1, \dots, c^n \Rightarrow \hat{w} = \underset{w}{\operatorname{argmax}} \sum_j \log P(c^j | x^j) \xrightarrow{\text{CMLF}}$$

## Lecture # 13: Context Free Parsing

CFGs

- Vocabulary of terminal symbols  $\Sigma$
- set of non-terminal symbols  $N$
- start symbol  $s \in N$
- Production rules  $X \rightarrow \alpha$   
 $X \in N, \alpha \in (N \cup \Sigma)^*$

Input: sentence  $w = (w_1, \dots, w_n)$  and CFG  $G$

Output (recog): true iff  $w \in \text{Language}(G)$

Output (pausing): one or more derivations for  $w$ , under  $G$

Agenda = {state 0}

while (Agenda not empty)

```

    s = pop a state from agenda
    if s is success-state return s
    elif s is not failure state:
        generate new states from s
        push new states onto agenda
    
```

return nil

### CFGs in CNF (Chomsky Normal Form)

- Vocabulary of terminal symbols  $\Sigma$
- Set of non-terminal symbols  $N$
- Special start symbol  $s \in N$
- Production rule of  $X \rightarrow \alpha$   
 $X \in N$   
 $\alpha \in N^* \cup \Sigma^*$

CFG

$\mathcal{L}_1$ Grammar		$\mathcal{L}_1$ in CNF
$S \rightarrow NP VP$		$S \rightarrow NP VP$
$S \rightarrow Aux NP VP$		$S \rightarrow X1 VP$ X1 → Aux NP
$S \rightarrow VP$		$S \rightarrow book   include   prefer$ $S \rightarrow Verb NP$ $S \rightarrow X2 PP$ $S \rightarrow Verb PP$ $S \rightarrow VP PP$ $NP \rightarrow I   she   me$ $NP \rightarrow TWA   Houston$ $NP \rightarrow Det Nominal$ $Nominal \rightarrow book   flight   meal   money$ $Nominal \rightarrow Nominal Noun$ $Nominal \rightarrow Nominal PP$ $VP \rightarrow book   include   prefer$ $VP \rightarrow Verb NP$ $VP \rightarrow Verb NP PP$ $VP \rightarrow Verb PP$ $VP \rightarrow VP PP$ $PP \rightarrow Preposition NP$
$NP \rightarrow Pronoun$		
$NP \rightarrow Proper-Noun$		
$NP \rightarrow Det Nominal$		
$Nominal \rightarrow Noun$		
$Nominal \rightarrow Nominal Noun$		
$Nominal \rightarrow Nominal PP$		
$VP \rightarrow Verb$		
$VP \rightarrow Verb NP$		
$VP \rightarrow Verb NP PP$		
$VP \rightarrow Verb PP$		
$VP \rightarrow VP PP$		
$PP \rightarrow Preposition NP$		

### Convert CFGs to CNF

For each rule,  
 $X \rightarrow ABC$   
Rewrite  
 $X \rightarrow A X_2$   
 $X_2 \rightarrow BC$

- Introduce a new non-terminal  
 $X_2$

CNF

$S \rightarrow Aux NP VP$   
 $\downarrow$

$S \rightarrow X1 VP$   
 $X1 \rightarrow Aux NP$

$VP \rightarrow Verb NP PP$

$VP \rightarrow X2 PP$

$X2 \rightarrow Verb NP$

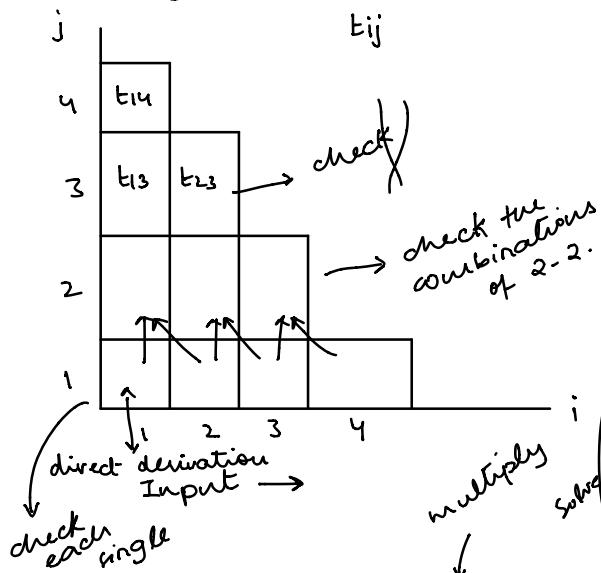
## CYK - Algorithm

- CFG parsing algorithm
- work only with CFG in CNF.

```

for i = 1 ... n
    C[i-1,i] = {V | V → w_i}
for l = 2 ... n
    for i = 0 ... n-l
        k = i + l
        for j = i+1 ... k-1
            C[i,j,k] = C[i,j] ∪ {V | V → VZ, V ∈ C[i,j], Z ∈ C[j,k]}
return true if S ∈ C[0,n]
    
```

## CYK Algorithm Dry Run



0,1 S, Nominal, VP, Verb, Noun	0,2 = [0,1][1,2] S Det = 0 Nom Det = 0 VP Det = 0 Verb Det = 0 N Det = 0	0,3 = [0,1][1,3] ∪ [0,2][2,3] S NP, X2	0,4 = [0,1][1,4] S NP, X2 Since S is here, it is correct.
Book	1,2 Det	1,3 = [1,2][2,3] Det N = 0 Det Nom = NP	1,4 = [1,2][2,4] ∪ [1,3][3,4] NP
	that	2,3 Noun, Nominal	2,4 = [2,3][3,4] N N = 0 N Nom = 0 Nom N = Nominal Nom Nom = 0
		meal	3,4 Noun, Nominal

Book that meal flight → n = 4

example

→ Bottom-up approach

SNP = 0 VP NP = 0 Verb NP = S, VP, X2 Nom NP = 0 Noun NP = 0	[1,2][2,4] Det Nominal = NP NP Noun = 0 NP Nominal = 0
--	---

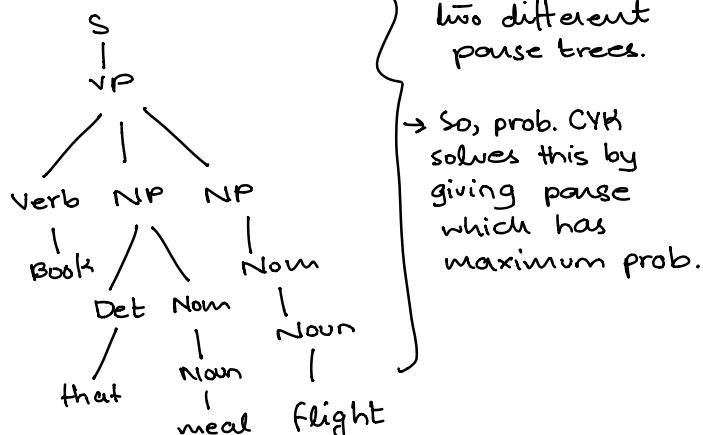
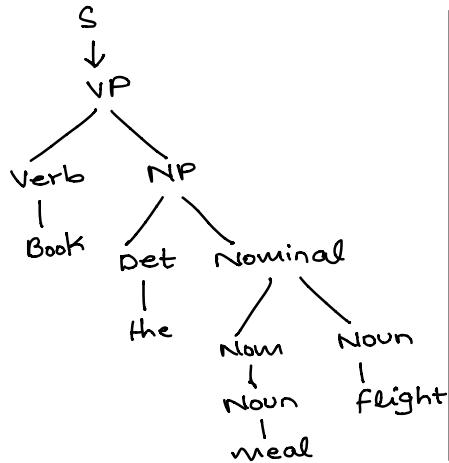
0,4 = [0,2][2,4] ∪ [0,3][3,4] ∪ [0,1][1,4]	S Noun = 0 S Nom = 0 VP Noun = 0 VP Nom = 0
--	--

Solve

</div

### Probabilistic CYK

→ used if ambiguous pause



two different pause trees.

So, prob. CYK solves this by giving pause which has maximum prob.

Example.

$$\begin{aligned}
 S \rightarrow NP\ VP & .80 \\
 NP \rightarrow Det\ N & .30 \\
 VP \rightarrow V\ NP & .20 \\
 V \rightarrow \text{includes} & .05 \\
 Det \rightarrow \text{the} & .50 \\
 Det \rightarrow a & .40 \\
 N \rightarrow \text{meal} & .01 \\
 N \rightarrow \text{flight} & .02
 \end{aligned}$$

0,1	0,2	0,3	0,4	0,5
det .50	$[0,1][1,2]$ $NP = 3 \times 5 \times 0.02 = 0.03$	0	0	$S$ 0.000 000 288
the	1,2 N .02	1,3 [1,2][2,3] 0	1,4 0	1,5 0
flight	2,3 V .05	2,4 [2,3][3,4] 0	2,5 VP .000012	

Sentence:

the flight includes a meal

$$\begin{aligned}
 [0,3] &= \underbrace{[0,1][1,3]}_0 \cup \underbrace{[0,2][2,3]}_0 \\
 &= 0 \\
 [1,4] &= \underbrace{[1,2][2,4]}_0 \cup \underbrace{[1,3][3,4]}_0 \\
 &= 0 \\
 [2,5] &= \underbrace{[2,3][3,5]}_0 \cup \underbrace{[2,4][4,5]}_0 \\
 VP &= .2 \times .05 \\
 &\quad \times .0012 \\
 VP &= .000012
 \end{aligned}$$

$[1,5] =$ $[1,2][2,5] \cup [1,3][3,5] \cup$ $[1,4][4,5]$ $= 0 \cup 0 \cup 0$ $= 0$	$3,4$ det .40	$3,5$ $[3,4][4,5]$ $NP = 3 \times .4 = .0012$
	$4,5$ N .01	
$[0,5] =$ $\underbrace{[0,1][1,5]}_S \cup \underbrace{[0,2][2,5]}_S \cup$ $\underbrace{[0,3][3,5]}_0 \cup \underbrace{[0,4][4,5]}_0$ $= .8 \times .03 \times .000012 = 0.000000288$	meal	

$$\begin{aligned}
 [0,4] &= \underbrace{[0,1][1,4]}_0 \cup \underbrace{[0,2][2,4]}_0 \cup \underbrace{[0,3][3,5]}_0 \\
 &= 0
 \end{aligned}$$

## Earley Parsing

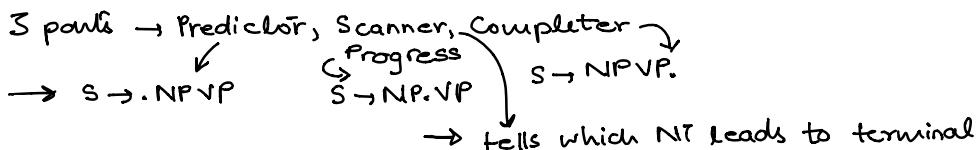
- Allows arbitrary CFGs
- Top down control
- Fills a table in a single sweep
- Table is length  $N+1$ ;  $N \rightarrow \#$  of words
- Table entries represent
  - ✓ Completed constituents and their locations
  - ✓ In-progress constituents
  - ✓ Predicted constituents

Table entries are called states and are represented with dotted rules.

$S \rightarrow .VP$	}	VP is predicted
$NP \rightarrow Det. Nominal$		NP is in-progress
$VP \rightarrow V NP$		VP is completed
$S \rightarrow .VP[0,0]$	}	VP is predicted at the start of sentence
$NP \rightarrow Det. Nominal[1,2]$		NP is in progress; det goes from 1 to 2.
$VP \rightarrow V NP, [0,3]$		A VP has been found at 0 & ending at 3.

## Steps

- 1) Predict all the states you can upfront
- 2) Read a word
  - Extend states based on matches
  - Generate new predictions
  - Go to step 2
- 3) When out of words, look at charts & see if a winner.



- $\rightarrow$  Last cell shows the parsing to be complete.  
 $\rightarrow$  each cell should have  $NT \rightarrow T$  transition.

Ex. Book that meal flight

$\rightarrow$  Since  $n=4$ , so table would be  $5 \times 5$

$\rightarrow$  Construct the table  $5 \times 5$ , the last  $0,4$  should show completion.

4,0	4,1	4,2	4,3	4,4
			3,3 (pred) PP → Prep NP	3,4
		2,2 (Pred) Nominal → Noun Nom → Nom Noun Nom → Nom PP	2,3 (Scan) Noun → meal Nom → Noun. Nom → Nom. Nom → Nom PP	2,4 (Scan) Noun → flight Nom → Nom. Nom → Nom.
	1,1 Pred NP → Pronoun X NP → Det Noun → NP → Proper Noun X PP → Preposition X	4,2 Scanner Det → that (Completer) NP → Det. Nominal	1,3 (Completer) NP → Det Noun.	1,4
0,0 (pred) Temp → S S → NP VP S → Aux NP VP S → VP NP → Pronoun NP → Proper Noun	0,1 (Scanning) Verb → book (Completer) VP → Verb. VP → Verb. NP VP → Verb. NP PP VP → VP. PP	0,2	0,3 (Comp) VP → Verb NP. VP → Verb NP. VP → Verb. PP VP → VP. PP S → VP.	0,4 VP → Verb NP. S → VP.

✗ disregard because book not aux. So, do not derive.

So, disregard all NP → cases.

write this.

$\{$   
 NP → Det Noun  
 VP → Verb  
 VP → Verb NP  
 VP → Verb NP PP  
 VP → Verb PP  
 VP → VP PP

Outcome:

Verb → Book [0,1]
Det → that [1,2]
Noun → meal [2,3]
Noun → flight [2,4]
Noun → Nominal Noun [2,4]
NP → Det Nominal [1,3]
VP → Verb NP [0,4]
S → VP [0,4]

→ Should have complete sentence

### Note

- When we are predicting, we go up in the table.
- When we have scanner or completer, we move right.

Runtime →  $O(n^3)$   
Memory →  $O(n^2)$

## Lecture #15: Lexical Semantics

Lexical semantics → linguistic study of meaning

Decomposition → what the components of meanings in a word are?

Ontological → How the meaning of word relates to meanings of other words.

Distributional → what contexts the word is found in, relative to other words.

### Semantic Relations

Synonymy → Equivalence

Antonymy → Opposition

Hyponymy → subset, is-a relation

Hypernymy → superset

Meronymy → part-of relation

Holonymy → has-a relation

✓ WordNet is a useful resource → requires many years

• There are intrinsic limit to this of manual effort

↳ The ontology is only as good as the ontologist.

### Word Similarity: Distributional Models

The intuition → meaning of a word is related to distribution of words around it.

• Represent word  $w$  as a feature vector.

• Suppose we had one binary feature  $f_i$  representing each of  $N$  words in lexicon  $v_i$ . The feature means  $w$  occurs in neighborhood of word  $v_i$  and hence takes value 1 if  $w$  and  $v_i$  are in same context window and 0 otherwise.

→ We could represent meaning of word  $w$  as the feature vector

$$\vec{w} = (f_1, f_2, f_3, \dots, f_N)$$

Example:

The corpus is

$s_1 = \text{The sky is blue \& beautiful}$

$s_2 = \text{Love this blue \& beautiful sky}$

$s_3 = \text{Today is sunny}$

$s_4 = \text{The sky is very blue \& sky is very beautiful today.}$

If  $w = \text{sky}$ ,  $v_1 = \text{blue}$ ,  $v_2 = \text{beautiful}$ ,  $v_3 = \text{love}$  &  $v_4 = \text{today}$

co-occurrence vector for  $w$  from corpus above would be:

$$\vec{w} = (1, 1, 0, \dots)$$

Co-occurrence Vector → Remove stopwords

$s_1 = \text{sky, blue, beautiful}$

$s_2 = \text{love, blue, beautiful, sky}$

$s_3 = \text{today, sunny}$

$s_4 = \text{sky, very, blue, beautiful, today}$

	sky	blue	beautiful	love	today	sunny	very
sky	1	1	1	1	1	0	1
blue	1	0	1	1	1	0	1
beautiful	1	1	0	1	1	0	1
love	1	1	1	0	0	0	0
today	1	1	1	0	0	1	1
sunny	0	0	0	0	1	0	0
very	1	1	1	0	1	0	0

#### 4 kinds of vector models

- Sparse vector representations
  - Mutual information weighted word co-occurrence matrix
- Dense vector representations
  - SVD (and Latent Semantic Analysis)
  - NN-inspired models (Skipgrams, cBOW)
  - Brown clusters

#### Distributional Similarity Measure

- Specifying a distributional similarity measure requires that we specify three parameters:
    - how co-occurrence terms are defined (neighbor?)
    - how these terms are weighted
    - what vector distance metric we use (cosine?)
  - ✓ If vocab size = N, each w (word) had N features specifying whether vocabulary element vj occurred in the neighbourhood.
  - ✓ Neighborhoods range from small window of words to very large windows of  $\pm 500$  words.
- Even with removal of stopwords, when used on large corpora, these co-occurrence vectors tend to be large. So choose words that occur in some sort of grammatical relation or syntactic dependency to target word.
- tea, water, coffee → drink etc

→ Since each word can be in a variety of different dependency relations with other words, augment the feature space. Each feature is now a pair of word and a relation. So we have vector of  $N \times R$  features where  $R = \#$  of possible relations.

Co-occurrence values are typically thought of as weights or measures of association b/w each target word  $w$  and given feature  $f$ .

Association measure can be:

- Binary → for each feature, 1 if relevant word occurred in context otherwise 0.
- Relative → with which the particular context feature had co-occurred with target word.

⇒ Freq, or prob are better measures of association.

### Measures of Association with context

For target word  $w$ , each element of its  $w$ -occurrence vector is a feature  $f$  consisting of a relation  $r$  and a related word  $w'$ ; we can say  $f = (r, w')$

The prob. of feature  $f$  given a target word  $w$  is

$P(f|w)$  for which MLE is

$$P(f|w) = \frac{\text{count}(f, w)}{\text{count}(w)}$$

Similarly for joint prob.  $P(f, w)$

$$P(f, w) = \frac{\text{count}(f, w)}{\sum_{w'} \text{count}(f, w')}$$

⇒ words like it, they, are ⇒ not good discriminative, no info.

Mutual information b/w two random variables  $X$  &  $Y$ :

$$I(X, Y) = \sum_x \sum_y P(x, y) \log_2 \left( \frac{P(x, y)}{P(x)P(y)} \right)$$

Pointwise MI ⇒ measure of how often two events  $x$  and  $y$  occur compared with what we would expect if they were independent.

$$I(x, y) = \log_2 \frac{P(x, y)}{P(x)P(y)}$$

We can use this to compute co-occurrence matrix target  $w$  (word), feature  $f$

$$\text{assoc}_{\text{PMI}}(w, f) = \log_2 \frac{P(w, f)}{P(w)P(f)}$$

Computing PPMI on a term-context matrix

	computer	data	pinch	result	sugar	count (w, context)	word	context
apricot	0	0	1	0	1			
pineapple	0	0	1	0	1			
digital	2	1	0	1	0			
information	1	6	0	4	0			
								= 11

$$p_{ij} = \frac{f_{ij}}{\sum_{i=1}^n \sum_{j=1}^c f_{ij}}$$

$$p(w_i) = \frac{\sum_{j=1}^c f_{ij}}{N}$$

$$p(c_j) = \frac{\sum_{i=1}^n f_{ij}}{N}$$

$$pmi_{ij} = \log_2 \frac{p_{ij}}{p_i * p_j}$$

$$ppmi_{ij} = \begin{cases} pmi_{ij} & \text{if } > 0 \\ 0 & \text{otherwise} \end{cases}$$

PPMI → positive pointwise mutual information

↳ replace all -ve values with zero.

$$N = \sum_{i=1}^n \sum_{j=1}^c f_{ij} = 19$$

$$p(w=\text{information}, c=\text{data}) = 6/19 = 0.32$$

$$p(w=\text{information}) = 11/19 = 0.58$$

$$p(w=\text{data}) = 7/19 = 0.37$$

	p(w, context)					p(w)
	computer	data	pinch	result	sugar	
apricot	0.00	0.00	0.05	0.00	0.05	0.11
pineapple	0.00	0.00	0.05	0.00	0.05	0.11
digital	0.11	0.05	0.00	0.05	0.00	0.21
information	0.05	0.32	0.00	0.21	0.00	0.58
p(context)	0.16	0.37	0.11	0.26	0.11	

Now we can compute pmii.

$$pmi(\text{info}, \text{data}) = \log_2 \left( \frac{0.32}{0.58 \times 0.37} \right) = 0.58$$

→ Similarly we can compute for all.

For PPMI table, 0 for all -ve values.

PPMI(w,context)					
	computer	data	pinch	result	sugar
apricot	-	-	2.25	-	2.25
pineapple	-	-	2.25	-	2.25
digital	1.66	0.00	-	0.00	-
information	0.00	0.57	-	0.47	-

→ To define similarity between two target words  $v$  &  $w$ , we need a measure of taking two such vectors and giving a measure of vector similarity.

Two simplest are:

$$\text{Manhattan (L1)} = \sum_{i=1}^n |v_i - w_i|$$

$$\text{Euclidean (L2)} = \sqrt{\sum_{i=1}^n (v_i - w_i)^2}$$

⇒ But both are sensitive to extreme values.

$$\text{dot-product}(v, w) = \vec{v} \cdot \vec{w} = \sum_{i=1}^n v_i w_i = v_1 w_1 + v_2 w_2 + \dots + v_N w_N$$

↳ High when two vectors in similar directions.

longer if vector is longer ⇒ higher values in each

$$|\vec{v}| = \sqrt{\sum_{i=1}^n v_i^2}$$

↓ dimension  
more frequent words will have higher dot products.  
so, divide by lengths.

$$\cos(v, w) = \frac{v \cdot w}{\|v\| \|w\|} = \frac{\sum_{i=1}^n v_i w_i}{\sqrt{\sum_{i=1}^n v_i^2} \sqrt{\sum_{i=1}^n w_i^2}}$$

0 → orthogonal

-1 → vectors in opposite direction

+1 → vectors in same direction

PPMI are positive, so range → 0-1

So,

$$\cosine(\text{apricot}, \text{information}) = \frac{1 \times 1 + 0 \times 6 + 1 \times 0}{\sqrt{1^2 + 0^2 + 0^2} \times \sqrt{1^2 + 6^2 + 1^2}} = 0.16$$

and so on.

## Lecture # 16: SVD

Alternative to PPMI for measuring association

tf-idf (term frequency - inverse document freq.)

freq. of a word

$df_i \rightarrow$  document frequency of word  $i$

$$idf_i = \log\left(\frac{N}{df_i}\right)$$

$w_{ij} = \text{word } i \text{ in document } j$

$$w_{ij} = tf_{ij} idf_j$$

→ PPMI vectors are sparse & long

dense vectors → short & easy to learn

↳ generalize better than storing counts

Three ways to get dense vectors:

- SVD → Singular Value Decomposition
- Neural Language Model → skip grams, CBOW
- Brown Clustering

### Dense Vectors via SVD

- Approx N-dimensional dataset using fewer dimensions.
- By first rotating the axes into new space  
→ Highest order dimension captures the most variance

Any rectangular  $w \times c$  matrix  $X$  equals the product of 3 matrices.

$W$ : rows corresponding to original but  $m$  columns represent dimensions in a new latent space such that

- $M$  column vectors are orthogonal to each other
- Columns are ordered by amount of variance in dataset each new dimension accounts for.

$S$ : diagonal  $m \times m$  matrix of singular values expressing importance of each dimension

$C$ : columns corresponding to original but  $m$  rows corresponding to singular values

contexts

$$\begin{matrix} \text{words } X = \\ \begin{matrix} w & S & C \\ \hline w \times c & w \times m & m \times c \end{matrix} \end{matrix}$$

→ Instead of keeping all  $m$  dimensions, keep top- $k$  singular values.

Each row of  $w$ :

→ A  $k$ -dimensional vector for word  $w$

contexts

$$\rightarrow \text{words } X = \begin{matrix} w \\ n \times c \end{matrix} \quad \begin{matrix} s \\ w \times k \end{matrix} \quad \begin{matrix} c \\ k \times k \end{matrix} \quad \begin{matrix} LSA \\ k \times c \end{matrix}$$

LSA →  $k=300$  commonly

The cells are weighted by product of two weights.

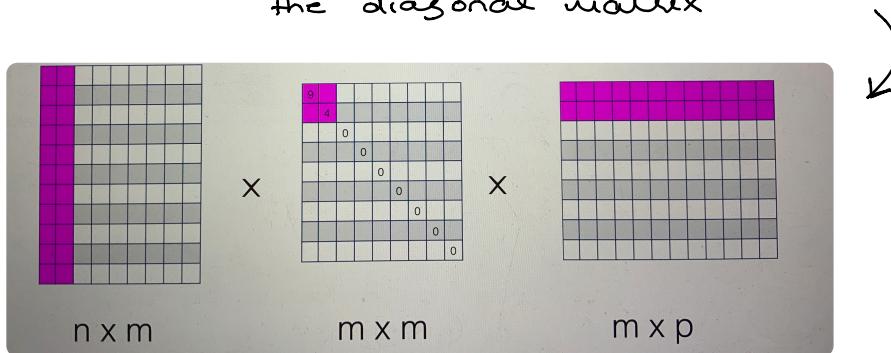
local weight = log freq

Global weight = Either idf or entropy measure

Any  $n \times p$  matrix  $X$  can be decomposed into product of three matrices (where  $m = \#$  of linearly independent rows)

$$X = \begin{matrix} n \times m \\ n \times p \end{matrix}, \begin{matrix} m \times m \\ \downarrow \end{matrix}, \begin{matrix} m \times p \end{matrix}$$

we can approx the full matrix  $X$  by considering the leftmost  $k$  terms in the diagonal matrix



LSA → applying SVD to term-document co-occurrence matrix

- Terms are typically weighed by tf-idf

- Dimensionality reduction → for terms, from a  $D$ -dimensionality sparse vector to a  $K$ -dimensionality dense one.  
 $K \ll D$

### Applying SVD to PPMI

$$X = \begin{matrix} w \\ |V| \times |V| \end{matrix} \quad \begin{matrix} s \\ |V| \times k \end{matrix} \quad \begin{matrix} c \\ k \times |V| \end{matrix}$$

→ Each row in  $w$  is  $k$ -dimensional rep. of each word  $w$ .  
→  $k$  ranges from 50 to 1000.

## Example

$A = mxn$ ,  $m = \text{terms}$ ,  $n = \text{documents/concepts}$

## d1: Romeo & Juliet

d2: Juliet: o happy dagger!

d3: Romeo died by dagger

d4: Live free or die, that's New Hampshire's motto

ds: Did you know, New Hampshire is in New England.

Query: dies, dagger

construct matrix A.

$A(\text{words}, \text{document})$

	d1	d2	d3	d4	d5
Romeo	1	0	1	0	0
Juliet	1	1	0	0	0
happy	0	1	0	0	0
dagger	0	0	1	0	0
live	0	0	0	1	0
die	0	0	1	1	0
free	0	0	0	1	0
New Hampshire	0	0	0	1	1

$$A = U S V^T$$

Eigen vector of B

$$B = AAT \quad (\text{term-term matrix})$$

$$\begin{array}{c} \downarrow \\ \left[ \begin{array}{cccccc} 1 & 0 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 \end{array} \right] \times \left[ \begin{array}{cccccc} 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 \end{array} \right] = \end{array}$$

	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>	C <sub>4</sub>	C <sub>5</sub>	C <sub>6</sub>	C <sub>7</sub>	C <sub>8</sub>
1	2	1	0	1	0	1	0	0
2	1	2	1	0	0	0	0	0
3	0	1	1	0	0	0	0	0
4	1	1	0	0	1	0	0	0
5	0	0	0	0	1	1	1	1
6	1	0	0	1	1	2	1	1
7	0	0	0	0	1	1	1	1
8	0	0	0	0	1	1	1	2

$A^T$        $(5 \times 8)$        $(8 \times 8)$

$B =$ 

	C <sub>1</sub>	C <sub>2</sub>	C <sub>3</sub>	C <sub>4</sub>	C <sub>5</sub>	C <sub>6</sub>	C <sub>7</sub>	C <sub>8</sub>
1	2	1	0	1	0	1	0	0
2	1	2	1	0	0	0	0	0
3	0	1	1	0	0	0	0	0
4	1	0	0	1	0	1	0	0
5	0	0	0	0	1	1	1	1
6	1	0	0	1	1	2	1	1
7	0	0	0	0	1	1	1	1
8	0	0	0	0	1	1	1	2

 $U =$ 

0.3902	-0.3497	0.0567	-0.1645	-0.4192	-0.2853	-0.5745	0.3285
-0.3902	0.3497	-0.0567	-0.1645	-0.1093	0.5812	-0.5745	0.1237
0.3902	-0.3497	0.0567	0.3971	0.5503	0.4518	-0.2380	0.0316
0.0018	0.3820	-0.5950	0.3971	0.2404	-0.4147	-0.2380	0.2365
-0.1995	-0.4766	-0.5867	-0.3971	0.1311	0.1665	0.2380	0.3601
-0.3920	-0.0323	0.5383	-0.0000	0.3715	-0.2482	0.0000	0.5966
0.5915	0.5089	0.0484	-0.3971	0.1311	0.1665	0.2380	0.3601
-0.0000	-0.0000	0.0000	0.5615	-0.5286	0.2959	0.3366	0.4522

$$C = A^T A \rightarrow (\text{doc-doc matrix})$$

$$C = \begin{bmatrix} 2 & 1 & 1 & 0 & 0 \\ 1 & 3 & 1 & 0 & 0 \\ 1 & 1 & 3 & 1 & 0 \\ 0 & 0 & 1 & 4 & 1 \\ 0 & 0 & 0 & 1 & 1 \end{bmatrix}$$

$$V^T =$$

-0.4298	0.3039	0.3039	-0.3039	0.7337
-0.5905	0.4926	0.2152	0.1173	-0.5905
-0.1957	-0.6832	0.6271	-0.2518	-0.1957
-0.6219	-0.4397	-0.4397	0.4397	0.1821
0.2040	0.0701	0.5241	0.7982	0.2040

↓  
eigen vectors of C

The S  $\rightarrow$  square root of eigen values of either  $A A^T(B)$  or  $A^T A(C)$ .

So, the S would be

0.7654	0	0	0	0
0	0.8952	0	0	0
0	0	1.5121	0	0
0	0	0	1.8478	0
0	0	0	0	2.2164

We can set k=2 (or any to get top k)

$$K=2 \quad A_{-2} = U_2 S_2 V_2^T \quad \left\{ \text{Truncated SVD} \right.$$

$$\text{So, } S_2 = \begin{bmatrix} 0.7654 & 0 \\ 0 & 0.8952 \end{bmatrix}$$

Similarly we can take 2 rows from U & 2 cols from V.

$\xrightarrow{\text{Term} = U \times S}$  doc-doc matrix's eigen vectors

$$\text{Romeo} = \begin{bmatrix} 0.7654 \times 0.3902 \\ 0.8952 \times -0.3497 \end{bmatrix} = \begin{bmatrix} 0.2986 \\ -0.3130 \end{bmatrix}$$

$$\text{Juliet} = \begin{bmatrix} 0.7654 \times -0.3902 \\ 0.8952 \times 0.3497 \end{bmatrix} = \begin{bmatrix} -0.2986 \\ 0.3130 \end{bmatrix}$$

$$\text{Doc} = SV^T$$

$$\text{Query} = \frac{\text{die} + \text{dagger}}{2} = \frac{\begin{bmatrix} -1.197 \\ -0.494 \end{bmatrix} + \begin{bmatrix} -1.001 \\ 0.742 \end{bmatrix}}{2} = \begin{bmatrix} -1.099 \\ 0.124 \end{bmatrix}$$

$\Rightarrow$  compute similarity using cosine similarity

$$\text{In } k=2, A_2 = U_2 S_2 V_2^T$$

0.3902	-0.3497	0.0567	-0.1645	-0.4192	-0.2853	-0.5745	0.3285
-0.3902	0.3497	-0.0567	-0.1645	-0.1093	0.5812	-0.5745	0.1237
0.3902	-0.3497	0.0567	0.3971	0.5503	0.4518	-0.2380	0.0316
0.0018	0.3820	-0.5950	0.3971	0.2404	-0.4147	-0.2380	0.2365
-0.1995	-0.4766	-0.5867	-0.3971	0.1311	0.1665	0.2380	0.3601
-0.3920	-0.0323	0.5383	-0.0000	0.3715	-0.2482	0.0000	0.5966
0.5915	0.5089	0.0484	-0.3971	0.1311	0.1665	0.2380	0.3601
-0.0000	-0.0000	0.0000	0.5615	-0.5286	0.2959	0.3366	0.4522

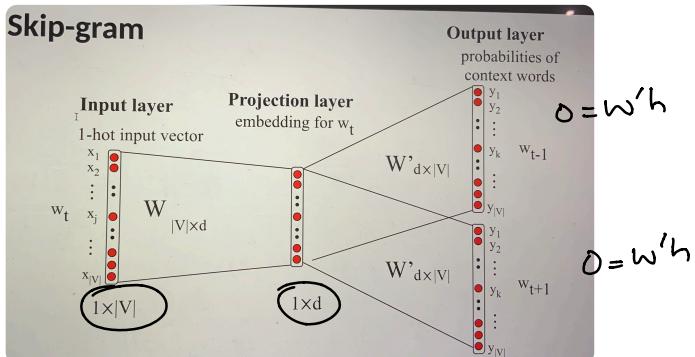
-0.4298	0.3039	0.3039	-0.3039	0.7337	$v_2$	$v_2^T$
-0.5905	0.4926	0.2152	0.1173	-0.5905		
-0.1957	-0.6832	0.6271	-0.2518	-0.1957		
-0.6219	-0.4397	-0.4397	0.4397	0.1821		
0.2040	0.0701	0.5241	0.7982	0.2040		

## Neural Language Models

→ Skip-gram & CBOW

↓  
faster, trainable, dense

### Skip-gram



$$O_k = v_k \cdot v_j$$

use softmax to turn into prob.

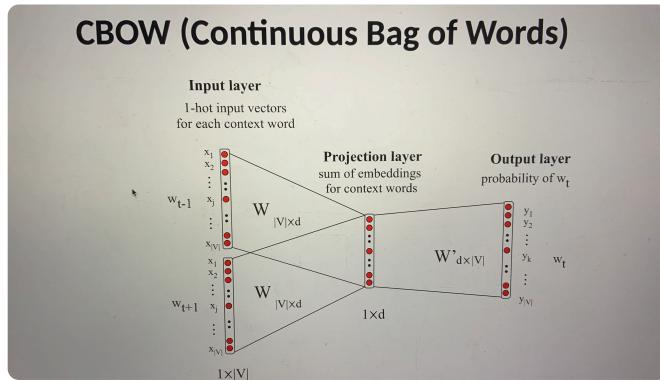
$$p(w_k | w_j) = \frac{\exp(v_k \cdot v_j)}{\sum_{i=1}^{|V|} \exp(v_i \cdot v_j)}$$

Since we have two embeddings  $v_j$  and  $v_j'$  for each word  $w_j$   
Perform linear combination.

$$\rightarrow \mathbf{W}\mathbf{W}'^T = \mathbf{M}^{PMI} - \log K$$

$\rightarrow$  So, skipgram is implicitly factoring a shifted version of the PMI matrix into the two embedding matrices.

### CBOW



### Brown Clustering

An agglomerative clustering algo that clusters based on which word precedes or follows them

### Evaluating Similarity

Intrinsic (task-based, end to end)

$\rightarrow$  QA, Spell check, Essay grading

Intrinsic

$\rightarrow$  Correlation b/w algo & human word sim.

### Lecture #9 & 10: Classification

#### Perceptron Learning

If  $n$ -th member of training set  $x(n)$  is correctly classified by weight vector  $w(n)$  computed at  $n$ -th iteration, no correction is made to weight vector.

$$w(n+1) = w(n) \text{ if } w^T x(n) > 0 \quad \# x(n) \in C_1 \\ w^T x(n) \leq 0 \quad \# x(n) \in C_2$$

Otherwise weight vector is updated

$$\rightarrow w(n+1) = w(n) - \gamma(n) x(n) \text{ if } w^T x(n) > 0 \quad \# x(n) \in C_1 \\ w(n+1) = w(n) + \gamma(n) x(n) \text{ if } w^T x(n) \leq 0 \quad \# x(n) \in C_2$$

Or

$$w(n+1) = w(n) + \gamma(n) e(n) x(n)$$

where

$$e(n) = d(n) - y(n)$$

$$\left\{ \begin{array}{l} \text{If } d=0 \rightarrow \text{subtract } (w_n - x(n)) \\ \text{If } d=1 \rightarrow \text{add } (w_n + x(n)) \end{array} \right\}$$

## Lecture # 11 : Syntax

Morphology deals with internal structure of the words.  
Syntax deals with combination of words.

↳ often irregular      ↳ often regular  
                                made up of general rules that apply across the board.

Semantics → about meaning

Syntax → about structure alone

A sentence can be syntactically well but semantically ill-formed.

Constituent: A group of words that go together

↳ larger than a word are called phrases.

→ Every Prep. P → contains a NP

### CFG

↙  
context-free because there is no context in LHS of rules - there is just one symbol.

CFGs → declarative programming

→ specify what is to be computed in terms of rules and let generalized computation mechanisms solve for particular cases.

Grammatical → a sentence in the language

→ A string is grammatical iff there exists a derivation for it.

A good grammar:

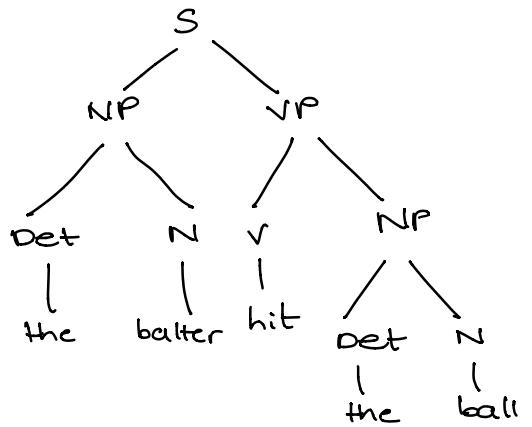
- Doesn't over-generate (high precision)
- Doesn't under-generate (high recall)

Parse tree → One step in discovering relations (grammatical) which can help discover semantic roles

### Myths

- (1) the subj is first noun-phrase in a sentence
  - (2) the subj is the actor in a sentence
  - (3) the subj is what the sentence is about
- often true but not always

constituency tree



dependency tree



## Advantages and Disadvantages

### Advantages of Constituency/Phrase Structure Grammar

- There are widely agreed-upon tests for constituency; there is little agreement above what constitutes a dependency relation
- Constituency maps more cleanly on to formal semantic representations than dependency
- This makes constituency useful in natural language understanding

### Advantages of Dependency Grammar

- It is easier to identify grammatical relations (like subject and object) in a dependency parse
- Dependency parses of sentences having the same meaning are more similar across languages that constituency parses
- Dependency parses are also useful for NLU (ask Google)
- Dependency trees are typically simpler

→ Advantages & disadvantages of constituent & dependency grammar.

CFGs → not prob.  
↳ not very convenient.

## Lecture # 12: CFG

A CFG consists of a set of rules or productions, each of which expresses the ways that symbols of language can be grouped or ordered together.

→ rules are hierarchically embedded, so we can combine the previous rules with others.

## Lecture # 14: Lexical Semantics

lexical semantics → linguistic study of meaning  
→ least useful approach

- Decompositional ⇒ what the components of meaning in a word are
- Ontological ⇒ How meaning of word relates to meaning of other words
- Distributional ⇒ what context the word is found in, relative to other words

Synonymy - equivalence  
 (small, little)  
 Antonymy - opposition  
 (small, large)  
 Meronymy - part-of  
 (liver, body)

Hyponymy - subset; is-a relation  
 (dog, mammal)  
 Hypernymy - superset  
 (mammal, dog)  
 Holonymy - has-a relation  
 (body, liver)

Propositional Meaning → substitutable without changing meaning

→ Define synonymy as a relation b/w senses other than between words.

→ group of antonyms

→ Reverses: some of change or movement in opposite directions

Antonyms → words with opposite meaning defining a binary opposition or are two ends of some scale.

### Ontology

→ describes common words, concepts and relationships b/w concepts used to describe and rep an area of knowledge.

WordNet → lexical resource that organizes words according to their semantic relations.

Synset → A set of words that are roughly synonymous for a particular sense.

In a wordnet, each sense of word is associated with a synset.

Each synset is associated with another through relations → antonymy, hyponymy etc.

consists of three databases → one for nouns  
→ one for verbs  
→ one for adverbs/adjectives

⇒ closed class words are not included.

#### Verb Relations in WordNet

Relation	Definition	Example
Hyponym	From events to superordinate events	$fly^9 \rightarrow travel^5$
Troponym	From a verb (event) to a specific manner elaboration of that verb	$walk^1 \rightarrow stroll^1$
Entails	From verbs (events) to the verbs (events) they entail	$snore^1 \rightarrow sleep^1$
Antonym	Opposites	$increase^1 \leftrightarrow decrease^1$

Figure 19.3 Verb relations in WordNet.

## Noun Relations in WordNet

Relation	Also called	Definition	Example
Hypernym	Superordinate	From concepts to superordinates	<i>breakfast</i> <sup>1</sup> → <i>meal</i> <sup>1</sup>
Hyponym	Subordinate	From concepts to subtypes	<i>meal</i> <sup>1</sup> → <i>lunch</i> <sup>1</sup>
Member Meronym	Has-Member	From groups to their members	<i>faculty</i> <sup>2</sup> → <i>professor</i> <sup>1</sup>
Has-Instance		From concepts to instances of the concept	<i>composer</i> <sup>1</sup> → <i>Bach</i> <sup>1</sup>
Instance		From instances to their concepts	<i>Austen</i> <sup>1</sup> → <i>author</i> <sup>1</sup>
Member Holonym	Member-Of	From members to their groups	<i>pilot</i> <sup>1</sup> → <i>crew</i> <sup>1</sup>
Part Meronym	Has-Part	From wholes to parts	<i>table</i> <sup>2</sup> → <i>leg</i> <sup>3</sup>
Part Holonym	Part-Of	From parts to wholes	<i>course</i> <sup>2</sup> → <i>meal</i> <sup>1</sup>
Antonym		Opposites	<i>leader</i> <sup>1</sup> → <i>follower</i> <sup>1</sup>

Figure 19.2 Noun relations in WordNet.

## Information content

$$IC(c) = \frac{-\log \# \text{ of words that are equivalent to or are hyponyms of } c}{\# \text{ of words in corpus}}$$

## Lecture # 17: Word2Vec

→ predict rather than count

Train a binary classifier on the task:

↓ Is w likely to show up near "apricot"? any word

We take the learned classifier weights as word embeddings

A word s near apricot

- Acts as a gold "correct answer" to the question.

W2V → provides a lot of options.

## Skipgrams with Negative Sampling

- Treat target & neighboring word as +ve examples
- Randomly sample other words to get -ve examples
- Use log regression to train a classifier to distinguish two cases
- Use weights as embeddings

Use sigmoid instead of dot prod. to get prob.

$$f(x) = \frac{1}{1 + e^{-x}}$$

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

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

target context

We have to maximize

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