

(*) compositional semantics :-

- if we understand individual meaning, then we can get overall meaning.
- "the meaning of an expression is function of and only of the meaning of its parts and the way in which the parts are combined."



challenges: (1) syntactic structure.
e.g. house rent &
rent house
Both have diff meanings.

- (2) Semantic Preference
- (3) Polysemy of constituent words.
- (4) Idiomatic & metaphoric usages.

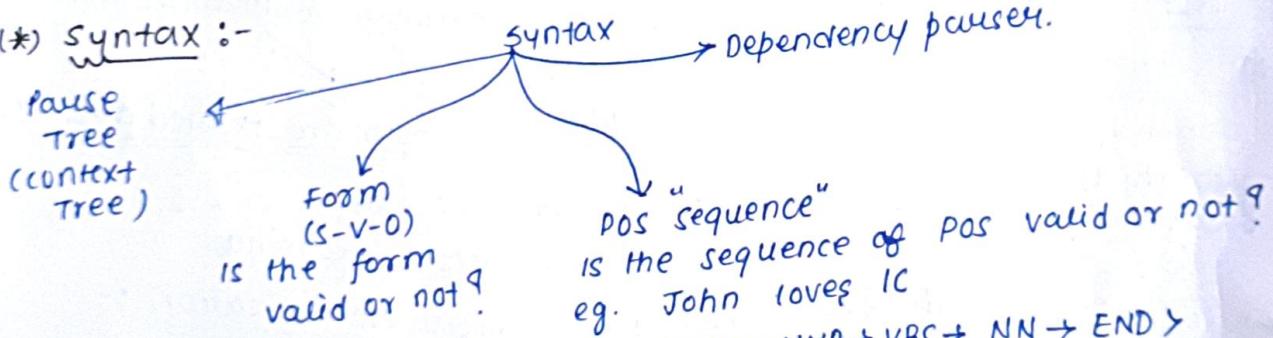
(*) distributional semantics :-

- "words that occur in similar contexts tends to have similar meaning".
- e.g. house and flat frequently occurs with context words
rent, bedroom, sale, etc.
- conditions for two entity to be similar acc. distribn sem:

- (1) Performing same activities / actions.
- (2) Performing same actions under certain condition.
- (3) Frequency of performing activities also same.

(****) "Even though two entity / words are semantically similar and if they does not satisfy above ③ conditions then they are not similar according to distributional semantics".

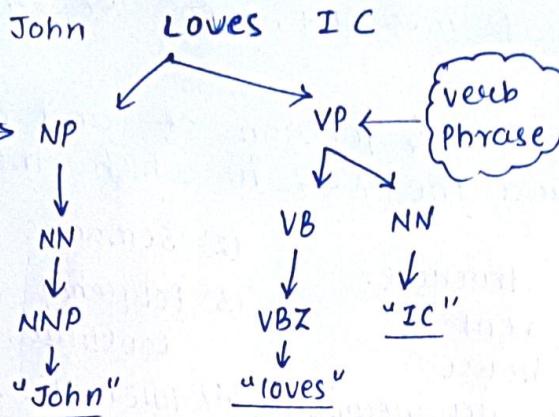
(*) Syntax :-



(****) A sentence without grammar can still have meaning and a sentence without meaning can still have grammar.

e.g. { colorless dreams sleeps furiously }
grammar ✓ but no meaning.

(*) Parse Tree :- (context - Tree)



(***) we can't extract subject, obj from parse tree.

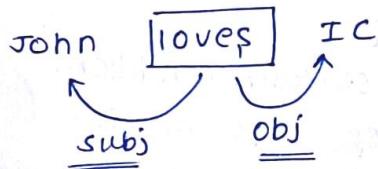
(***) parse tree doesn't tell us the relationship.

→ Parse tree is not necessarily a binary tree.

→ According to Chomsky, any parse tree can be converted to a Binary tree called Chomsky Normal Form (CNF).

(*) Dependency Parser :-

(***) DP tells direct relationship between John and loves & loves and I C.



(***) DP takes longer time than parse tree.

(**) If we want to get subject and object then only we will use dependency parsing else we will use parse tree.

(*) NLP Preprocessing :-

(1) stop words Removal :

- depends on situation whether to remove stop words or keep as it is.

(2) Punctuation removal :

Machine-learning & machine learning

★ "de-hyphenation required".

Both are treated differently

(3) De-abbreviation :

ML → machine learning.

(4) Numerical Normalization :-

50 - fifty

(5) Lemmatization / stemming :-

(**) Lemmatization takes more time than stemming.

saying → say plays → play
playing → play

(**) In neural models,

we don't remove stop words.

in neural model we doesn't take into account the count of words. Thus, it doesn't affect much. whereas in statistical model it can affect.

Neural → compositional + distributional.

- (=) Domain specific Pre-processing :-
- (6) Named Entity Recognition (NER): -
- considering New York, machine Learning collectively, rather than considering two different words...
- (7) Tense Normalization: [2]
- converting everything to present or past tense.
- (8) Plural Normalization:
- convert all words to plural or singular.
- (9) Active passive Normalization:-
- convert all sentences to passive.

(*) Vector Space Modelling :-

→ Distributional semantics claims that they understand meaning because they understand context.

$$\begin{bmatrix} w_{j-2} & w_{j-1} & \text{house} & w_{j+1} & w_{j+2} \end{bmatrix} \quad \begin{bmatrix} w_{i-2} & w_{i-1} & \text{flat} & w_{i+1} & w_{i+2} \end{bmatrix}$$

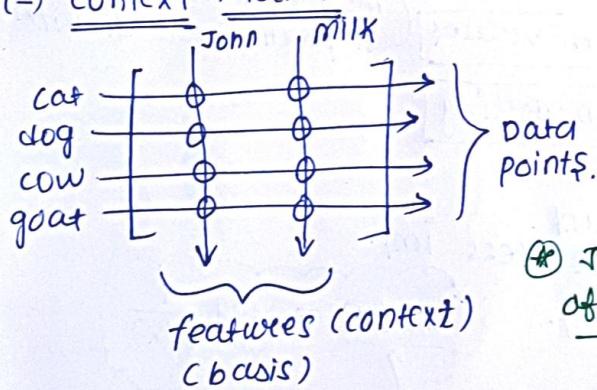
} if these two words satisfies ③ conditions of similarity of distribution semantics then house and flat are similar.

• diff. Distribution: if there is diff distribution then acc. to distribn semantics they will be not considerd similar even though they are similar in reality. ↗

Drawback of Distributional Sem.

- "Represent a linguistic unit" as a vector."

(=) Context Matrix :-



- Given a set of features and set of datapoints

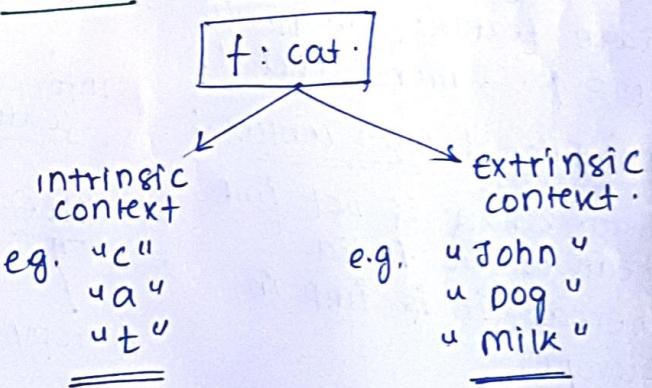
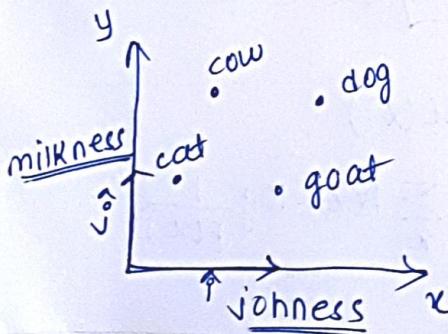
$$\hookrightarrow (f_1, f_2, \dots, f_m)$$

$$\hookrightarrow (w_1, w_2, \dots, w_{i+1})$$

$f_1: \text{John}$ } words associated
 $f_2: \text{Milk}$ } with each datapoint.

$\text{John} \neq \text{John-ness}$.

④ John-ness is unique characteristic of John. initially, $\underline{\text{John}} = \underline{\text{John-ness}}$



(=) Context Engineering:

for word: EXT:- word
INT:- characters.

For phrase:

EXT: phrase
INT: words.

For sentence:

EXT: sentence
INT: words / phrases.

Q. When should we use extrinsic and when should we use intrinsic features?

→ On document level, we use "intrinsic".

On document level, extrinsic context makes less sense.
e.g., document created on same date and time are in same collection but not similar.

→ If two doc are in same folder we cannot claim that they are similar. As we don't know how the folder was created.

→ Folders are not always organized based on content.

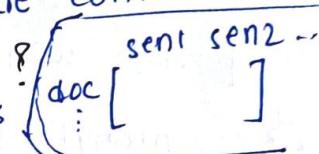
→ The moment we go outside doc. we go outside content.

Q. Using sentences as features is preferable?

- We are computing similarity b/w two docs and asking that 2 doc should contain

same sentences, chances of which are very low.

- ∵ usually words / phrases are used as features...



(*) How to Initialize Mc matrix?

- Linear transformation is performed by Mc, hence it needs to be learned well.
- Initialize using
 - Random values (maybe belonging to a particular distribn)
 - TF-IDF modelling

• TF-IDF Modelling :-

- Rare features are more important.
- Features occurring in all docs are less imp.
- Rare features helps in clustering.
- IDF ↑ → more special.

$$IC(x) = \log \frac{1}{P(x)}$$

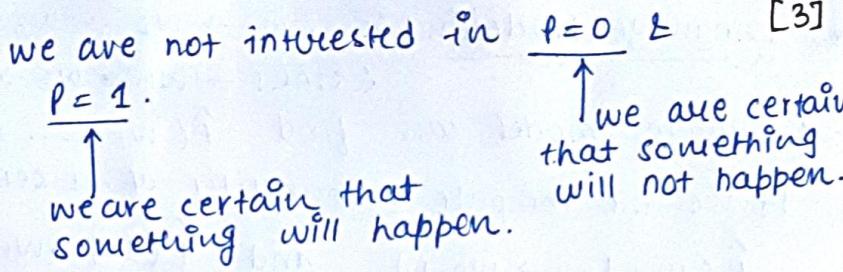
information content

• conditional Info^n content :-

- Autorikshaw is not imp. feature in India. whereas it is imp in Europe.

$$H(x) = E(IC(x))$$

$$\text{Entropy} = \sum_x P(x) \cdot \log \frac{1}{P(x)}$$



Inverse Document Frequency :-

Frequency :-

$IDF \leftarrow \log \frac{|ID|}{\# di}$

$\# di \leftarrow \text{no. of docs in which } f_i \text{ is present}$

$$TF(\text{cat}) = \sum_{di} \frac{|W_{\text{cat}} \in di|}{|W \in di|}$$

information content.

$$f_i = \text{term} \quad w_j^o = \text{cat.}$$

$$P(f_i | w_j^o) = \frac{\text{count}(f_i^o, w_j^o)}{\text{count}(w_j^o)}$$

for doc, w_j^o becomes d_i^o

- while finding $P(f_i | w_j^o)$, we don't consider all other w_j^o 's and f_i^o 's.

- Each calculation is independent in terms of rows and columns.

(**) Assumption: All features are independent of each other.

- for this assumption,

If we have seq. "John loves IC". If we want to predict "IC" after "John loves", then if we assume that they are all completely independent of each other then model will not be able to predict "IC".

- All of them must have some dependence. Hence unigram models are worst language models.

(***) But it is a good way to initialize, rather to initialize with random values. It is a kind of inductive bias. doing across all w_j^o 's

$$H(f_i | w_j^o) = P(f_i^o | w_j^o) * \log \frac{1}{P(f_i)}$$

$$H(f_i) = \sum_j H(f_i | w_j^o)$$

This assumption is also known as "bag of words" assumption.

(*) Language Modelling :- $w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \dots \rightarrow w_n >$
 ↳ START $\rightarrow w_1 \rightarrow w_2 \rightarrow w_3 \rightarrow \dots \rightarrow w_n \rightarrow END$

- Language model can find $\hat{P}(w_1 \rightarrow \dots \rightarrow w_n) = ?$
 If we can compute probability of a sentence then we can

$$\frac{\hat{P}(w_n | w_1 \rightarrow w_{n-1})}{\uparrow} \quad \text{and} \quad \frac{\hat{P}(w_i \rightarrow w_n | w_1 \rightarrow w_{i-1})}{\uparrow}$$

Predicting next word of the sequence.

Gram Modelling

Predicting next sequence given a sequence of words.

$$P(w_1 \rightarrow w_2 \rightarrow w_3) \neq P(w_1 w_2 w_3)$$



$$= P(w_1) * P(w_2 | \leftarrow w_1) \\ * P(w_3 | \leftarrow w_2 \leftarrow w_1)$$

$$= P(w_1) * P(w_2 | w_1) * P(w_3 | w_2 w_1)$$

sequence does not matter
in this case.

start from begin

$$\frac{P(w_3 | w_2 w_1)}{\uparrow}$$

$$\frac{P(w_3 | \leftarrow w_2 \leftarrow w_1)}{\uparrow}$$

search in a doc where w_1 and w_2 occurs together and w_3 occurs after w_1 .

start from end.

search in a doc where w_1 and w_2 occurs together.

Trigram modelling.

$P(w_3 | \leftarrow w_2 \leftarrow w_1) \Rightarrow$ gram modelling.

$P(w_3 | w_2) \Rightarrow$ Bigram modelling.

$P(w_3) \Rightarrow$ Unigram modelling.

$$\frac{\hat{P}(w_n | \leftarrow w_{n-1} \leftarrow \dots \leftarrow w_1)}{\uparrow}$$

n-gram modelling.

(computationally expensive)

(*) Markov Assumption :-

(***) Default degree of markov assumption = 1.

$$\hat{P}(w_n | \leftarrow w_{n-1} \leftarrow \dots \leftarrow w_1) \approx P(w_n | w_{n-1})$$

$P(c | \leftarrow b \leftarrow a)$ approxes

$$P(c | \leftarrow b)$$

$$P(A \rightarrow B \rightarrow C) = P(A) * P(B | \leftarrow A) * P(C | \leftarrow B \leftarrow A)$$

$$P(A \rightarrow B \rightarrow C) = P(A) * P(B | \leftarrow A) * P(C | \leftarrow B) \quad [\because \text{Bigram modelling}]$$

Problem with Bigras:

e.g. John who lives in NY and works at Google loves IC.

'IC' predicted only based on loves...

we can use Maximum Likelihood Estimate (MLE),

[4]

$$P(IC | loves) = \frac{\text{count}(IC, loves)}{\text{count}(loves)}$$

(***) what if 'IC' and 'loves' have not come together yet.

$$\text{count}(IC, loves) = 0.$$

model will predict that IC can never occur with loves.

SOLN: use various smoothing techniques...

we need,

$$\hat{P}(IC | \leftarrow loves) = \frac{\text{count}(IC \leftarrow loves)}{\text{count}(loves)}$$

Q. what if there is out of vocabulary situation?

- solution: use smoothing techniques...

(**) 'BOW' contradicts with language model.

↳ can be used where no seq. required.

(*) Skip-Gram Modelling :- $\hat{P}(w_1 \rightarrow w_2 \rightarrow \dots \rightarrow w_i \rightarrow \dots \rightarrow w_n)$
In trigram modelling,
e.g. John → loves → IC

depends on previous
as well as future words.

$$P(IC | \leftarrow loves \leftarrow John) = \frac{\text{count}(John \rightarrow loves \rightarrow IC)}{\text{count}(John \rightarrow loves)}$$

 NO significant no. of
NO significant no. of
It may happen that
than 'loves'.
John → loves sequence.
John → loves → IC.
'John' has more influence on 'IC'

strong influence 
Boost whenever John comes
Supress whenever John comes
IC comes.
IC won't come.

(*) Skip-Gram (word2vec) :-
 $P(w_{t+j}^c | w_t)$ w_t : target word w_c : context word.

e.g.: John → loves → IC

$$\{ w_{John}^c, \underbrace{w_{loves}^c}_{\text{center word.}}, w_{IC}^c \}$$

$m=1$ ← context window.

$$\prod_{j=-m}^m (P(w_{t+j}^c | w_t))$$

$$P(John \rightarrow loves \rightarrow IC)$$

$$= \prod_{-m}^n P(w_{t+j}^c | loves)$$

$$= P(John | loves) \cdot P(IC | loves)$$

(**) skip-gram assumption:
conditional independence
of w_c^c .

$$P(\text{seq}) = \prod_{w^t=1}^T \left(\prod_{j=-m}^m P(w_{t+j}^c | w^t) \right)$$

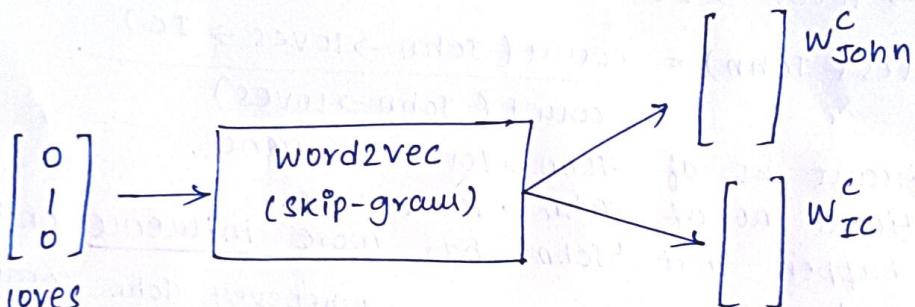
\rightarrow In word2vec we initially starts with word of vector space and moves to probability space.

initially, we don't have the co-ordinates of "John", "loves" and "IC".

Vocab = {"John", "loves", "IC"}

Initially, we assume that

$$\begin{aligned} \text{Initial } \left\{ \begin{array}{l} \text{John} \equiv \text{Johnness} \\ \text{loves} \equiv \text{love ness} \\ \text{IC} \equiv \text{IC-ness} \end{array} \right. & \rightarrow \begin{array}{l} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \\ \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \\ \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \end{array} \left. \begin{array}{l} \text{one hot vectors.} \\ \text{length of} \\ \text{one-hot vector} = |\text{V}| \end{array} \right. \end{aligned}$$



w^t

- Word2vec gathers all of its friends.
- Take i/p as target word and gather all the content words.
- While training model learns the ideal position of target words.

\Rightarrow Represent w_1, w_2, w_3 context matrix as one-hot vector matrix.

	w_1	w_2	w_3	w_{all}
w_1	1	0	0	0
w_2	0	1	0	0
w_3	0	0	1	0
\vdots	\vdots	\vdots	\vdots	\vdots
w_{all}	0	0	0	1

- same words acts as data point as well as a feature.

Data points.

feature