

Stanford - CS224n - Lecture 1 - Introduction and word vectors

*) language and social interactions make humans different.

Mankind is like a human network and language is an important means of communication. (knowledge → intelligence)



The way we acquire knowledge is from our ancestors, through books (their observations), through language.



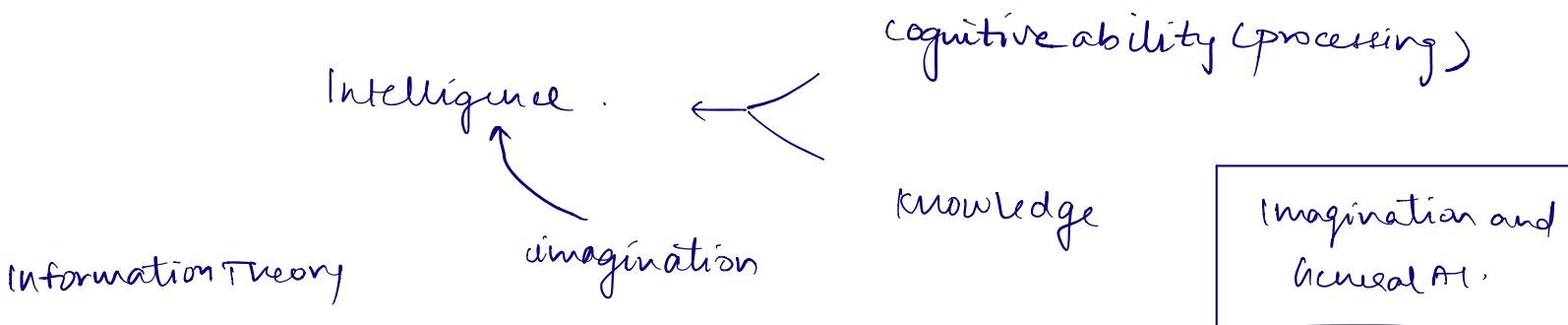
Language can be seen as one of the most important differentiator between humans and other animals.

Even animals with a good sense of control and vision.



Again, language and social interaction can be seen as the key difference, the way of passing knowledge, the way of acquiring intelligence.

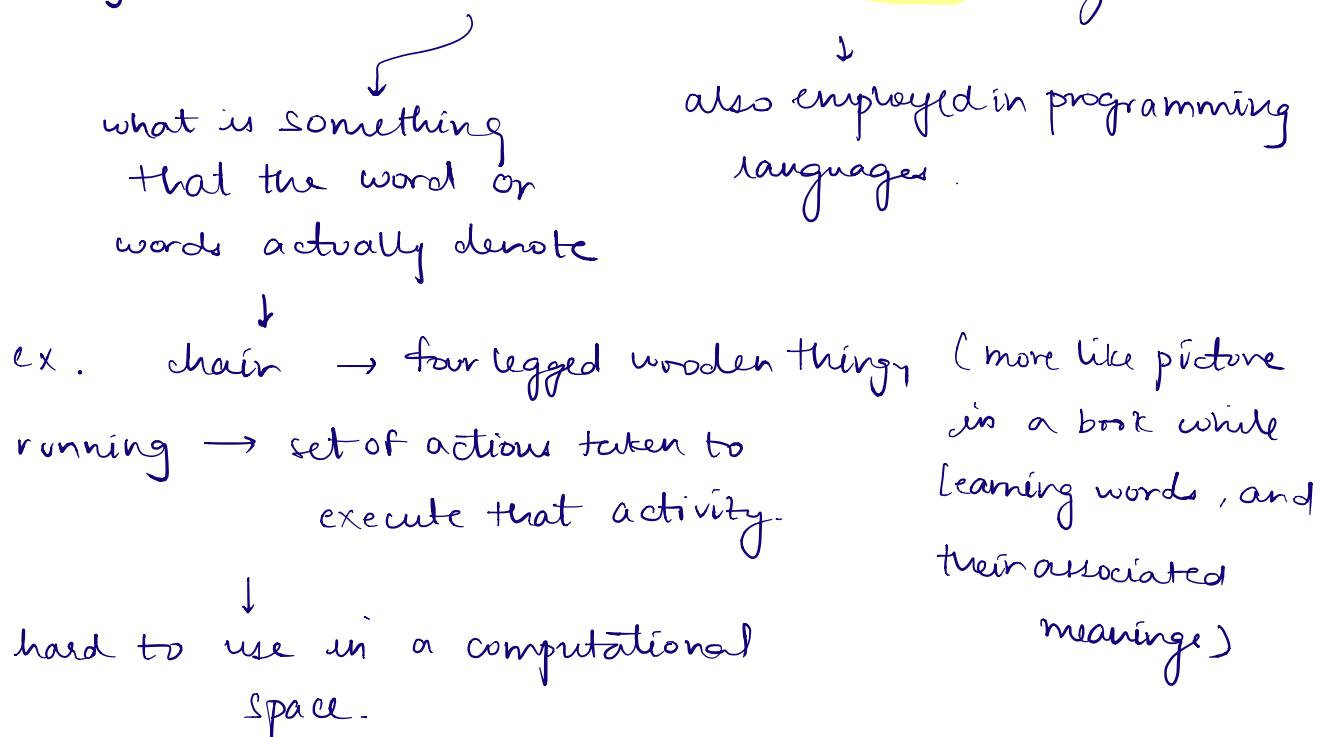
sides



-
- *) A gamechanger was writing, that allowed knowledge to be transferred temporally and spatially. (5000 yrs old)
Very recent on scale of evolution. Enable major human progress.

- * The human communication network is pathetically slow as compared to what computer networks have been able to achieve
- * But we have a sort of adaptive compression mechanism, wherein, we need a very few bit sequences to convey information. For example, by passing a basic explanation about a scene, we are able to imagine the picture in our minds. This also works as we assume the person on the other hand to have a similar amount of acumen to understand what we're trying to convey.

* meaning → idea → denotational semantics (linguists)



* One way of getting a sense of meaning has been through the use of Wordnet, but it didn't give good results for a set of reasons. It contains set of synonyms for words to get a sense of meaning along with pos tags based on different contexts. Hard to maintain, and also doesn't give a good sense of similarity b/w words in two different

similarity sets.

- *) Traditional NLP (pre 2012 period, before Neural Networks
(localised approach)
were being used) used one-hot-encoded vectors, where
1 was assigned to the word corresponding to its index in
a vocabulary. Issues:

-) very large vocabularies \rightarrow very large vectors
-) understanding relationships b/w words.

for example if we search for "Seattle motel", we would
like to match documents containing "Seattle hotel",
since "motel" and "hotel" represent basically the
same thing.

↓

they have no similarity notion.

↓

People tried to rely on wordnet building similarity tables,
and owing to the large size of vocabulary would be
even bigger.

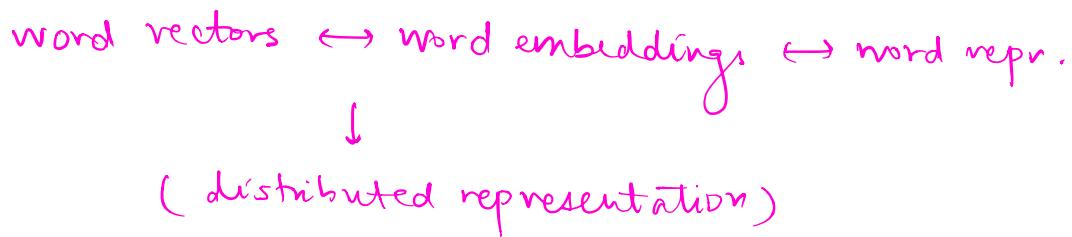
- *) Distributional semantics \rightarrow If we can say on which context,
the word can be used, and in which context it can't be,
then we can say, we have an understanding of the meaning
of the word.

↓

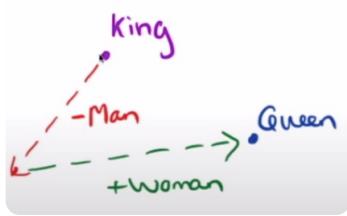
context represents meaning *(word vectors)*

- *) Owing to the above idea, we build a **dense vector** for
each word, chosen so that similar words appear in

similar contexts.



- * Numeric vector, smallish, dense. meaning distributed along the dimension of the vector. min: 50,



on personal mle: 300

max out performance: 1000s.

↓

Dimensions of meaning
in this
vector
space.

many magnitudes smaller than the vectors.

- i) Now, since words are being represented using vectors., they can be plotted in the particular vectorspace.

The closer words will correspond to being more similar, as compared to others. This can also be visualized by projecting words in a 2D space.

Note : NLP people call large pile of text a corpus.

- * word2vec: an algorithm to learn the dense vector representations of words. Basically done by going over a large corpus of text. for each position t in the text, we can consider a center word ' c ', and context words represented as ' o '.

↓

we use the similarity between the word vectors of ' o '

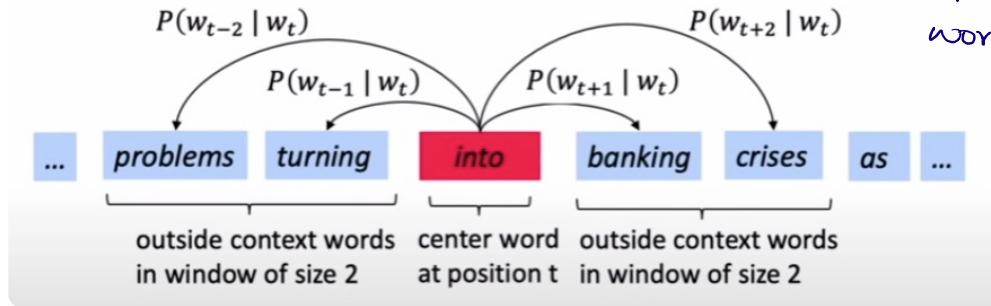
and 'c' to calculate the probability of o given c .

↓

we keep adjusting the word vectors to maximise this probability.

context of $into \rightarrow$ meaning of $into$

represented by the word vector.



*) objective fn :

for each position $t = 1 \dots T$, predict context words within a window of fixed size, m , given center word, w_j .

$$L(\theta) = \prod_{t=1}^T \prod_{\substack{-m \leq j \leq m \\ j \neq 0}} P(w_{t+j} | w_j)$$

Likelihood function represents parameters to be optimized.

→ going through all the words one by one in the corpus

conditional probability of context word wrt center word

dependent on the parameters.

↓

vector representation of words.

$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \log P(w_{t+j} | w_j)$$

objective fn

minimizing objective fn \leftrightarrow maximizing predictive accuracy.

Log \rightarrow summation instead of product

(solves the problem of very small values,
due to multiplication of values b/w of 1)

* How to calculate $P(w_{t+j} | w_t; \theta)$?

Two vectors per word ; one representation as center word ,
other as context word .

v \rightarrow center

u \rightarrow context

Then for center word c and a context word o

$$P(o|c) = \frac{\exp(v_o^T \cdot v_c)}{\sum_{w \in V} \exp(v_w^T \cdot v_c)}$$

measure of similarity
large dot product, large probability

$\underbrace{\qquad\qquad\qquad}_{\text{Softmax Distribution}}$

normalize over entire vocabulary.

↓
maps numbers to a probability distribution

* Softmax because it assigns max probability for largest individual value , but does so in a manner that even the smallest values have some probability assigned to them.

* Optimization of parameters to minimize the objective vector. θ represents all model parameters, in one long vector. for d dimensional word representations,

V size of vocabulary and 2 vectors for each word,
the long vector is of the dimension $(2V+d)$.

$$\theta = \begin{bmatrix} v_{aardvark} \\ v_a \\ \vdots \\ v_{zebra} \\ u_{aardvark} \\ u_a \\ \vdots \\ u_{zebra} \end{bmatrix} \in \mathbb{R}^{2dV}$$

$$= \frac{\partial}{\partial v_c} \log \frac{\exp(v_0^T \cdot v_c)}{\sum \exp(v_w^T \cdot v_c)}$$

$$= \frac{\partial}{\partial v_c} \log \exp(v_0^T \cdot v_c) - \frac{\partial}{\partial v_c} \log \sum \exp(v_w^T \cdot v_c)$$

(property of the log)

$$\frac{\partial}{\partial v_{c_1}} (v_{0,1}v_{c_1} + v_{0,2}v_{c_2} + \dots + v_{0,n}v_{c_n}) = v_{0,1}$$

↓
Similarly for complete v_c

$$\frac{\partial}{\partial v_c} (v_0^T \cdot v_c) \rightarrow \text{numerator}$$

↓
multivariate calculus

$$\frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^\top \cdot v_c) = \sum_{x=1}^V \frac{\exp(u_x^\top \cdot v_c) \cdot u_x}{\sum_{w=1}^V \exp(u_w^\top \cdot v_c)}$$

↓

denominator

↓

$$\frac{\partial}{\partial v_c} \log P(x/c) = v_0 - \sum_{x=1}^V \frac{\exp(u_x^\top \cdot v_c) \cdot u_x}{\sum_{w=1}^V \exp(u_w^\top \cdot v_c)}$$

$$= v_0 - \sum_{x=1}^V \frac{\exp(u_x^\top \cdot v_c)}{\sum_{w=1}^V \exp(u_w^\top \cdot v_c)} \cdot u_x$$

$$= v_0 - \sum_{x=1}^V p(x/c) \cdot u_x$$

context word

expected context ↗

convex sum

↓

difference b/w the actual context word

and expected context word