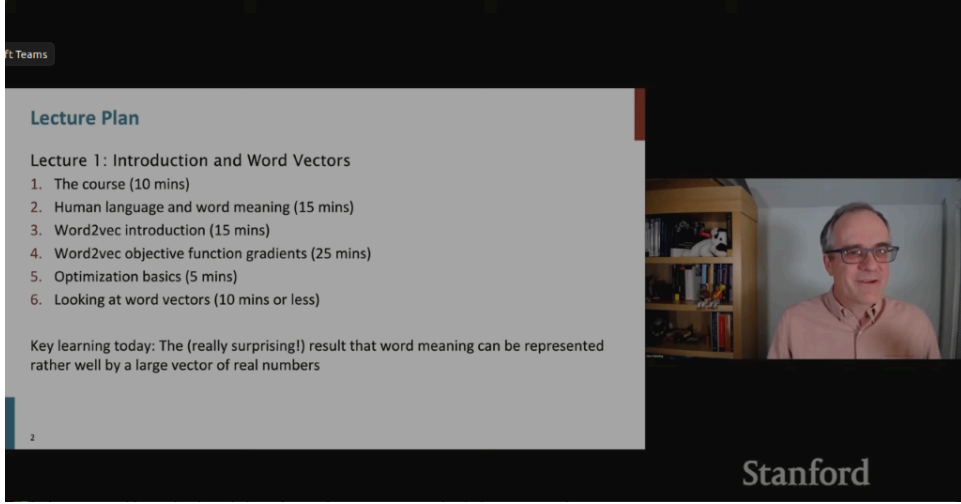


Lecture Outline:



Intro to NLP

Human Languages:

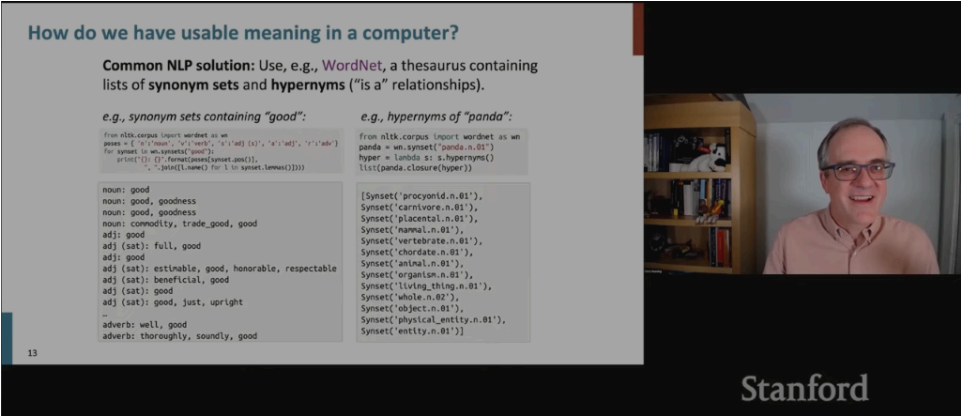
- Languages are means of sharing ideas, facts, intents etc.
- Languages have evolved/ use of languages have evolved to be useful for communication.
- Children vs Machines: Children learn a language very easily by interacting with the (multi-modal) world as compared machine. The Machines (advances in NLP) have no-where near the language-acquisition ability of Children.
- How to represent a language for better understanding of a Machine?? -> Deep learning provides the tools.
- The most important question to be answered by this course is -> How do we represent words? (to be better understood by the machines)

Uses of NLP:

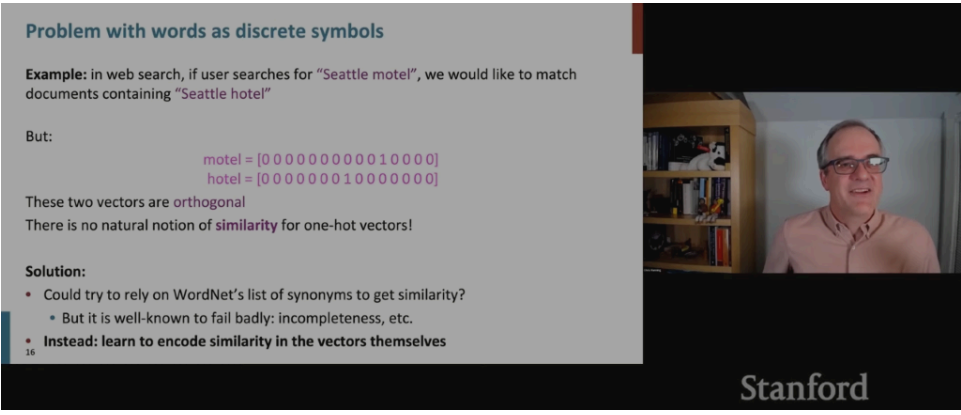
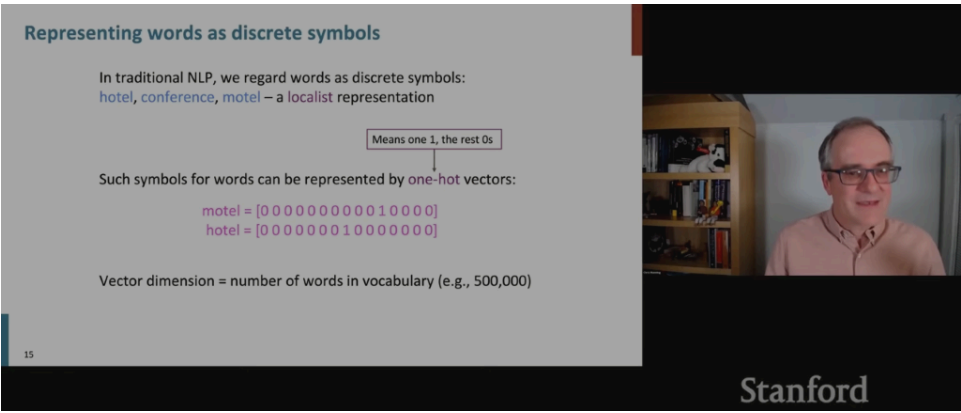
- Machine Translation
- GPTs – Question answers, text generation and information retrieval
- Speech-to-text

Representing meaning of the word:

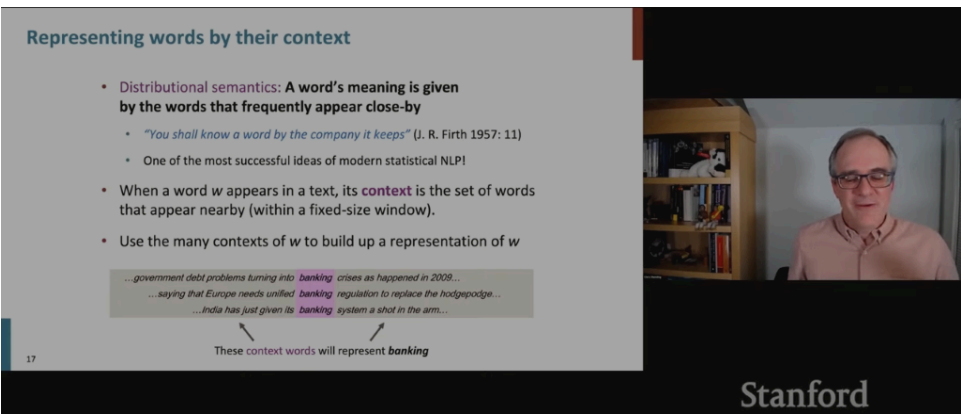
- What is Meaning? -> an idea behind a physical/imaginary concept
- Linguistic way of thinking about meaning?
  - A bond between signifier (i.e. the word, e.g. a chair) and signified (an idea behind the word, e.g. sitting)
- How to have usable meaning in computers? -> A dictionary



- Problem with WordNet:
  - Hard to always keep it up-to-date.
  - Details are missing: e.g. Good and Proficient have same meaning in some context but not always same.
- Problems in traditional NLP ways:
  - Representing word as a discrete symbol and associated problems:
    - In short, the discrete vectors do not include context of the word and thus vectors of the similar words in context can be orthogonal



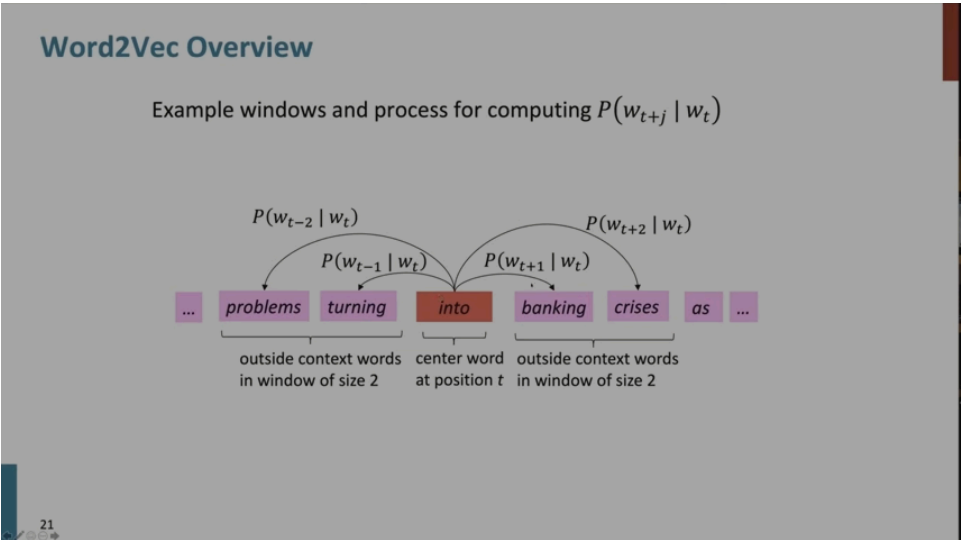
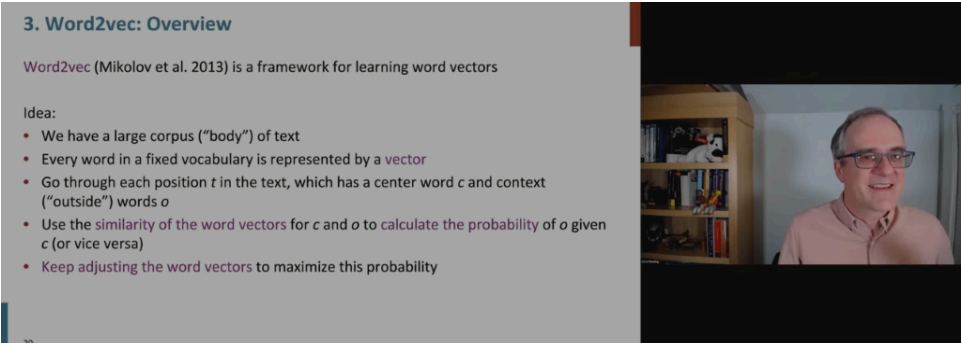
- Representing words by their context:



Word vectors: meaning of the word in context so that a word w1 has similar vector to word w2 if both of them are used in similar context.

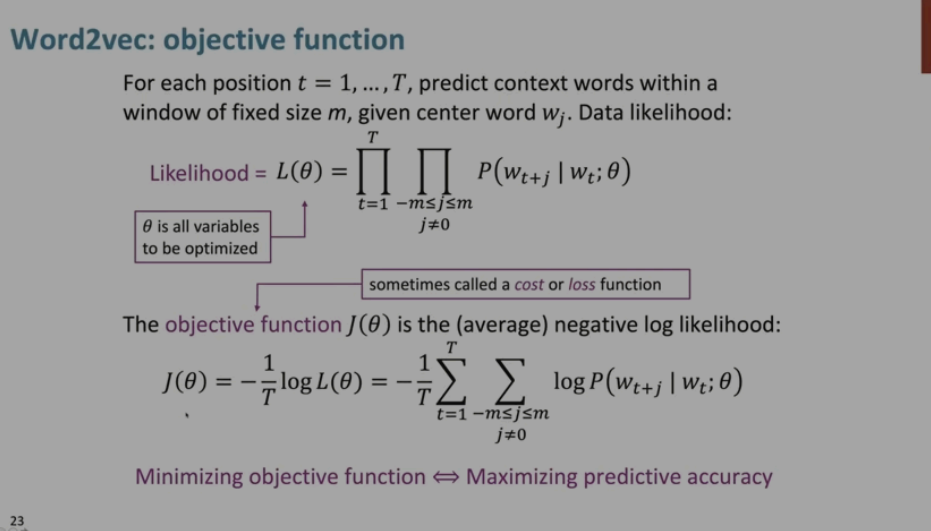
Word2Vec:

- Overview:
  - Train a model that generates similar vector embeddings for words having similar in-context meaning.
  - During training, each word plays two roles. It becomes center word and also a context word.
  - Simple dot-product similarity is used to calculate similarity score between center-word and context-words vectors.
  - This similarity score is then converted to a probability distribution by using softmax

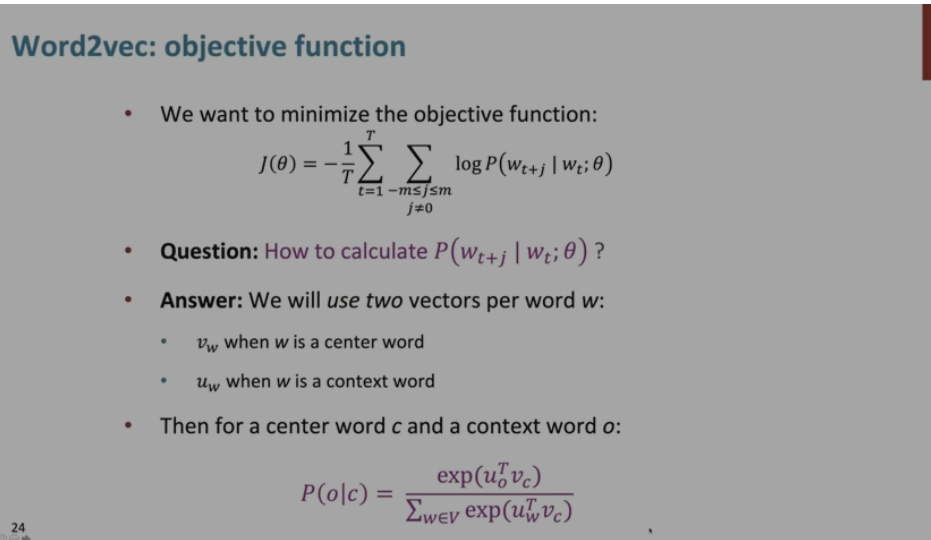


- Likelihood Function:
  - Pick each word in the corpus as center word (w) and then predict the probability of context words (w<sub>i</sub>) within a window of size m.
  - Repeat this process for all words and then take product.
- Loss function:
  - Likelihood has to be maximized but the loss (the difference between model predicted likelihood and Ground truth likelihood) has to be minimized)
  - Hence take log of likelihood function and minimize the negative average of that log.

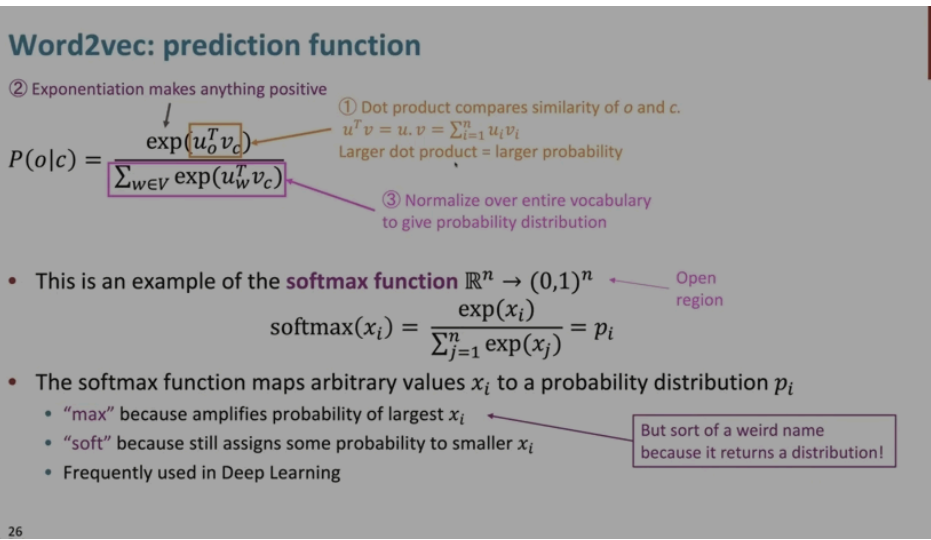
Hint:  $\text{Log}(a \cdot b) = \text{log}(a) + \text{log}(b)$



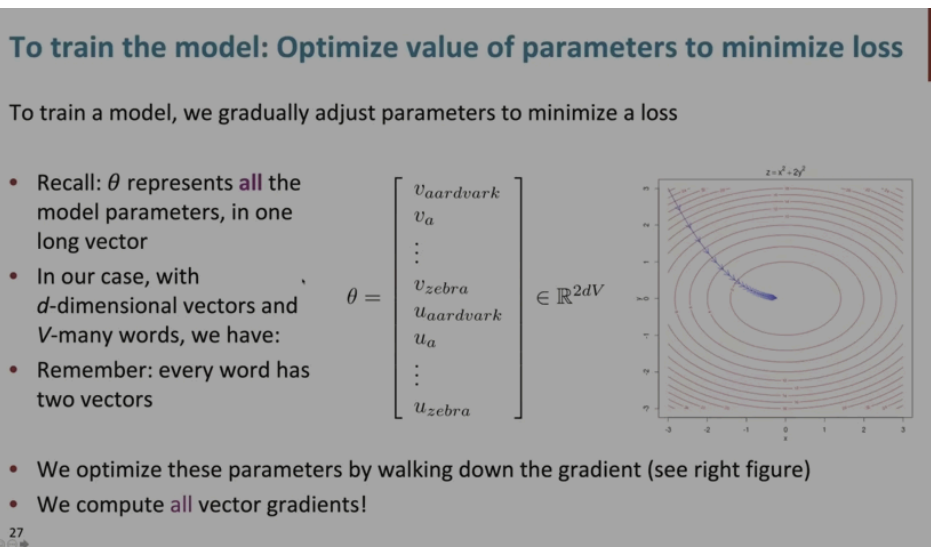
- Calculating Probability of Context words given a center words
  - Why to assign two different vectors to a word when it is used a center word and a context word?



- Use softmax to convert the word vectors into probabilities:



- Training of Word2Vec model (just like any other model)



Partial derivative of Word2Vec objective function.

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{j \neq 0} \log p(w_{t+j} | w_t; \theta)$$

$$\therefore p(w_{t+j} | w_t) = \frac{\exp(u_0^T \cdot v_c)}{\sum_{w=1}^V \exp(u_w^T \cdot v_c)}$$

$$\frac{\partial}{\partial v_c} \log \frac{\exp(u_0^T \cdot v_c)}{\sum_{w=1}^V \exp(u_w^T \cdot v_c)} \rightarrow \frac{\partial \left( \log \frac{a}{b} \right)}{\partial z} = \log(a) - \log(b)$$

$$\frac{\partial}{\partial v_c} \log (\exp(u_0^T \cdot v_c)) - \log \left( \sum_{w=1}^V \exp(u_w^T \cdot v_c) \right)$$

(I)                      (II)

$$\textcircled{I} \log \exp(a) = a.$$

↓

$$\frac{\partial}{\partial v_c} (u_0^T \cdot v_c)$$

↓

applying product Rule

↓

$$u_0^T \cdot \frac{\partial}{\partial v_c} v_c + v_c \cdot \frac{\partial}{\partial v_c} u_0^T$$

↓

$$u_0^T$$

$$\frac{\partial}{\partial v_c} \log \left( \sum_{w=1}^V \exp(u_w^T \cdot v_c) \right) \left. \begin{array}{l} \text{chain rule} = \frac{df(z)}{dz} = \frac{df}{dz} \cdot \frac{dz}{dv_c} \\ \uparrow \\ f(z) \end{array} \right\}$$

$$\frac{1}{\sum_{w=1}^V \exp(u_w^T \cdot v_c)} \cdot \frac{\partial}{\partial v_c} \sum_{w=1}^V \exp(u_w^T \cdot v_c) \left. \begin{array}{l} \frac{d}{da} \log(a) = \frac{1}{a} \leftarrow \\ \text{applying chain rule again.} \end{array} \right\}$$

$$\sum_{w=1}^V \frac{\partial}{\partial v_c} \exp(u_w^T \cdot v_c)$$

$$\sum_{w=1}^V \exp(u_w^T \cdot v_c) \cdot \frac{\partial}{\partial v_c} u_w^T \cdot v_c$$

$$\sum_{w=1}^V \exp(u_w^T \cdot v_c) \cdot u_w$$

$$\frac{\sum_{w=1}^V \exp(u_w^T \cdot v_c) \cdot u_w}{\sum_{w=1}^V \exp(u_w^T \cdot v_c)}$$

Putting (I) & (II) together,

$$u_0 = \frac{\sum_{w=1}^V \exp(u_w^T \cdot v_c) \cdot u_w}{\sum_{w=1}^V \exp(u_w^T \cdot v_c)}$$

$$u_0 = \sum_{w=1}^V \frac{\exp(u_w^T \cdot v_c)}{\sum_{w=1}^V \exp(u_w^T \cdot v_c)} \cdot u_w \rightarrow \text{softmax from input.}$$

$$u_0 = \text{expected}$$

$$\therefore \text{observed} = \text{expected}$$