

<https://web.stanford.edu/class/cs224n/>

<https://www.youtube.com/watch?v=rmVRLeJRkl4&list=PLoROMvody4rOSH4v6133s9LFPRHjEmbmJ&index=2>

for word embedding if we use one hot ---→ no

در نمایش "one-hot" هر واژه به یک بردار دودویی تبدیل می‌شود، که تنها یک عضو آن برابر با ۱ است و سایر عناصر برابر با ۰ هستند. به عبارت دیگر، هر واژه در یک فضای بلند بُعدی نمایش داده می‌شود و تنها یک بُعد از آن یک است.

مشکلات این نمایش:

عدم درک ارتباطات معنایی: این نمایش اطلاعات معنایی بین واژگان را نمی‌فهمد. دو واژه با این نمایش تماماً مستقل از یکدیگر هستند و هیچ اطلاعاتی در مورد شباهت یا تفاوت‌های معنایی آنها در این نمایش نهفته نیست.

اندازه بالای فضای ویژگی‌ها: اگر تعداد واژگان زیاد باشد، این نمایش باعث ایجاد یک فضای بسیار بزرگ می‌شود که نه تنها می‌تواند مشکلات محاسباتی ایجاد کند، بلکه به اندازه کافی اطلاعات مفید برای یادگیری ارائه نمی‌دهد.

عدم درک تفاوت‌ها و شباهت‌ها: این نمایش تفاوت‌ها و شباهت‌های معنایی بین واژگان را نمی‌فهمد. مثلاً فاصله بین دو بردار "one-hot" برابر با ۲ خواهد بود، بدون اینکه به ما بگوید این تفاوت ناشی از چه معناست.

اضافه‌وزنی اطلاعات: این نمایش به هر واژه یک بُعد دارد و اطلاعات معنایی بسیار محدودی ارائه می‌کند. این در حالی است که بسیاری از واژگان در زبان‌ها ارتباطات پیچیده‌تری دارند که با نمایش "one-hot" به درستی نمی‌توان آنها را ادراک کرد.

به همین دلیل، مدل‌های مبتنی بر نمایش‌های جدیدی مانند "word embeddings" (تعبیر واژگان) استفاده می‌شوند که این مشکلات را حل کرده و واژگان را در یک فضای چگال و معنایی بهتر نمایش می‌دهند.

ی روش دیگر همیشه featurized word embedding که با یک سری فیچر بیایم هر word رو نمایش بدیم اینطوری مفهوم sim را داریم

wordnet

WordNet [Miller, 1995]:

WordNet یک پایگاه داده لغت‌نامه است که ارتباطات معنایی بین واژگان را ارائه می‌دهد.

برای واژگان، مترادف‌ها، زیرمجموعه‌ها (hyponyms)، و ارتباطات معنایی دیگر را برچسب‌گذاری می‌کند.

UniMorph [Batsuren et al., 2022]:

UniMorph یک منبع دیگر است که برای اطلاعات مربوط به ساختار زیر واژگانی (morphology) در بسیاری از زبان‌ها برچسب‌گذاری می‌کند.

اطلاعات مربوط به ساختار زیر واژگانی شامل جزئیاتی مانند پیشوند، پسوند، یا تغییرات دیگر در واژه‌ها می‌شود.

مشکلات:

کمبود منابع انسانی:

منابع با برچسب انسانی همواره در حجم واژگان نسبت به روش‌هایی که می‌توانند واژگان را از منابع متنی طبیعی استخراج کنند، کمبود دارند.

به روزرسانی این منابع هزینه‌بر است و همواره ناقص می‌مانند.

تضاد بین ابعاد و کارایی تعبیه (Embedding):

برای نمایش همه دسته‌ها نیاز به بردار با ابعاد بسیار بالا (بیشتر از اندازه واژگان) وجود دارد. روش‌های نورونی مدرن که عمدتاً با بردارهای فشرده (dense vectors) کار می‌کنند، با این بردارهای بزرگ به خوبی هماهنگی ندارند.

کمبود دقت در نمایش‌های ایده‌آل برای متن:

اشاره به این نکته شده که دیدگاه انسانی در مورد نمایش مناسب برای متن، به عنوان یک پیشنهاد مکرر در دوره، معمولاً عملکرد کمتری دارد نسبت به روش‌هایی که به داده امکان می‌دهند بیشتری از جزئیات را مشخص کنند، حداقل زمانی که داده به میزان زیادی برای یادگیری وجود دارد.

similarity

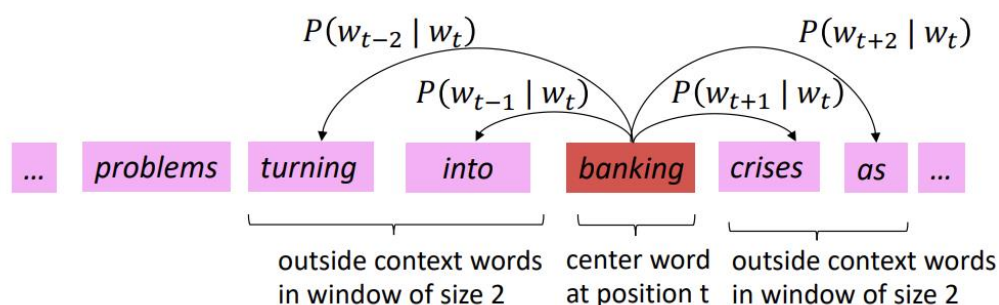
context word

distributional semantic

You shall know a word by the company it keeps.

word2vec 2013

Example windows and process for computing  $P(w_{t+j} | w_t)$



## Word2vec: objective function

For each position  $t = 1, \dots, T$ , predict context words within a window of fixed size  $m$ , given center word  $w_t$ . Data likelihood:

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

$\theta$  is all variables to be optimized

sometimes called a *cost* or *loss* function

The **objective function**  $J(\theta)$  is the (average) negative log likelihood:

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

Minimizing objective function  $\Leftrightarrow$  Maximizing predictive accuracy

- **Question:** How to calculate  $P(w_{t+j} | w_t; \theta)$ ?
- **Answer:** We will use two vectors per word  $w$ :
  - $v_w$  when  $w$  is a center word
  - $u_w$  when  $w$  is a context word
- Then for a center word  $c$  and a context word  $o$ :

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

28

## Word2vec: prediction function

② Exponentiation makes anything positive

① Dot product compares similarity of  $o$  and  $c$ .  
 $u^T v = u \cdot v = \sum_{i=1}^n u_i v_i$   
 Larger dot product = larger probability

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

③ Normalize over entire vocabulary to give probability distribution

- This is an example of the **softmax function**  $\mathbb{R}^n \rightarrow (0,1)^n$  ← Open region

$$\text{softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^n \exp(x_j)} = p_i$$

- The softmax function maps arbitrary values  $x_i$  to a probability distribution  $p_i$ 
  - “max” because amplifies probability of largest  $x_i$
  - “soft” because still assigns some probability to smaller  $x_i$
  - Frequently used in Deep Learning

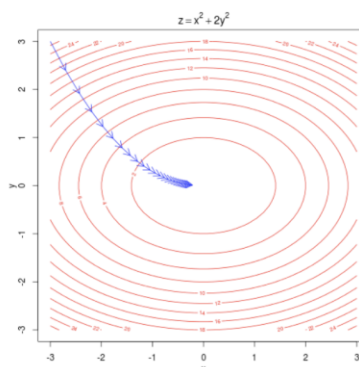
But sort of a weird name because it returns a distribution!

## To train the model: Optimize value of parameters to minimize loss

To train a model, we gradually adjust parameters to minimize a loss

- Recall:  $\theta$  represents **all** the model parameters, in one long vector
- In our case, with  $d$ -dimensional vectors and  $V$ -many words, we have  $\rightarrow$
- Remember: every word has two vectors

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



- We optimize these parameters by walking down the gradient (see right figure)
- We compute **all** vector gradients!

Remember: every word has two vectors  $\rightarrow u$  and  $v$

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

$$= \underbrace{\frac{\partial}{\partial v_c} \log \exp(u_o^T v_c)}_{(1)} - \underbrace{\frac{\partial}{\partial v_c} \log \sum_{w=1}^V \exp(u_w^T v_c)}_{(2)}$$

$$(1) \quad \frac{\partial}{\partial v_c} \log \exp(u_o^T v_c) = \frac{\partial}{\partial v_c} u_o^T v_c = u_o$$

Vector!  
Not high  
school  
single  
variable  
calculus

You can do things one variable at a time,  
and this may be helpful when things  
get gnarly.

$$\forall j \quad \frac{\partial}{\partial (v_c)_j} u_o^T v_c = \frac{\partial}{\partial (v_c)_j} \sum_{i=1}^d (u_o)_i (v_c)_i$$

$$= (u_o)_j$$

Each term is zero except when  $i=j$

7

-- $\rightarrow$  hint:

Apply the chain rule where  $u = e^x$

$$\frac{1}{\ln(10)} \cdot \frac{d}{dx} (\ln(e^x))$$

- Replace  $e^x$  with  $u$
- Apply the chain rule where  $u = e^x$
- Take the derivative of a logarithm expression
- Replace  $u$  with  $e^x$

$$\frac{1}{\ln(10)} \cdot \frac{1}{e^x} \cdot \frac{d}{dx} (e^x)$$

$$\begin{aligned}
 & \textcircled{2} \quad \frac{\partial}{\partial v_c} \log \underbrace{\sum_{w=1}^V \exp(u_w^T v_c)}_{\substack{f \\ z = g(v_c)}} \\
 &= \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \frac{\partial}{\partial v_c} \sum_{x=1}^V \exp(u_x^T v_c) \quad \text{Important to change index} \\
 & \quad \frac{\partial}{\partial v_c} f(\underbrace{z}_{g(v_c)}) = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial v_c} \quad \text{Use chain rule} \\
 &= \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \left( \sum_{x=1}^V \frac{\partial}{\partial v_c} \underbrace{\exp(u_x^T v_c)}_{\substack{f \\ z = g(v_c)}} \right) \quad \text{Move deriv inside sum} \\
 & \quad \left( \sum_{x=1}^V \exp(u_x^T v_c) \frac{\partial}{\partial v_c} u_x^T v_c \right) \quad \text{Chain rule} \\
 & \quad \left( \sum_{x=1}^V \exp(u_x^T v_c) u_x \right)
 \end{aligned}$$

$$\begin{aligned}
\frac{\partial}{\partial v_c} \log(p(o|c)) &= u_o - \frac{1}{\sum_{w=1}^V \exp(u_w^T v_c)} \cdot \left( \sum_{x=1}^V \exp(u_x^T v_c) u_x \right) \\
&= u_o - \sum_{x=1}^V \frac{\exp(u_x^T v_c)}{\sum_{w=1}^V \exp(u_w^T v_c)} u_x \quad \text{distribute term across sum} \\
&= u_o - \underbrace{\sum_{x=1}^V p(x|c)}_{\text{this is an expectation: average over all context vectors weighted by their probability}} u_x \\
&= \text{observed} - \text{expected}
\end{aligned}$$

This is just the derivatives for the center vector parameters  
 Also need derivatives for output vector parameters  
 (they're similar)  
 Then we have derivative w.r.t. all parameters and can minimize

29

به این مدل میگویند "Bag of words" چون به همه اعضای context میخواند prob  
 بالایی بده حالا مهم نی کجای context باشه

گرفتاری:

- **Problem:**  $J(\theta)$  is a function of **all** windows in the corpus (potentially billions!)
  - So  $\nabla_{\theta} J(\theta)$  is **very expensive to compute**
- You would wait a very long time before making a single update!
- **Very bad idea** for pretty much all neural nets!

حل:

## Mini Batch Gradient Descent

- **Solution: Stochastic gradient descent (SGD)**
  - Repeatedly sample windows, and update after each one
- Algorithm:

```
while True:
    window = sample_window(corpus)
    theta_grad = evaluate_gradient(J, window, theta)
    theta = theta - alpha * theta_grad
```

دوتا روش داریم

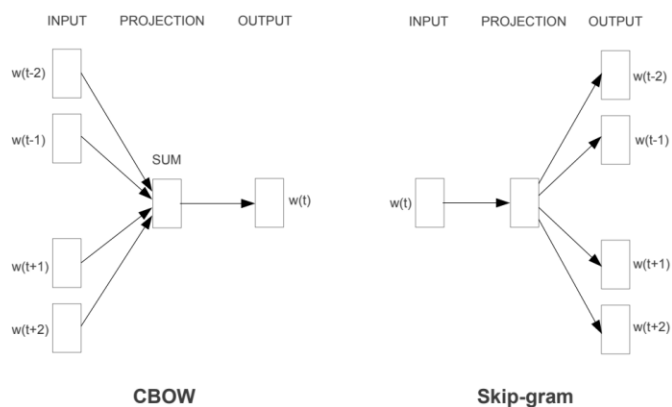
### 1. Skip-grams (SG)

Predict context ("outside") words (position independent) given center word

### 2. Continuous Bag of Words (CBOW)

Predict center word from (bag of) context words

معماری در مقاله:

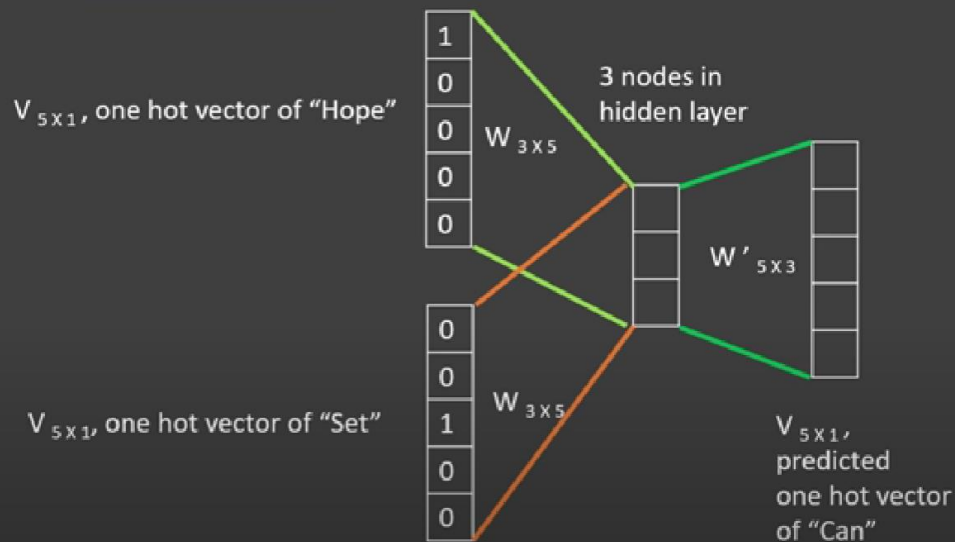


مثال :



## CBOW - Working

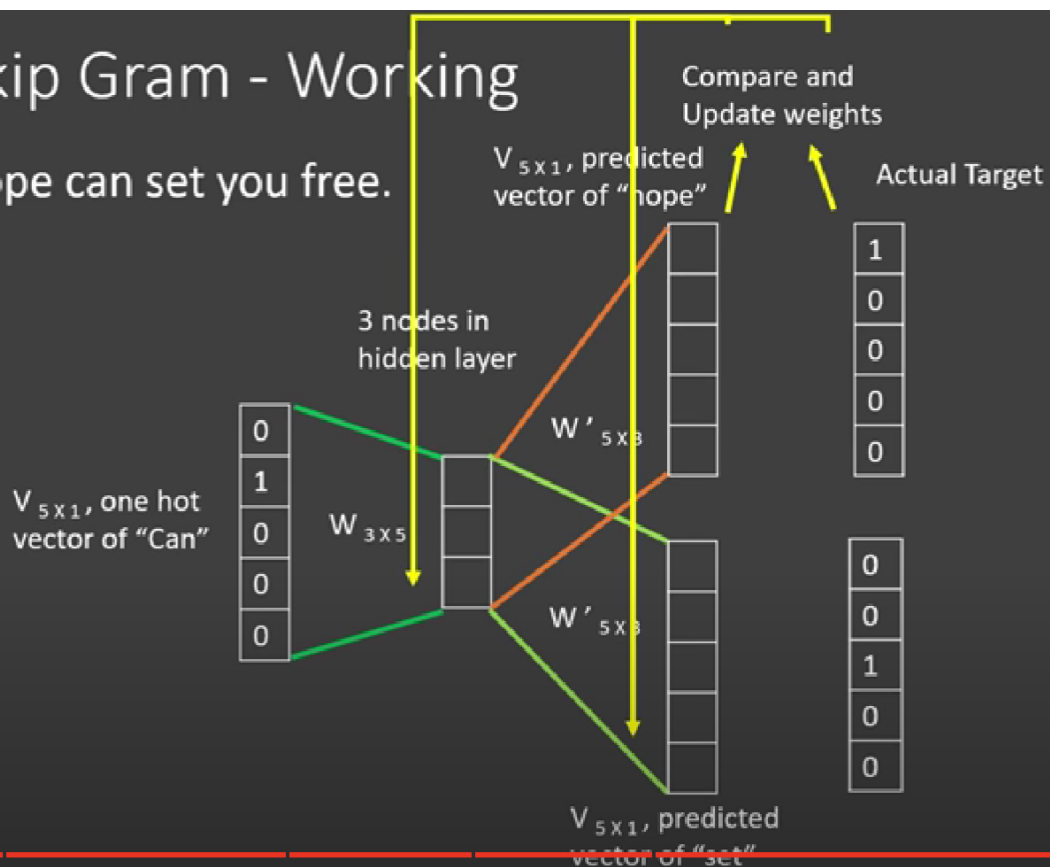
Hope can set you free.



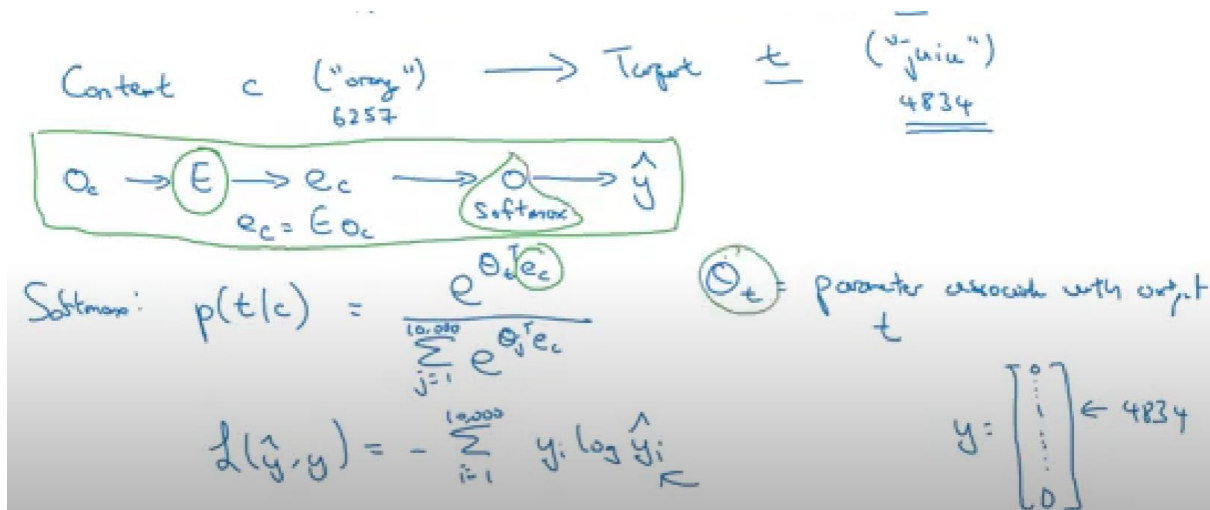
مثال :

## Skip Gram - Working

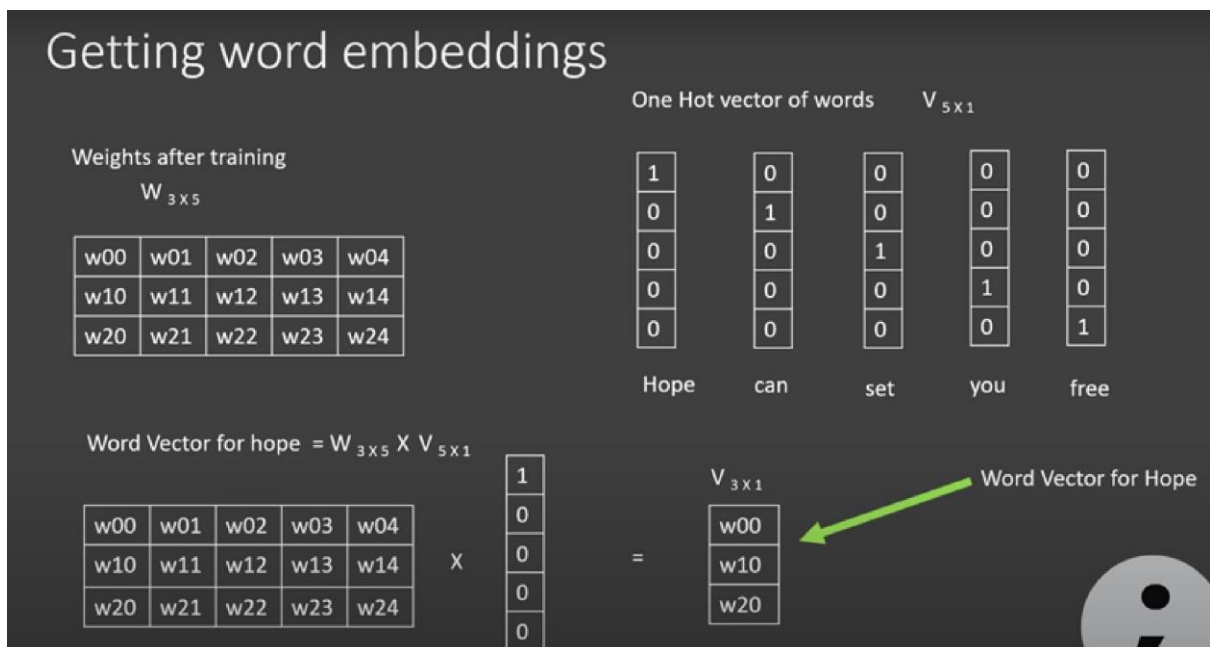
Hope can set you free.



البته یک چیز جالب دیگه هم فهمیدم ما اینجا میخوایم کلمه بعدی رو حدس بزنیم بنابراین میام میگی به ازای یک context word که داریم بیا و target word منو پیدا کن و اندرو میگف ما میایم context word رو تبدیل میکنیم به one hot میدیم به word embedding تبدیل میشه به e ما حالا اینو میدیم یک شبکه و میگی کلمه بعدی رو حدس بزن و اون میاد نزدیک ترین کلمه رو حدس میزنه



در واقع اینطوری e را پیدا میکنیم:



الان قضیه ای که وجود داره اینه در skip gram ما اخرش ی soft max داریم که خیلی هزینه بره

یک راه اینه مدل سلسله مراتبی **soft max** رو داشته باشیم --- < مثلا بگیم جواب در  
نصفه اوله بعد بگیم ایا در یک چهارو اوله و ....

ی حرف جالبی که زده اینه که مسئله ای که برای ما اهمیت داره اینه که ایا دوتا **word** باهم  
همسایه هستند یا نه ؟

یعنی بیایم مسئله چند کلاسمونو به مسئله دو کلاسه تبدیل کنیم

**Soft max ----- > sigmoid**

اینجاست که **negative sampling** میاد

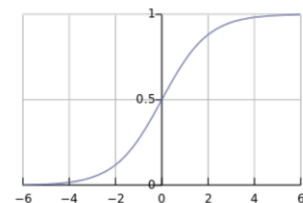
- Introduced in: "Distributed Representations of Words and Phrases and their Compositionality" (Mikolov et al. 2013)
- Overall objective function (they maximize):

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J_t(\theta)$$

$$J_t(\theta) = \log \sigma(u_o^T v_c) + \sum_{i=1}^k \mathbb{E}_{j \sim P(w)} [\log \sigma(-u_j^T v_c)]$$

sigmoid  
rather than softmax

- The logistic/sigmoid function:  $\sigma(x) = \frac{1}{1+e^{-x}}$   
(we'll become good friends soon)
- We maximize the probability of two words  
co-occurring in first log and minimize probability  
of noise words in second part



در واقع ما کزیمم میکنیم احتمال اون کلمه ای که کنارشه و مینیمم میکنیم احتمال **noise**  
ها رو

$$J_{neg-sample}(\mathbf{u}_o, \mathbf{v}_c, U) = -\log \sigma(\mathbf{u}_o^T \mathbf{v}_c) - \sum_{k \in \{K \text{ sampled indices}\}} \log \sigma(-\mathbf{u}_k^T \mathbf{v}_c)$$

فانشن **loss** میشه این ریختی

Sample with  $P(w)=U(w)^{3/4}/Z$ , the unigram distribution  $U(w)$  raised to the 3/4 power  
(We provide this function in the starter code).

The power makes less frequent words be sampled more often

احتمال هم اینطوری انتخاب میکنیم ----- < چون کلمات انگلیسی **of** و **the** اینا زیاد داره

لینک مقاله ها:

<https://arxiv.org/pdf/1301.3781.pdf>

[https://proceedings.neurips.cc/paper\\_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2013/file/9aa42b31882ec039965f3c4923ce901b-Paper.pdf)

ی قضیه دیگ داریم co-occurrence matrix

این ریختی:

### Example: Window based co-occurrence matrix

- Window length 1 (more common: 5–10)
- Symmetric (irrelevant whether left or right context)

- Example corpus:

- I like deep learning
- I like NLP
- I enjoy flying

counts	I	like	enjoy	deep	learning	NLP	flying	.
I	0	2	1	0	0	0	0	0
like	2	0	0	1	0	1	0	0
enjoy	1	0	0	0	0	0	1	0
deep	0	1	0	0	1	0	0	0
learning	0	0	0	1	0	0	0	1
NLP	0	1	0	0	0	0	0	1
flying	0	0	1	0	0	0	0	1
.	0	0	0	0	1	1	1	0

در واقع دو نوع رویکرد وجود داره یکی window بیس یکی روی کل doc

Building a co-occurrence matrix  $X$

- 2 options: windows vs. full document
- Window: Similar to word2vec, use window around each word → captures some syntactic and semantic information ("word space")
- Word-document co-occurrence matrix will give general topics (all sports terms will have similar entries) leading to "Latent Semantic Analysis" ("document space")

بدبختی: یک ماتریس گنده sparse و دارای اطلاعات کم و نویزی

راه حل: Low-dimensional vectors

فان فکتر

واریانس دامنه تغییرات هر متغیر نسبت ب خودش

کوواریانس دامنه تغییرات نسبت به بقیه متغیر ها

بردار ویژه هم برداریه ک میاد حاصل ضرب دو تا ماتریس رو خطی میکنه

به اون عدده هم که تو بردار ویژه ضرب میشه میگن مقدار ویژه

دترمینان هم معکوس یک ماتریسه

حالا ترانهاده چیه میاد جای سطر و ستون رو عوض میکنه

ایده pca اینه اون فیچر های correlated رو با ایده بردار ویژه اینا تبدیل به یک ویژگی کن

اما svd میاد ماتریسه ما رو تبدیل میکنه به حاصل ضرب ی ماتریس متفبقارن و قطری و متقارن

The diagram shows the SVD decomposition of a matrix  $X$ . On the left is a 3x5 matrix  $X$  with stars. This is equal to the product of three matrices: a 3x3 matrix  $U$  with stars, a 3x3 diagonal matrix  $\Sigma$  with dots on the diagonal, and a 5x5 matrix  $V^T$  with stars. Braces are used to label each matrix.

اعداد در ماتریس قطری با یک order هستند پس ما میتونیم اینطوری حذف کنیم

The diagram shows the truncated SVD decomposition. The matrix  $X$  is equal to the product of a truncated  $U$  matrix (3x3, with the third column highlighted in blue), a truncated  $\Sigma$  matrix (3x3, with the third row and column highlighted in blue), and a truncated  $V^T$  matrix (5x5, with the first three rows highlighted in blue). Braces are used to label each matrix.

اما اگر ما یک Co-occurrence vector داشته باشیم بیایم روش ی  $svd$  بزنینم اونقد خوب کار نمیکنه چون ما کلمات خیلی پرتکراری همچون *the* و *and* و ... داریم

راه حل COALS :

- Scaling the counts in the cells can help **a lot**
  - Problem: function words (*the, he, has*) are too frequent  $\rightarrow$  syntax has too much impact. Some fixes:
    - log the frequencies
    - $\min(X, t)$ , with  $t \approx 100$
    - Ignore the function words
- Ramped windows that count closer words more than further away words
- Use Pearson correlations instead of counts, then set negative values to 0
- Etc.

Glove:

The image shows handwritten mathematical definitions for GloVe word embeddings on a dark background. The first line defines  $X_{ij}$  as the number of times word  $j$  appears in the context of word  $i$ . The second line defines the conditional probability  $P(j|i)$  as the probability that  $j$  is in the context of  $i$ , which is equal to the ratio of  $X_{ij}$  to the sum of  $X_{ik}$  over all  $k$ .

$$X_{ij} = \# j \text{ appears in the context of } i$$
$$P(j|i) = P(j \text{ is in the context of } i)$$
$$= \frac{X_{ij}}{X_i \leftarrow \sum_k X_{ik}}$$

$$P(j|i) = P(j \text{ is in the context of } i)$$

$$= \frac{X_{ij}}{X_i \leftarrow \sum_k X_{ik}}$$

۴

	k = solid	k = gas	k = water	k = (random)
P(k ice)	high	low	high	low
P(k steam)	low	high	high	low
P(k ice)/P(k steam)	>1	<1	~ 1	~ 1

برای حل بزرگ بودن ماتریس:

Most general way ...

$$F(w_i, w_j, \tilde{w}_k) = \frac{P(k|i)}{P(k|j)} \quad \leftarrow \text{scalar}$$

$\nwarrow$  Some function       $\nwarrow$  word vectors       $\nwarrow$  word vectors

حالا چطوری F رو بدست آوریم:

Assuming Homomorphism

$$F(\underbrace{w_i^T \cdot \tilde{w}_k}_u - \underbrace{w_j^T \cdot \tilde{w}_k}_v) = \frac{F(w_i^T \cdot \tilde{w}_k)}{F(w_j^T \cdot \tilde{w}_k)} = \frac{P(k|i)}{P(k|j)}$$

So,

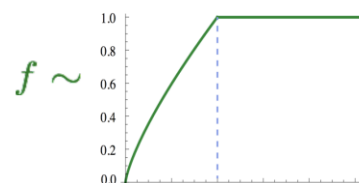
$$F(w_i^T \cdot \tilde{w}_k) = c P(k|i) \quad \text{we can ignore } c \text{ safely}$$

$F(x) = e^x$  is a solution!

پس loss function

Loss:  $J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$

- Fast training
- Scalable to huge corpora



$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

$$f(x) = \begin{cases} (x/x_{\max})^{100} & \text{if } x < x_{\max} \\ 1 & \text{otherwise} \end{cases}$$

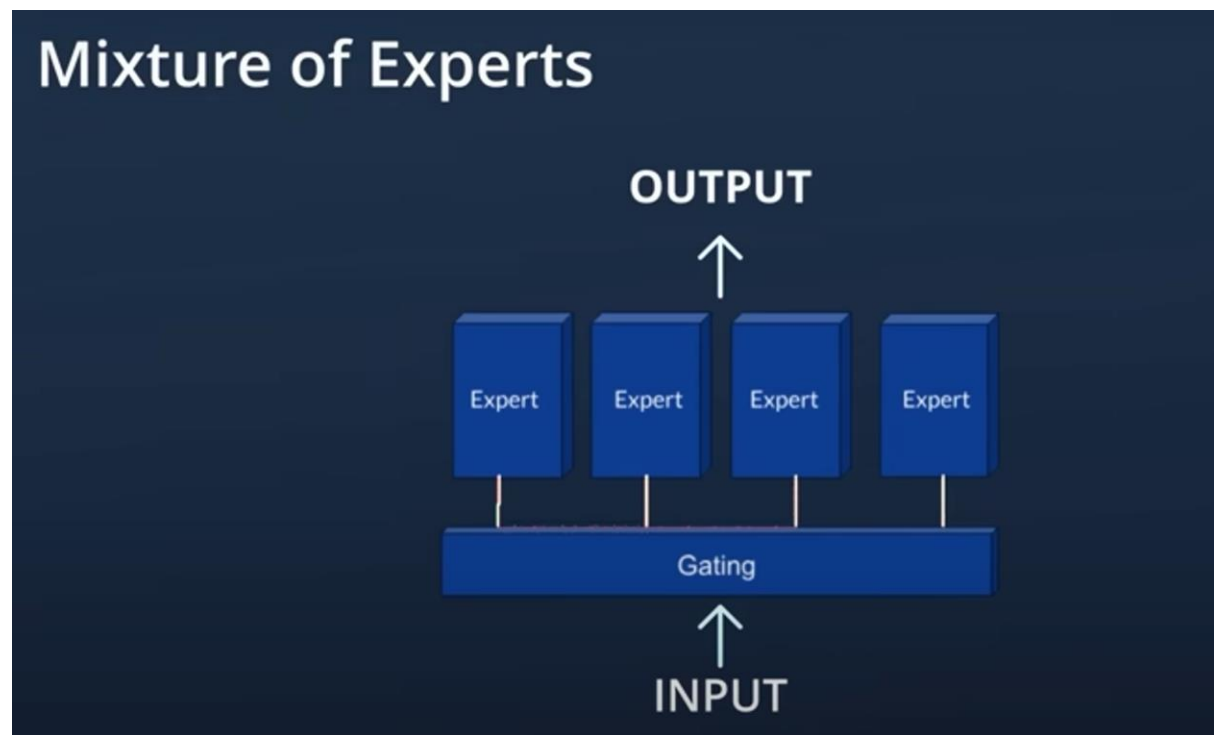
مدل‌های زبانی را می‌توان به دو صورت Intrinsic یا Extrinsic ارزیابی کرد.

Intrinsic: شامل مقایسه مدل Word Embeddings با یک مدل مرجع (مانند lexical database) می‌شود.

Extrinsic: برای ارزیابی مدل، نتایج آن را در یک تسک یا پروژه‌ای مانند text classification, machine translation و summarization استفاده کنیم

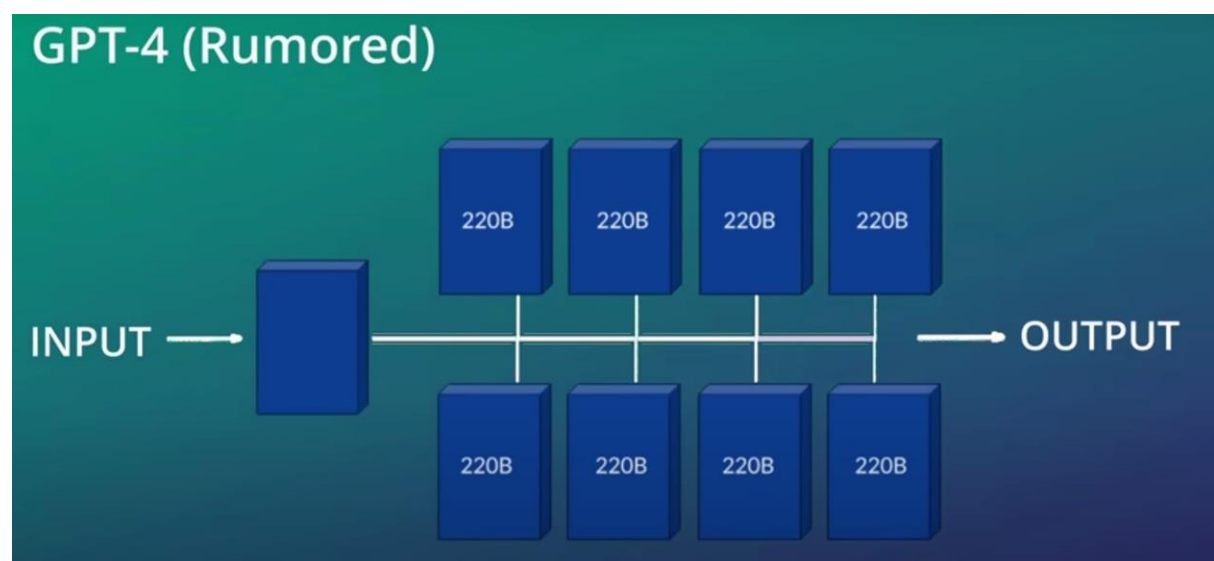


میسترال :



یک تعدادی expert داریم اینجا هشت تا هفت میلیاردی و هر کدوم متخصص یک چیزی هستند

و وقتی یک تسک میاد gating تصمیم میگیره بره سمت کدوم expert  
شایعه :



# Fine Tuning LLMs with Instruction

## *What is instruction Tuning?*

Instruction tuning represents a specialized form of fine-tuning in which a model is trained using pairs of input-output instructions

### **Prepare training data:**

- Use existing datasets or create new ones.
- Use prompt template libraries to format data as instructions.
- Divide the data into training, validation, and test sets.

### **Fine-tune the LLM:**

- Select prompts from the training set.
- Pass prompts to the LLM and generate completions.
- Compare completions to responses in training data.
- Calculate loss using the cross-entropy function.
- Update model weights using backpropagation.
- Repeat for multiple batches and epochs.

### **Evaluate LLM performance:**

- Use validation data set to calculate validation accuracy.
- Use test data set to calculate test accuracy.

## **catastrophic forgetting**

Catastrophic forgetting happens because the full fine-tuning process modifies the weights of the original LLM. While this leads to great performance on a single fine-tuning task, it can degrade performance on other tasks.

<https://medium.com/@veer15/the-hitchhikers-guide-to-instruction-tuning-large-language-models-d6441dbf1413#:~:text=Instruction%20tuning%20represents%20a%20specialized,Output%3A%20%E2%80%9CEnglish%2C%20French%E2%80%9D>

لینک باحال برای fin ai

<https://github.com/Sahaj777/Mastering-AI-in-Finance?tab=readme-ov-file#llms>