

L8 Word Vectors

+ Self-Supervised Learning

Feb 3, 2025

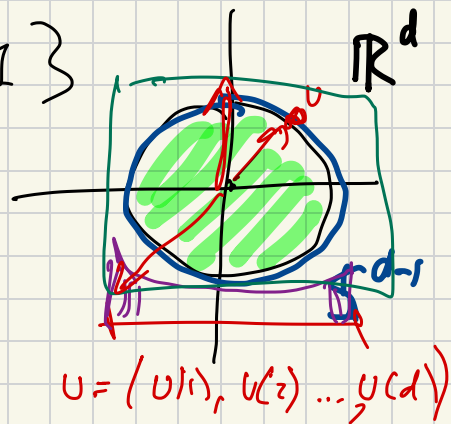


Data Mining:

Jeff M. Phillips

Generate set iid $u_1, u_2, \dots, u_d \sim \text{Unif}(\mathbb{S}^{d-1})$

$$\mathbb{S}^{d-1} = \{x \in \mathbb{R}^d \text{ s.t. } \|x\|=1\}$$



Step 1: generate $v \in [-1, 1]$

$$v \in \text{Unif}(0, 1)$$

$$= 0. b_1 b_2 b_3 \dots b_{\log_2 N}$$

$(0, 1)$

Generate angle $\theta \sim \text{Unif}(0, 2\pi)$

works $d=2$, hard for $d \geq 3$

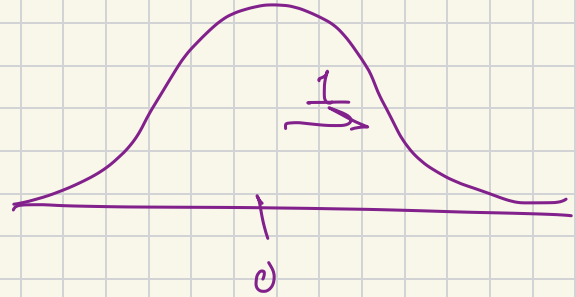
Generate $u_1, u_2, \dots, u_d \sim [\text{Unif}(0, 1)]^d \rightarrow \underline{\text{normalize}}$

Generate points g_1, g_2, \dots, g_d

$\sim N(0, 1)$
iid
normal

vector $g = (g_1, \dots, g_d)$

return $v = g / \|g\|$



Box-Muller Transform

$u_1, u_2 \sim \text{Unif}(0, 1)$

$$g_1 = \sqrt{-2 \ln(u_1)} \cos(2\pi u_2) \sim N(0, 1)$$
$$g_2 = \sqrt{-2 \ln(u_1)} \sin(2\pi u_2) \sim N(0, 1)$$

$$N(0, 1) = f(x) = \frac{1}{\sqrt{2\pi}} e^{-\|x\|^2/2}$$

Word Embeddings

A quick brown fox jumped over the
lazy dog.

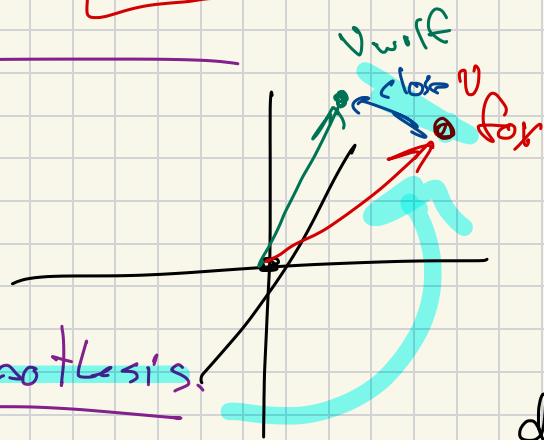
subject

verb

Newer map each word fox

premise: words are similar if
they appear in similar
context.

→ Distribution Hypothesis



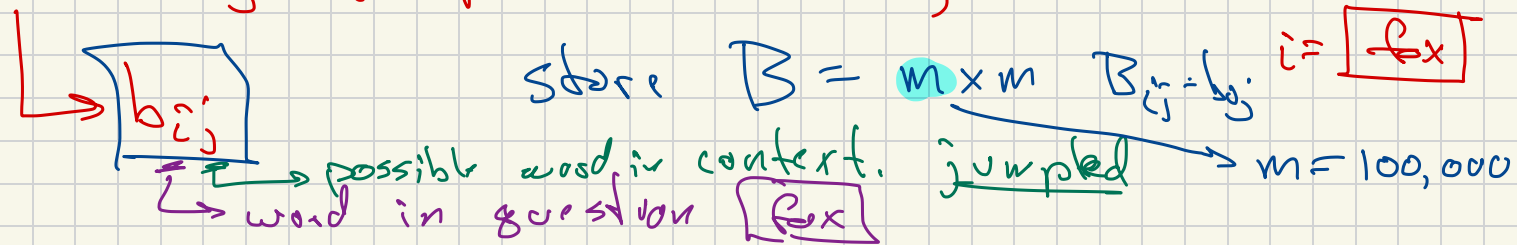
\mathbb{R}^d
 $d=300$

First Gen PPMI vectors. CBOW

A quick brown fox jumped over a (a) dog

- context window
 - k words on each side
 - k words before word

aggregate bag of words : # instances of
 $j = \text{jumped}$ in any context window of



Is frequency-context vector

$$b_i = (b_{i,1}, b_{i,2}, \dots, b_{i,m}) \in \mathbb{R}^m$$

a good representation?

$$p(i,j) = \frac{b_{i,j}}{N} \quad N = \# \text{ words in corpus}$$

$$p(i) = \frac{n_i}{N} \quad n_i = \# \text{ times word } i \text{ occurs}$$

$$v_{ij} = \max \left\{ \log \left(\frac{p(i,j)}{p(i) \cdot p(j)} \right), 0 \right\}$$

PPMI vec

$$V_i = (v_{i,1}, v_{i,2}, \dots, v_{i,m}) \in \mathbb{R}^m$$

↑ useful vectors rep word i .

Second Generation

want vectors

$$v_i \in \mathbb{R}^d$$

$$d = 300.$$

work like PPMI vectors

1. Get PPMI vectors \mathbb{R}^m

2. Dim Reduction $\mathbb{R}^m \rightarrow \mathbb{R}^d$.

Better idea: learn mapping to \mathbb{R}^d

using Self-supervised learning

→ Word2vec, GloVe

2013!

Self-Supervised Learning

Supervised Learning

$$x_i \in \mathcal{X} (= \mathbb{R}^d)$$

Input

$$X, y = \{(x_i, y_i)\}$$

$$y_i \in \{-1, +1\} \text{ classification}$$
$$y_i \in \mathbb{R} \text{ regression}$$

Label y_i .

Sol $X = \{x_1, \dots, x_n\}$
independent variables

↳ make prediction
 x_i on y_i

learns \rightarrow

$$f(x_i) \approx y_i$$

$$y = (y_1, y_2, \dots, y_n)$$

$$X = (x_1, x_2, \dots, x_n)$$
$$\in \mathbb{R}^{n \times d}$$

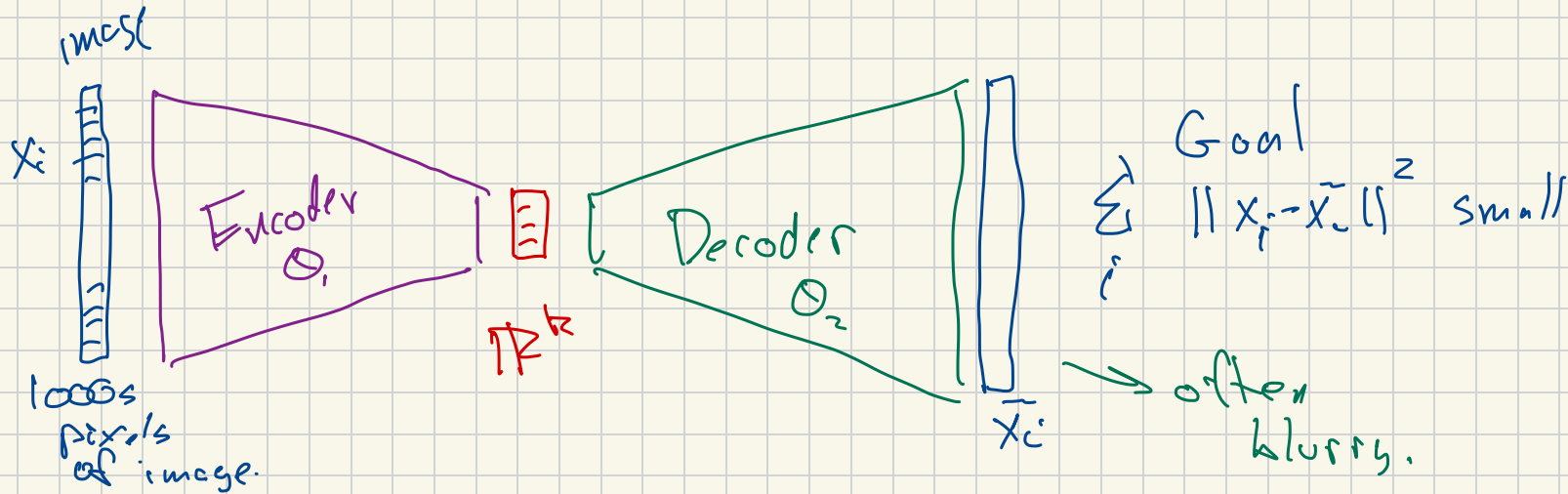
Problem

Labels hard to get.

Self-supervised learning

Train on data x_i to predict
itself, or part of itself.

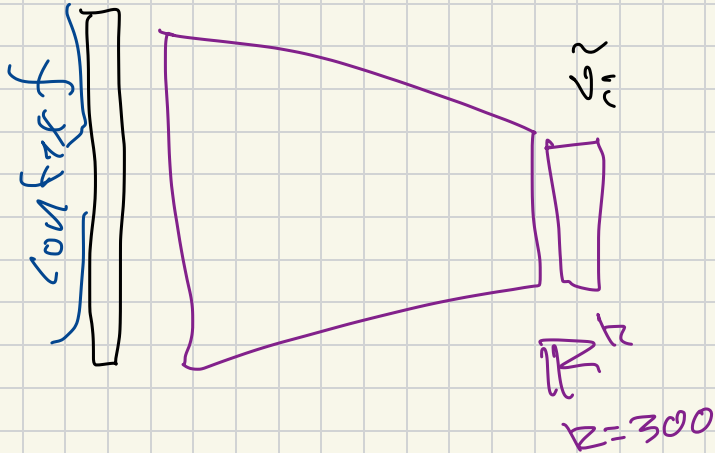
Ex. Auto-encoder



For text: masked model

A quick brown fox jumped over the lazy dog,

A - quick brown jumped over the ...

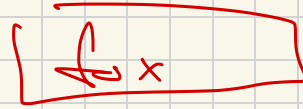


Data set

$$v_1, v_2, \dots, v_m \in \mathbb{R}^{300}$$

Guess $i \Rightarrow \mathcal{N}(\vec{v}_i)$

context window

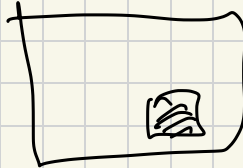


next word prediction.

↳ apples other data

graph embedding

images



spatial data

Third generation

ELMo
(Owens, AI)

BERT → RoBERTa
(Google)

① Contextual Embeddings

homonyms: apple → fruit
→ company.

ELMo: learn function $f: [\text{context} + \text{word}] \rightarrow \mathbb{R}^d$

② Transformer architecture, w/ attention
much larger context.
selected most influential words in context.