

Unsupervised Pretraining: We have a lot of text data, but mostly unlabeled - can still learn "language representation"

local, non-contextual
word embeddings
sentence embedding

global, contextual
pretrained LMs.

↳ "meaningful" approach (s.t. similar words are closer in rep) = word2vec (low-dim projections)
essentially, "meaning" of word - what words are close to it?

↳ then, can we predict word's neighbors from embedding value? $p(o|c) = \frac{\exp(u_o^T v_c)}{\sum_{w \in \text{vocab}} \exp(u_w^T v_c)}$
↑ "context word" ↑ "center word" ↑ log. res. vocabulary

$$\text{Train: } \arg \max_{u_1, \dots, u_n, v_1, \dots, v_n} \sum_{c, o} \log p(o|c)$$

all c, o combinations

Issue: huge vocabulary + really costly to compute all combs

↳ can reframe as binary classification problem. (right/wrong word?)

$p(o \text{ is the right word} | c) = \sigma(u_o^T v_c) \leftarrow$ sigmoid, now output is one of 2 classes.

↳ but we also need to push different words apart!

$p(w \text{ is the wrong word} | c) = \sigma(-u_w^T v_c)$

randomly chosen "negatives"

$$\arg \max_{u_1, \dots, u_n, v_1, \dots, v_n} \sum_{c, o} (\log p(o \text{ is right word} | c) + \sum_w \log p(w \text{ is wrong} | c))$$

$$\sum_{c, o} (\log \sigma(u_o^T v_c) + \sum_w \log \sigma(-u_w^T v_c))$$

word 2 vec in a nutshell

"push v_c toward u_o and away from other u_w "

→ can then determine, in theory, real algebraic relations!

→ this idealized word rep. \Rightarrow basic one-hot encoding scheme

Issue w/ word embeddings: vectors do not change given context. "local context"

↳ use pretrained LMs hidden state after training on sentence

"Sesame Street" models: ELMO: bidirectional LSTM model

BERT: huge transformer LM

→ context-dependent embeddings

ELMO: trains forward + backward LM on same data. "Add" all hidden states.

↳ "Embeddings from Language Models"

large corpus of unlabeled text. IN PAPER

simple: $ELMO_t = [h_{t, \text{last}}^{\text{fwd}}, h_{t, \text{last}}^{\text{bwd}}]$
complex: $ELMO_t = \gamma \sum_{i=1}^L w_i [h_{t, i}^{\text{fwd}}, h_{t, i}^{\text{bwd}}]$
↑ learned

→ for a downstream app: feed existing rep + ELMO rep to model (after running ELMO).

really good at generation (GPT et al.)

pros: context-specific, semantically meaningful (very)
cons: computationally expensive, storage issue.

BERT (industry-leading) → "Bidirectional Encoder Representations from Transformers"

↳ basically, ELMO that is purely attention-based.

Base model for BERT: decoder-only transformer (specified direction)

"We can do better": full self-attention layers (no mask), no input-output shift.

but randomly mask out some tokens → so some words are 1:1, but keeps model actively "on its toes" but model cannot directly overhear.

fill in the blanks.

TLDR: encoder-only transformer w/ input mask,

Training BERT: pass in pairs of sentences (Q/A, machine translation, etc.).

↳ apply 15% mask

↳ coin flip A/B or B/A sentence order

→ very first token: trained on binary classification (does A follow B or B follow A)

Using BERT (unlike FlanT5): For classification, 1. replace A/A classifier w/ desired cross-entropy loss, 2. finetune end-to-end

→ can also be used to extract representations (like ELMO).
for per-portion, -SQUAD: highlight span of paragraph
-NER: output class for each token

↳ best config: concat last k hidden state layers.

change loss for each token