

Language Embeddings

Plan

Review

Autoencoders

Contrastive Learning

PCA

Logistics

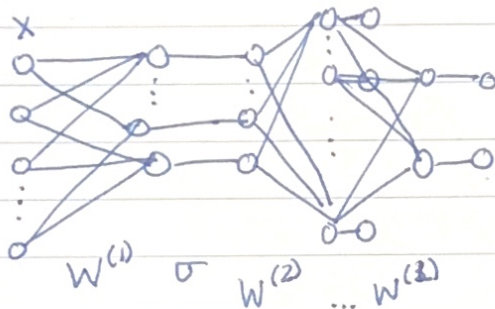
Games tonight!

Check in form!

Scribed notes!

Zoom!

Review



1. Model $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$

2. Loss $\mathcal{L}(w)$

3. Optimizer

$$w^{(t+1)} = w^{(t)} - \alpha \nabla \mathcal{L}(w^{(t)})$$

momentum $v^{(t+1)} = (1-\beta)v^{(t)} + \beta \nabla_w \mathcal{L}(w^{(t)})$
 adaptivity $s^{(t+1)} = (1-\beta)s^{(t)} + \beta [\nabla_w^2 \mathcal{L}(w^{(t)})]^{-1}$

Backprop:

Forward:

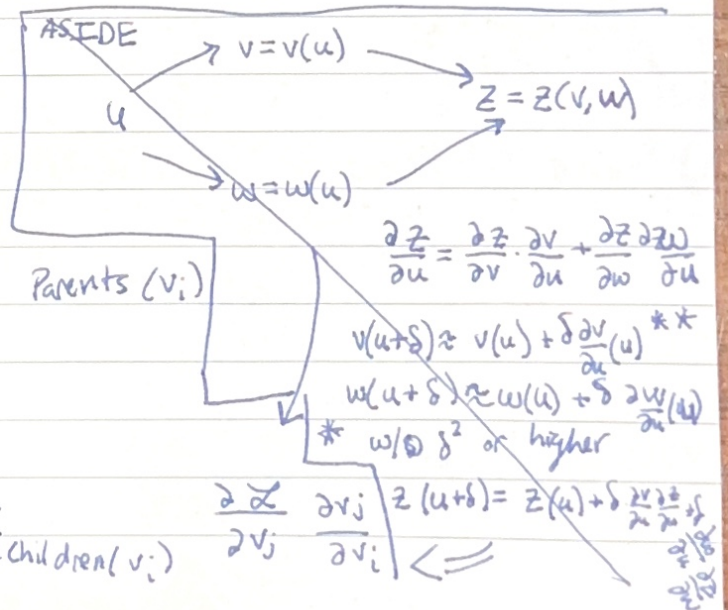
for i in $\{1, \dots, N\}$:

compute v_i from Parents (v_i)

Backward:

for i in $\{N, \dots, 1\}$:

compute $\frac{\partial \mathcal{L}}{\partial v_i} = \sum_{j \in \text{children}(v_i)}$



$$\frac{\partial \mathcal{L}}{\partial v_i} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial v_i} + \frac{\partial \mathcal{L}}{\partial w} \frac{\partial w}{\partial v_i}$$

$$\frac{\partial \mathcal{L}}{\partial v_i} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial v_i} + \frac{\partial \mathcal{L}}{\partial w} \frac{\partial w}{\partial v_i}$$

$$\frac{\partial \mathcal{L}}{\partial v_i} = \frac{\partial \mathcal{L}}{\partial z} \frac{\partial z}{\partial v_i} + \frac{\partial \mathcal{L}}{\partial w} \frac{\partial w}{\partial v_i}$$

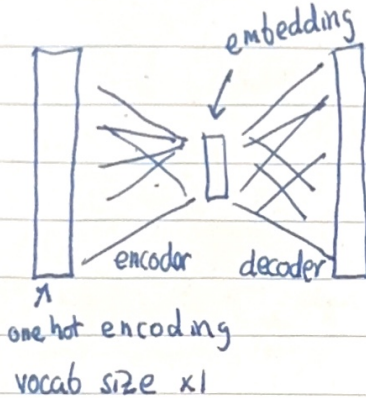
Language Embeddings

Our models take vectors... how can we convert words to vectors?

Approach #1: One-hot encodings

- ⊖ not meaningful
- ⊖ large

Approach #2: Autoencoders



Questions for Activity:

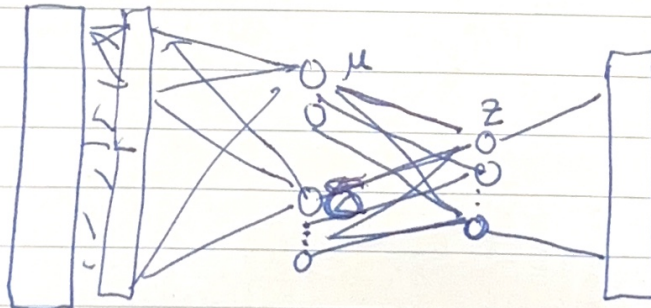
- What loss?
- How to make similar words close?

Why train decoder?

Approach #3: Variational Autoencoders

What if we want embeddings to be nice distributed in latent space?

$$z \sim \mathcal{N}(0, I) \quad z \in \mathbb{R}^n \quad z_i = \mu_i + \sigma_i \epsilon_i \quad \text{for } \mathcal{N}(0, 1)$$



$$\mathcal{L}(w) = \|f(x) - x\|_2^2 + \lambda \text{KL}(\mathcal{N}(0, I), \mathcal{N}(\mu, \sigma^2))$$

↑

distance between distributions

$$\text{KL}(P||Q) = H(P, Q) - H(P)$$

Contrastive Learning

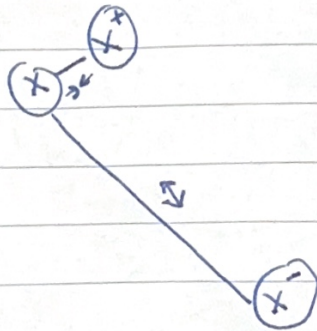
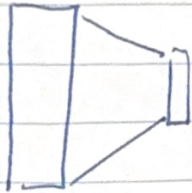
We're working on unsupervised task i.e., no labels

positive (word, next-word) hopefully close
 $\rightarrow (x, x^+)$ $f(x)^T f(x^+)$ large

negative (word, unrelated-word) probably far
 $\rightarrow (x, x^-)$ $|f(x)^T f(x^-)|$ small

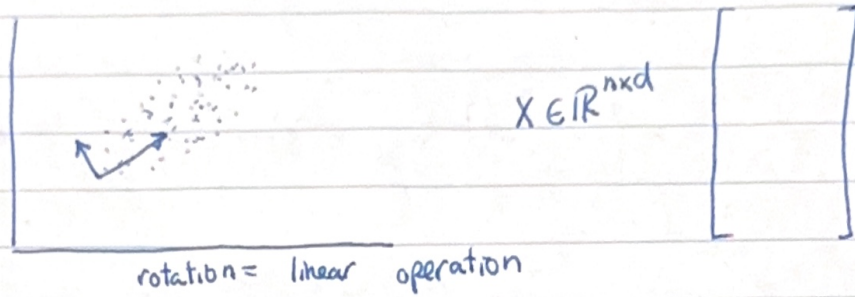
$$\mathcal{L}(w) = -\sum_{x, x^+} f(x)^T f(x^+) + \sum_{x, x^-} [f(x)^T f(x^-)]^2$$

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^r$$



Principal Component Analysis

Motivation: visualize points in \mathbb{R}^d meaningfully



Find $v \in \mathbb{R}^d$: Xv is meaningful i.e., captures variation of data!
 $n \times d \quad d \times 1$

$$\Leftrightarrow \|Xv\|_2^2 \text{ is large} \Leftrightarrow v^T X^T X v \text{ is large}$$

$$\max_{v: \|v\|_2^2=1} v^T X^T X v \Leftrightarrow \text{largest eigenvalue of } X^T X$$

$$\max_{v: \|v\|_2^2=1, v \perp v^{(1)}} v^T X^T X v \Leftrightarrow \text{2nd largest eigenvalue}$$

$$X^T X = \sum_{i=1}^r \lambda_i v^{(i)} v^{(i)T} \quad \text{where } \|v^{(i)}\|_2^2=1 \text{ and } v^{(i)} \cdot v^{(j)}=0$$

$d \times d \quad 1 \times d$

maximize eigenvalue corresponding to largest λ_i

$$\text{sanity check: } X^T X v^{(j)} = \sum_{i=1}^r \lambda_i v^{(i)} v^{(i)T} v^{(j)} = \lambda_j v^{(j)} v^{(j)T} v^{(j)} = \lambda_j v^{(j)}$$