# Language Embeddings

Logistics

Games tonight!

Check in form!

Scribed notes!

Zoom!

## Review



$W^{(1)} \quad \sigma \quad W^{(2)} \quad \dots \quad W^{(3)}$

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

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

3. Optimizer

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

momentum $\quad v^{(t+1)} = (1-\beta) v^{(t)} + \beta \nabla_w \mathcal{L}(w^{(t)})$

adaptivity $\quad s^{(t+1)} = (1-\beta) g^{(t)} + \beta \left[\nabla_w \mathcal{L}(w^{(t)})\right]^2$

**Backprop:**

ASIDE

$v = v(u)$

$u$

$z = z(v, w)$

$w = w(u)$

$$\frac{\partial z}{\partial u} = \frac{\partial z}{\partial v} \cdot \frac{\partial v}{\partial u} + \frac{\partial z}{\partial w} \cdot \frac{\partial w}{\partial u}$$

**Forward:**

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

compute $v_i$ from Parents($v_i$)

$v(u+\delta) \approx v(u) + \delta \frac{\partial v}{\partial u}(u)$ **

$w(u+\delta) \approx w(u) + \delta \frac{\partial w}{\partial u}(u)$

\* w/o $\delta^2$ or higher

**Backward:**

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

compute $\dfrac{\partial \mathcal{L}}{\partial v_i} = \displaystyle\sum_{j \in \text{Children}(v_i)} \frac{\partial \mathcal{L}}{\partial v_j} \frac{\partial v_j}{\partial v_i}$

$z(u+\delta) = z(u) + \delta \frac{\partial v}{\partial u}\frac{\partial z}{\partial v} + \delta \dots$

[Left margin sideways derivation:]

$\frac{\partial z}{\partial u} = \frac{\partial z}{\partial v}\frac{\partial v}{\partial u} + \frac{\partial z}{\partial w}\frac{\partial w}{\partial u}$

$= \left[\frac{\partial z}{\partial v}, \frac{\partial z}{\partial w}\right] \begin{bmatrix} \frac{\partial v}{\partial u} \\ \frac{\partial w}{\partial u} \end{bmatrix}$

$z(u+\delta) \approx z(v(u), w(u)) + \frac{\partial z}{\partial v} \delta \frac{\partial v}{\partial u} + \frac{\partial z}{\partial w} \delta \frac{\partial w}{\partial u}$

$\approx z(v(u), w(u)) + \delta \left[s_v + s_w\right]$

$\approx z(v+\delta s_v, w+\delta s_w)$
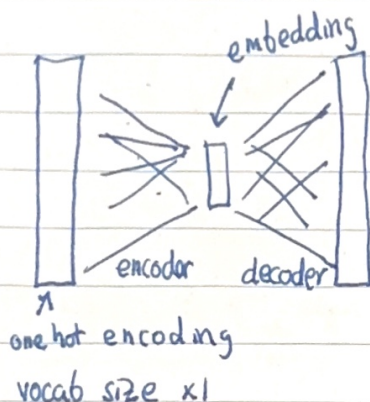
# Language Embeddings

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

Approach #1: One-hot encodings
- ⊖ not meaningful
- ⊖ large

Approach #2: Autoencoders

embedding



encoder  decoder

↑
one hot encoding
vocab size ×1

Why train decoder?
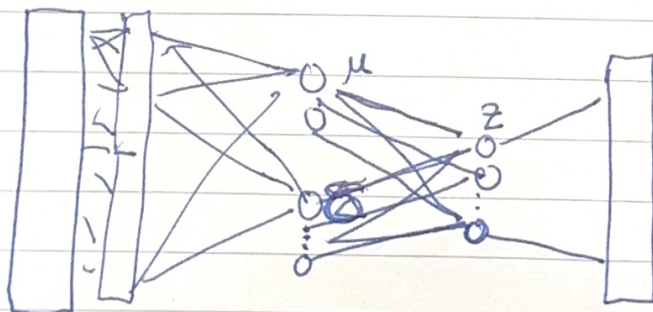
Questions for Activity:
- What loss?
- How to make similar words close?

Approach #3: Variational Autoencoders

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

$$Z \sim N(0, I) \quad Z \in \mathbb{R}^r \quad Z_i = \mu_i + \sigma_i \epsilon_i \quad \text{for } N(0,1)$$

x



$\mu$

$Z$

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

↑
distance between distributions
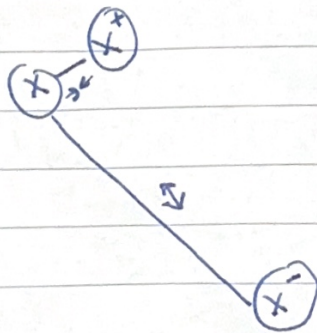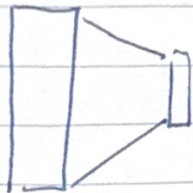
$$KL(p \| q) = H(p, q) - H(p)$$

## Contrastive Learning

We're working on unsupervised task ie., no labels

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

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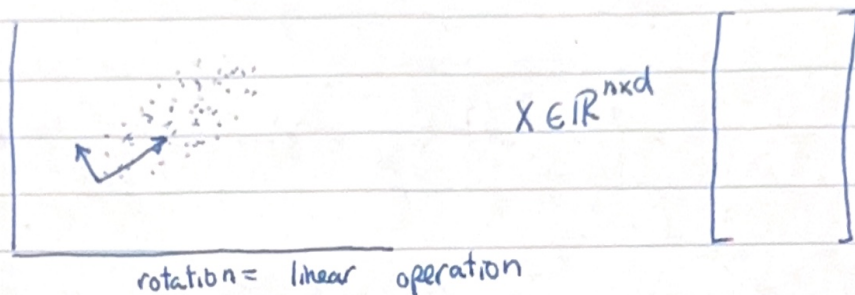$\hookrightarrow (x, x^-)$          $|f(x)^T f(x^-)|$ small

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

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

# Principal Component Analysis

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



$X \in \mathbb{R}^{n \times d}$

rotation = linear operation

Find $v \in \mathbb{R}^d$: $\underset{n \times d \; d \times 1}{X v}$ is meaningful ie, captures variation of data!

$\iff \| X v \|_2^2$ is large $\iff v^T X^T X v$ is large

$\underset{v: \|v\|_2^2 = 1}{\max} \quad v^T X^T X v \iff$ largest eigenvalue of $X^T X$

$\underset{v: \|v\|_2^2 = 1, v \perp v^{(1)}}{\max} \quad v^T X^T X v \iff$ 2nd largest eigenvalue

$X^T X = \overset{r \; \swarrow \text{rank}}{\underset{i=1}{\sum}} \lambda_i \underset{d \times 1}{v^{(i)}} \underset{1 \times d}{v^{(i)T}}$ where $\| v^{(i)} \|_2^2 = 1$ and $v^{(i)} \cdot v^{(j)} = 0$

maximize w/eigenvalue corresponding to largest $\lambda_i$

sanity check: $X^T X v^{(j)} = \overset{r}{\underset{i=1}{\sum}} \lambda_i v^{(i)} v^{(i)T} \; v^{(j)} = \lambda_j v^{(j)} v^{(j)} \cdot v^{(j)} = \lambda_j v_j$