

Matrix factorization

distributional hypotheses?

\underline{w} \therefore which words does w appear together with within some distance?

• $\phi(w)$: a feature vector of w

$\phi(w)_i$: the count of the word i appearing in the context of w .
the co-occurrence vector of w .

• A co-occurrence matrix A

$$A = \begin{bmatrix} \phi(w_0) \\ \phi(w_1) \\ \vdots \\ \phi(w_{|V|}) \end{bmatrix}^T \in \mathbb{R}^{|V| \times |V|}$$

• compute/subtract the mean vector.

$$\mu = \frac{1}{|V|} \sum_{i=1}^{|V|} \phi(w_i)$$

$$\tilde{A} = \begin{bmatrix} \phi(w_0) - \mu \\ \phi(w_1) - \mu \\ \vdots \\ \phi(w_{|V|}) - \mu \end{bmatrix}^T : \text{centered co-occurrence matrix}$$

$$\begin{aligned} U &\in \mathbb{R}^{|V| \times d} \\ V &\in \mathbb{R}^{|V| \times d}, \quad d \ll |V| \end{aligned}$$

$$\text{PCA: } \min_{U, V} \| \tilde{A} - UV^T \|_F^2$$

$$= \sum_{i=1}^{|V|} \sum_{j=1}^{|V|} \left(\tilde{\phi}(w_i)_j - u_i^T v_j \right)^2$$

word embedding \swarrow
context embedding \nwarrow

extension to n-grams

~~$$A = \begin{bmatrix} | & & | \\ \text{ } & & \text{ } \\ | & & | \end{bmatrix}$$~~

$$A = \begin{bmatrix} | & & | \\ \text{ } & & \text{ } \\ | & & | \end{bmatrix}$$

$$\tilde{A} = \tilde{A} - \mu \mathbb{1}$$

$$\|\tilde{A} - UV^T\|_F^2$$

$$U \in \mathbb{R}^{|\mathcal{V}| \times d}$$

$$V \in \mathbb{R}^{|\mathcal{V}|^n \times d}$$

← impractical to estimate!

BoW

$$V_{i\bar{j}} = V'_{i\bar{j}_1} + V'_{i\bar{j}_2} + \dots + V'_{i\bar{j}_n}, \text{ where } \bar{j} = (\bar{j}_1, \bar{j}_2, \dots, \bar{j}_n)$$

$$V' \in \mathbb{R}^{|\mathcal{V}| \times d}$$

$$\min_{U, V} \|\tilde{A} - UV^T\|_F^2$$

$$\sum_{i=1}^{|\mathcal{V}|} \sum_{\bar{j}=1}^{|\mathcal{V}|^n} \left(\underset{\tilde{\mu}_{\bar{j}}}{\phi(w_i)_{\bar{j}}} - u_i^T \left(\sum_{n'=1}^n V'_{i\bar{j}_{n'}} \right) \right)^2$$

beyond BoW

n-gram

$$\tilde{\phi}(w_i)_{\bar{j}} = \phi(w_i)_{\bar{j}-n}$$

$$\sum_{i=1}^{|\mathcal{V}|} \sum_{\bar{j}=1}^{M^n} \left(\tilde{\phi}(w_i)_{\bar{j}} - u_i^T F_{\theta}(\bar{j}) \right)^2$$

$$= \sum_{i=1}^{|\mathcal{V}|} \sum_{\bar{j}=1}^{M^n} \left(\phi(w_i)_{\bar{j}} - (u_i^T F_{\theta}(\bar{j}) + \mu_{\bar{j}}) \right)^2$$

A generative story

$v_D \in \mathbb{R}^d$: document embedding

$u_w \in \mathbb{R}^d$: word embedding of w

for the i th location in the document, which word should it write?

$$p(w_i = w | D) = \frac{\exp\{v_D^T u_w\}}{\sum_{w' \in V} \exp\{v_D^T u_{w'}\}}$$

vocabulary

$$\frac{1}{N_D} \sum_{n=1}^{N_D} \sum_{t=1}^{T_n} \log p(w_t | D)$$

$v_D \in \Delta^{d+1} \sim \text{Dirichlet}(\rho)$: mixture coefficients

$u_w \in \mathbb{R}^k$: word embedding $|V|$

$u_d \in \mathbb{R}^k$: topic embedding d topics.

$$p(w_1 \dots w_T, v_D, u_d) = \prod_{t=1}^T \sum_{d'=1}^d v_D^{d'} p(w_t | u_{d'})$$

$$p(w_t | u_{d'}) = \frac{\exp\{u_{w_t}^T u_{d'}\}}{\sum_{w' \in V} \exp\{u_{w'}^T u_{d'}\}}$$

(c) autoregressive model

• Probabilistic modeling

• observations X

• hidden Z

parameters Θ

$$\left\{ \begin{array}{l} p(X|z, \theta) \\ p(z) \\ p(\theta) \\ \dots \end{array} \right\}$$

\Rightarrow

long

$$\frac{p(\theta | D_{\text{train}} = \{x_1, \dots, x_N\})}{p(z_1, \dots, z_N | D_{\text{train}})}$$

$$\int p(X_{\text{test}} | \theta, z) p(\theta, z | D_{\text{train}}) d\theta, z = p(X_{\text{test}} | D_{\text{train}})$$

\uparrow
predictive distrib.