

# **Some Remarks on Natural Language Processing**

Nils Kujath

v2026.03

# Degenerate Geometry of One-Hot Word Representations

- Let  $V = \{w_1, \dots, w_{|V|}\}$  be a finite vocabulary. The canonical (standard) basis of  $\mathbb{R}^{|V|}$  is the indexed set  $\mathcal{B}_{\text{can}} = \{\mathbf{e}_1, \dots, \mathbf{e}_{|V|}\}$ , where  $\mathbf{e}_i$  denotes the  $i$ -th basis vector, i.e. the unique (one-hot) vector with a 1 in coordinate  $i$  and zeros elsewhere. The BoW construction  $\boxed{\iota : V \rightarrow \mathbb{R}^{|V|}, \quad w_i \mapsto \mathbf{e}_i}$  identifies each word  $w_i \in V$  with the canonical basis vector  $\mathbf{e}_i \in \mathbb{R}^{|V|}$ .
- Under the standard inner product on  $\mathbb{R}^{|V|}$ , all distinct word types are mutually orthogonal:

$$\mathbf{e}_u^\top \mathbf{e}_v = \delta_{uv} = \begin{cases} 1 & \text{if } u = v, \\ 0 & \text{if } u \neq v, \end{cases} \quad \forall u, v \in \{1, \dots, |V|\}.$$

Fmr., the metric restricted to  $\iota(V)$  is trivial; all word types are equidistant in  $\mathbb{R}^{|V|}$ :  $\forall u, v \in \{1, \dots, |V|\}$ ,

$$d|_{\iota(V) \times \iota(V)}(\mathbf{e}_u, \mathbf{e}_v) = \|\mathbf{e}_u - \mathbf{e}_v\| = \begin{cases} 0 & \text{if } u = v, \\ \sqrt{\underbrace{(1-0)^2 + (0-1)^2}_{\text{pos. } u} + \underbrace{(0-0)^2 + \dots + (0-0)^2}_{\text{pos. } v} + \dots + \underbrace{(0-0)^2}_{|V|-2 \text{ pos.}}} = \sqrt{2} & \text{if } u \neq v \end{cases}$$

- Consequently, the geometry induced by  $\iota$  is degenerate: all distinct words are equally dissimilar, and no notion of graded semantic proximity can be expressed under the BoW representation.

# From the Distributional Hypothesis to Word Embeddings

- Let  $V = \{w_1, \dots, w_{|V|}\}$  be a finite vocabulary. Let  $\mathcal{D} = (t_1, t_2, \dots, t_N)$  be a corpus of  $N$  tokens from  $V$ . Fix a context window size  $k \in \mathbb{N}^+$ . Define the context map  $\mathcal{C}_k : \{n \in \mathbb{N} : k < n \leq N - k\} \rightarrow V^{2k}$  by:

$$\mathcal{C}_k(n) = (t_{n-k}, \dots, t_{n-1}, t_{n+1}, \dots, t_{n+k}).$$

Note that positions  $n \leq k$  and  $n > N - k$  are excluded since  $k$  tokens of context are required on each side.

- For each  $w_i \in V$ , define the distributional profile of  $w_i$  in  $\mathcal{D}$  as the multiset:

$$\Delta_k(w_i) = \{\{\mathcal{C}_k(n) : n \in \{k+1, \dots, N-k\}, t_n = w_i\}\}.$$

- The Distributional Hypothesis (see esp. Harris 1954 and Firth 1957) asserts that  $w_i$  and  $w_j$  are semantically similar iff they appear in similar contexts, that is, iff  $\Delta_k(w_i) \approx \Delta_k(w_j)$ . The previous slide has shown that the degenerate geometry of BoW representations precludes any graded notion of similarity. However, comparing multisets over  $V^{2k}$  directly also seems intractable. The goal is therefore to find a map:

$$\phi : V \rightarrow \mathbb{R}^m (m \ll |V|) \quad \text{s.t.} \quad \Delta_k(w_i) \approx \Delta_k(w_j) \quad \text{is operationalised as} \quad \phi(w_i) \approx \phi(w_j) \text{ in } \mathbb{R}^m.$$

That is,  $\phi$  embeds the discrete set  $V$  into (the so-called embedding space)  $\mathbb{R}^m$  such that distributional similarity in  $\mathcal{D}$  is faithfully compressed into geometric proximity.

## References

- Firth, John R. 1957. A synopsis of linguistic theory, 1930–1955. In *Studies in Linguistic Analysis*, 1–32. Oxford: Basil Blackwell.
- Harris, Zellig S. 1954. Distributional structure. *Word* 10(2–3). 146–162.