

# Edit Distance

---

## 1. The Edit Distance Problem

### Definition of a ref-word

Given a finite set of variables  $V \subseteq \text{SVars}$  we define the alphabet of ref words as:  $\Gamma_V := \{x \vdash, \neg x\}$ . And given an alphabet  $\Sigma$  such that  $\Sigma \cap \Gamma_V = \emptyset$  we can define the set of ref words over  $\Sigma$  and  $V$  as:  $\mathbf{r} \in (\Sigma \cup \Gamma_V)^*$ . Next, a ref-word is valid if and only if, every occurrence of a variable in the ref-word is opened exactly once and closed afterwards, exactly once.

### Functions on ref-words

We can define the projection of a ref word over a set  $S$ ,  $r \uparrow S$ , recursively as:<sup>1</sup>

1.  $r \in S \rightarrow r \uparrow S = r$
2.  $r \notin S \rightarrow r \uparrow S = \epsilon$
3.  $(r_1 \cdot r_2) \uparrow S = (r_1 \uparrow S) \cdot (r_2 \uparrow S)$

$\text{Vars}(\mathbf{r})$  is the set of variables  $x \in V$  that occurs in the ref-word:

$$\text{Vars}(r) := \{x \in V \mid \exists r_x^{pre}, r_x, r_x^{post} \in (\Sigma \cup \Gamma_V)^* \text{ such that } r = r_x^{pre} \cdot x \vdash \cdot r_x \cdot \neg x \cdot r_x^{post}\} \quad (1)$$

$\text{tup}(r)$  are the positions each ref-word is referencing, and is defined as:

$$\text{tup}(r) := \{x \mapsto [i_x, j_x] \mid x \in \text{Vars}(r), i_x = |r_x^{pre} \uparrow \Sigma|, j_x = i_x + |r_x \uparrow \Sigma|\} \quad (2)$$

**Postulate:**

$$\text{valid}(r) \rightarrow |\text{tup}(r)| = |\text{Vars}(r)| \quad (3)$$

### Definition of ref-word tuple

### Distance between two ref words

Next, given two ref words  $r_1, r_2 \in (\Sigma \cup \Gamma_V)^*$  and a distance function  $\mathbf{d} : \Sigma^* \times \Sigma^* \rightarrow \mathbb{R}$  the distance  $\mathbf{d}_{\mathbf{r}}$  between  $r_1$  and  $r_2$  is defined as:  $\mathbf{d}_{\mathbf{r}}(r_1, r_2) = \mathbf{d}(r_1 \uparrow \Sigma, r_2 \uparrow \Sigma)$

### Ref-word distance languages

Given a ref-word language reference  $L_{ref} \subseteq (\Sigma \cup \Gamma_V)^*$  and a distance  $k \in \mathbb{R}$ , the k-distance ref-word language reference is defined as

$$\llbracket L_{ref} \rrbracket_{ref}^k = \{r \in (\Sigma \cup \Gamma_V)^* \mid \text{valid}(r), \exists r' (r' \in L_{ref} \wedge \mathbf{d}_{\mathbf{r}}(r, r') \leq k)\} \quad (4)$$

Given a document  $d$ , the k-distance ref-word language is:

$$\llbracket L_{ref} \rrbracket_d^k = \{\text{tup}(r) \mid r \uparrow \Sigma = d, \exists r' \in L_{ref} (\mathbf{d}_{\mathbf{r}}(r, r') \leq k)\} \quad (5)$$

### Postulates

$$\llbracket \llbracket A \rrbracket_{ref}^{k_1} \rrbracket_{ref}^{k_2} = \llbracket A \rrbracket_{ref}^{k_1+k_2} \quad (6)$$

$$\llbracket \llbracket A \rrbracket_{ref}^k \rrbracket_d^0 = \llbracket A \rrbracket_d^k \quad (7)$$

---

<sup>1</sup>In the paper (Doleschal, 2021) this operation is defined for  $\Sigma$  as  $\text{doc}(\sigma)$