

Expected Overlap of Words in Context (EOWC)

Probability that a randomly picked word from d1 is w_i

Count of word w_i in d1

$$d1 = (x_1, \dots, x_N) \quad x_i = c(w_i, d1) / |d1|$$

$$d2 = (y_1, \dots, y_N) \quad y_i = c(w_i, d2) / |d2|$$

Total counts of words in d1

$$Sim(d1, d2) = d1 \cdot d2 = x_1 y_1 + \dots + x_N y_N = \sum_{i=1}^N x_i y_i$$

Probability that two randomly picked words from d1 and d2, respectively, are identical.