

$$L = \{so, soo, sooo, \dots\}, |L| = \infty$$

the regular expression (RE, re) that describes L is $so+$

a meta-character, 1 or more repetitions.

A, B two tokens WordPiece

$$pmi(A, B) = \log_2 \left(\frac{P(AB)}{P(A)P(B)} \right)$$

$$= \frac{\frac{\#AB}{N}}{\frac{\#A}{N} \frac{\#B}{N}} = N \cdot \frac{\#AB}{\#A \cdot \#B}$$

$$P(AB) = \frac{\# \text{times } AB \text{ appears in text } T}{\# \text{ tokens in } T}$$

$$P(A) = \frac{\# \text{ times } A \text{ appears in } T}{\# \text{ tokens in } T}$$

$$P(B) = \frac{\# \text{ times } B \text{ appears in } T}{\# \text{ tokens in } T} \rightarrow N$$

\Rightarrow ranking pairs AB based on $pmi(A, B) \Leftrightarrow$ ranking using $\frac{\#AB}{\#A \cdot \#B}$