# Lecture 04: Empirical Laws

Tuesday, 28. January 2025   03:50

| Concept | Instance of concept | eg: Will Will |
|---|---|---|
| Type | Token | Type: 1  Token: 2 |

**Type-Token Ratio (TTR):** average a new word appears in text

high TTR: tendency to use new words
Low TTR: tendency to use repeated

Problem with length:                    repeated
                              30000 tks ) 60000 tks
[no. of token ↑ linearly but words do not]
                    ← what can be done ?
Runij average *^

First: (1-1000) + (3-1003) average:

Word Frequency, freq of freq:
    1          3993        3993 words in corpus with freq 1
  >100         102         102 words in corpus with freq more than 1.

**Zipf's Law:**
   (i) Count the frequency of each word     eg: the on and Tom      (i)      (ii)
   (ii) List them in decreasing order          4000 1000 2000 1500    [the, and, Tom, on] words
                                                                      [4000, 2000, 1500, 1000] Freq ) R^n by Zipf's Law
$$f \propto \frac{1}{r} \quad \therefore \quad f \cdot r = K \quad (\text{Frequency} \times \text{rank} = K)$$      [1, 2, 3, 4] Rank

                                                          "Let $P_r$ denote the probability
| R | F |        $f_1 \cdot 50 = K$     $f_1 \cdot 50 = f_2 \cdot 150$    of word of rank $r$ &
|---|---|        $f_2 \cdot 150 = K$    $\boxed{f_1 = 3f_2}$              N denotes the total number
| 50 | $f_1$ |                                                           of word occurences. "
| 150 | $f_2$ |

$$P_r = \frac{f_r}{N} \quad \left(\frac{\text{No. of word}}{\text{Total}}\right)$$

Zipf tell that the:
$$f \propto \frac{1}{r} \qquad \frac{f}{N} = \frac{A}{r} \qquad \boxed{K = A \cdot N} \quad A \approx 0.1//$$

**Sage's Law:**
The number of meaning m of a word obeys the law:
$$m \propto \sqrt{f} \quad \Rightarrow \quad m \propto \sqrt{\frac{1}{r}} \quad \text{or} \quad m \propto \frac{1}{\sqrt{r}}$$

⇒ Rank ≈ 1000 , avg 2.1 meanings
   Rank ≈ 5000 , avg 3 meanings

**Zip Law:**
Word Frequency is inversely proportional to their length.
$$f \propto \frac{1}{l} \qquad \text{short word} \rightarrow \text{more frequency}$$

**Heaps Law:**
How the size of overall vocab grow with the size of the corpus?
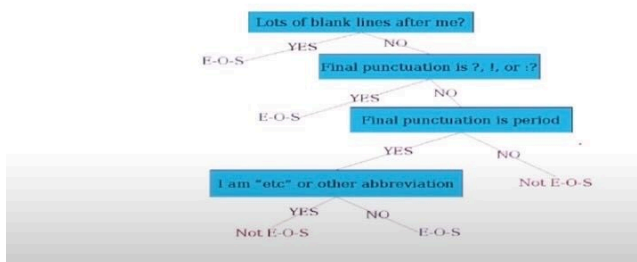Let |V| be the size of vocab & N be the numb of tokens.
$$|V| = kN^{\beta}$$

**Tokenization:** Process of segmenting a string of char → words

Prob 1. Where words start & end.
Approach: End of sent / Not End of sent   (Binary classifier)

Decision Tree: Is this word the end-of-sentence (E-O-S)?

Lots of blank lines after me?
  YES → E-O-S      NO
Final punctuation is ?, !, or :?
  YES → E-O-S      NO
Final punctuation is period
  YES             NO → Not E-O-S
I am "etc" or other abbreviation
  YES → Not E-O-S   NO → E-O-S

(i) EOL parameters: get unique token
(ii) EOL Hyphen & Lexical Hyphen (co-)
(iii) Problem with compound words

**Normalization:**
U.S.A and USA should be matched.

**Case Folding:**
All cases to lower.

**Lemmatization:**
can't → cant  } How  →
won't → wont  }

**Morphology:**
Stem: The core meaning unit
Affix: Bits & pieces adhering to stems

Porter's Algorithm :   Step 1 :

Porter's Algorithm :    Step 1a:
- $sses \longrightarrow ss$ (caresses → caress)
- $ies \longrightarrow i$ (ponies → poni)
- $ss \longrightarrow ss$ (caress → caress)
- $s \longrightarrow \phi$ (cats → cat)

Size of unique word :

$$V = K \times N^{\beta}$$

$N$: total words
$K$: const. $\Big\}$ Heap's Law