# Exercise 1 – IRRS

- CAROL AZPARRENT ESTELA
- MOHANNA FATHOLLAHI

## Exercise 4

We have a document collection with a total of $10^6$ term occurrences. Supposing that terms are distributed in the texts following a power law of the form

$$f_i \cong \frac{c}{(i+10)^2}$$

give estimates of (1) the number of occurrences of the most frequent term; (2) the number of occurrences of the 100-th most frequent term; (3) the number of words occurring more than 2 times. *Hint:* $\sum_{i=11}^{\infty} \frac{1}{i^2} \cong 0.095$.

$N = 10^6$ } total # of tokens

**#1)** the number of occurrence of the most frequent term

# of tokens    of    rank = 1

rank = 1

★ $f(1) = token_1$

rank = 2 $\quad f(2) = token_2$

rank = V $\quad f(V) = token_V$

$\sum_1^V token = N = 10^6$

$\sum_{r=1}^V f(r) = 10^6$

$f_i \cong \frac{c}{(i+10)^2} \Rightarrow \sum_{r=1}^V \frac{c}{(r+10)^2} = 10^6$

$C \cdot \sum_{r=1}^V \frac{1}{(r+10)^2} = 10^6$

$$\boxed{\sum_{r=1}^\infty \frac{1}{(r+10)^2} = \sum_{r=11}^\infty \frac{1}{r^2} \cong 0.095}$$

$\Rightarrow C * 0.095 = 10^6$

$C = \dfrac{10^6}{0.095}$

$C = 1.05 * 10^7$

$$\rightarrowtail \quad f_i = \frac{1.05 * 10^7}{(i + 10)^2}$$

$$\text{rank} = 1 \quad \longrightarrow \quad f(1) = \frac{1.05 * 10^7}{(1 + 10)^2} = 868423 \; \ell$$

**#2)** <u>the number of occurrences</u> of 100-th most
$$\underbrace{\hspace{3cm}}_{\text{\# of tokens}} \qquad \text{of} \quad \underbrace{\text{frequent term}}_{\text{rank} = 100}$$

$$* \quad f(100) = \text{token}_{100} \; *$$

$$\text{rank} = 100 \quad \longrightarrow \quad f(100) = \frac{1.05 * 10^7}{(100 + 10)^2} = 8684$$

**#3)** the number of words occurring more than 2 times.
$$= \qquad f(r) > 2$$

$$f(r) > 2 \quad \Rightarrow \quad \frac{1.05 * 10^7}{(r + 10)^2} > 2$$

$$1.05 * 10^7 > 2 * (r + 10)^2$$
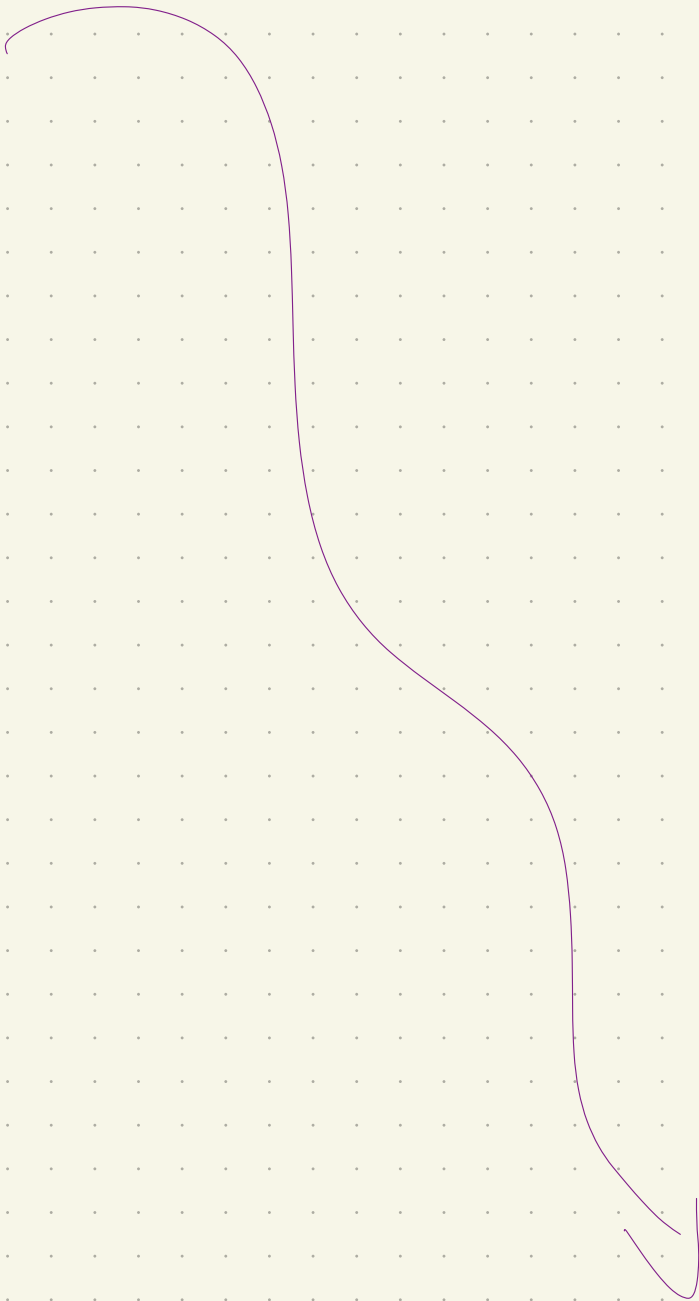$$\frac{1.05 * 10^7}{2} > (r + 10)^2$$
$$(r + 10)^2 < \frac{1.05 * 10^7}{2}$$
$$r^2 + 20r - \frac{1.05 * 10^7}{2} < 0$$
$$r^2 + 20r - 5.25 * 10^6 < 0$$

$\hookrightarrow \quad r_1 = -2301.3$

$$\boxed{r_2 = 22\,81}$$

$\therefore$ the number of word occuring
more than 2 times is 2281.

**Exercise 5**

We are given a random sample of 10,000 documents from a collection containing 1,000,000 documents. We count the different words in this sample, and we find 5,000. Supposing that the collection satisfies Heaps' law with exponent 0.5, give a reasoned estimate of the number of different words you expect to find in the whole collection.

$$d = 5000$$
$$\beta = 0.5$$

$$5000 = K * N_1^{0.5}$$

$$k = \frac{5000}{N_1^{0.5}}$$

total amount of documents : $d = \frac{5000}{N_1^{0.5}} * N_2^{0.5}$

✳ Assuming : $N_2 = 100 * N_1$

$$d = \frac{5000}{N_1^{0.5}} * \left(100 * N_1\right)^{0.5}$$

$$d = \frac{5000}{N_1^{0.5}} * 100^{0.5} * N_1^{0.5}$$

$$d = 5000 * \sqrt{100}$$

$$d = \underline{50000}$$

∴ there are 50000 amount of words in the whole collection approximately

Let us deduce Heaps' law from Zipf's law.

- Let a collection have $N$ word occurrences, with the frequence $f_i$ of the $i$-th most common word proportional to $i^{-\alpha}$, $\alpha > 1$.

- Figure out (from previous exercises) the proportionality constant.

- Estimate the rank $i$ such that $f_i$ is likely to be less than 1.

- Explain why this should roughly be the number of distinct words we expect to see in the collection.

- Deduce that this number is $k \cdot N^\beta$. Tell the values of $k$ and $\beta$ as a function of $\alpha$.

[Note: The given formulation of Zipf's law cannot, for obvious reasons, be taken too literally: If for some large $i$ we have $c \cdot i^{-\alpha} = 0.03$, it makes no sense to say that the $i$th word appears 0.03 times in the collection. More abstractly, one could imagine texts generated by some random process which assigns probability $P(w)$ to the event that a random position in the text contains the word $w$. Then the word with rank 1 is the $w$ with highest $P(w)$, etc. Zipf's law is a statement about the form of the probability distribution $P$. One can then compute rigorously the expected number of distinct words in

a text of length $N$ according to this probabilistic model. Let us just say that we this way we obtain the same $\beta$ but a different $k$.]
[Note 2: It is also possible but a bit more involved to deduce a power law (generalizing Zipf's law) from Heap's law]

$$* \quad \alpha > 1 \longrightarrow f(i) = \frac{c}{i^\alpha}$$

$$* \quad V: \text{total \# of distinct words}$$

$$\sum_{i=1}^{V} f(\ell) = N$$

$$\sum_{i=1}^{V} \frac{c}{i^\alpha} = N$$

$$c * \sum_{i=1}^{V} \frac{1}{i^\alpha} = N$$

$$\boxed{\text{Zeta Riemann} \\ \zeta(\alpha) \approx \sum_{i=1}^{V} \frac{1}{i^\alpha}}$$

$$c = \frac{N}{\overbrace{\sum_{i=1}^{V} \frac{1}{i^\alpha}}^{\zeta(\alpha)}}$$

$$c = \frac{N}{\zeta(\alpha)}$$

$*\quad f(i) < 1$

$$\frac{N}{\varsigma(\alpha)} * \frac{1}{i^\alpha} < 1$$

$$i^\alpha > \frac{N}{\varsigma(\alpha)}$$

$$i > \left(\frac{N}{\varsigma(\alpha)}\right)^{\frac{1}{\alpha}} \qquad \Bigg\{ \quad$$

the words that have a rank "$i$" higher than $\left(\frac{N}{\varsigma(\alpha)}\right)^{\frac{1}{\alpha}}$ have frequency $f(i)$ greater than 1.

these words appear more than 1 time at the collection

$\rightarrow i$: the number of distinct word of the collection. and in Heap's law is equal to $d$.

$$d \approx \left(\frac{N}{\varsigma(\alpha)}\right)^{\frac{1}{\alpha}}$$

$$d \approx \frac{N^{1/\alpha}}{\varsigma(\alpha)^{1/\alpha}}$$

$$\Bigg\{ \quad k = \frac{1}{\varsigma(\alpha)^{1/\alpha}}$$

$$\quad \beta = \frac{1}{\alpha}$$

$$d = k * N^\beta \quad \} \text{ Heaps' law}$$