

Relembrando...

# Domínio

## Smartphones

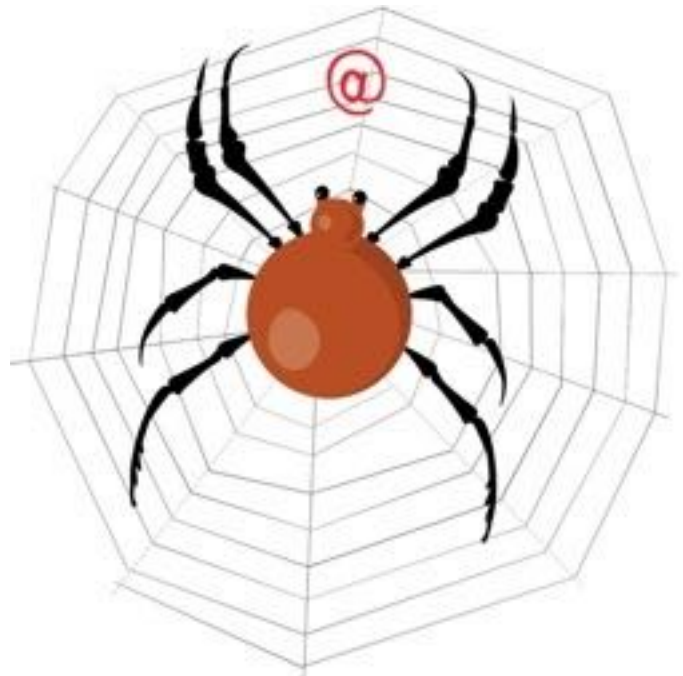


# Sites

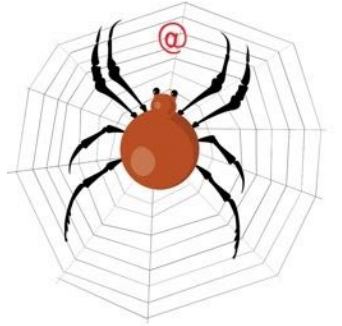
- Americanas
- Extra
- Submarino
- Nagem
- Saraiva
- Casas Bahia
- Ricardo Eletro
- Ponto Frio
- Cissa Magazine
- ShopTime



# Crawler



# Roteiro

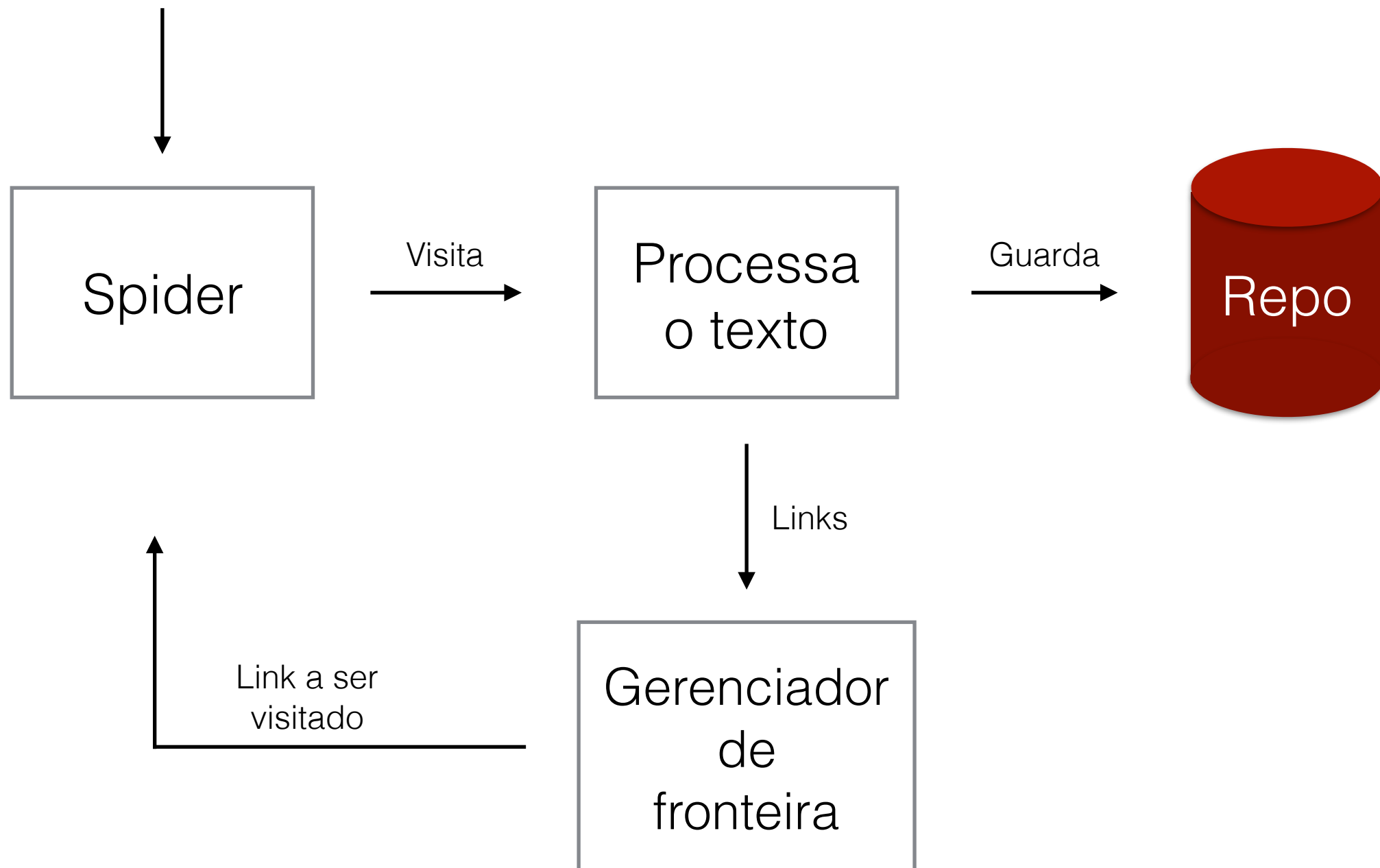


- Crawlers implementados
- Comparativo e conclusões

Crawlers implementados

# Crawler

Semente (url inicial)



# Crawler Busca em Largura



Gerenciador  
de  
fronteira

A rectangular box with a thin black border containing the text "Gerenciador de fronteira" centered within it.

- Adiciona novos links no final da lista



# Crawler Heurístico Sem Peso



Gerenciador  
de  
fronteira

The diagram consists of a single rectangular box with a thin black border. Inside the box, the text 'Gerenciador de fronteira' is centered and arranged in three lines: 'Gerenciador' on the top line, 'de' on the middle line, and 'fronteira' on the bottom line.

- Adiciona novos links que contém alguma palavra relevante na âncora ou na URL no início da lista
- Os que não tiverem, vão pra o final

# Crawler Heurístico Com Peso



Gerenciador  
de  
fronteira

A diagram showing a rectangular box with a thin black border. Inside the box, the text "Gerenciador de fronteira" is centered and arranged in three lines: "Gerenciador" on the top line, "de" on the middle line, and "fronteira" on the bottom line.

- Calcula um peso pra o link baseado na presença de palavras heurísticas na âncora ou na URL
- A lista é ordenada e o próximo link visitado é o que tem maior peso

# Crawler Heurístico Com Peso

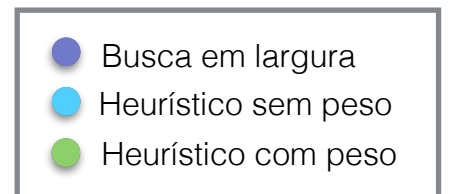
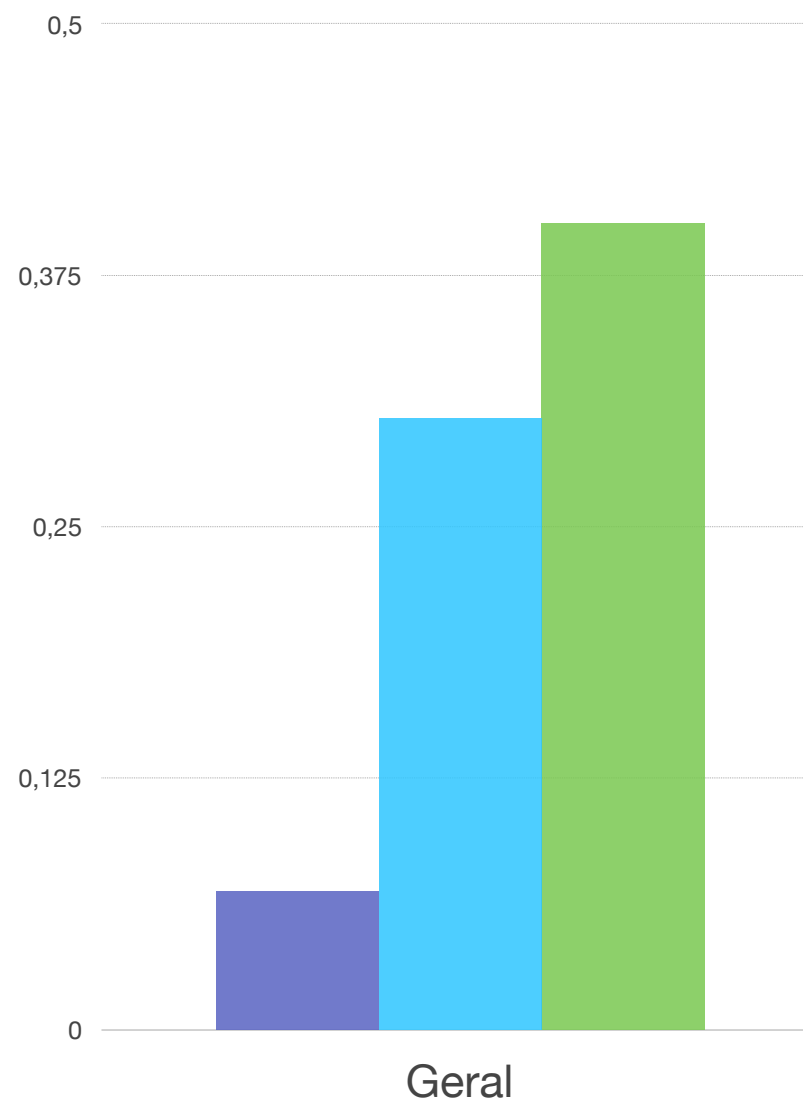
- Conjunto de palavras e pesos utilizado

```
1 samsung 4
2 apple 4
3 motorola 4
4 positivo 4
5 produto 4
6 desbloqueado 5
7 smartphone 2
8 celular 2
9 review -2
10 categoria -3
11 derivacao -3
12 order -3
13 filtro -5
14 avaliacao -7
15 page -7
16 capa -8
17 acessorio -8
18 linha -10
```

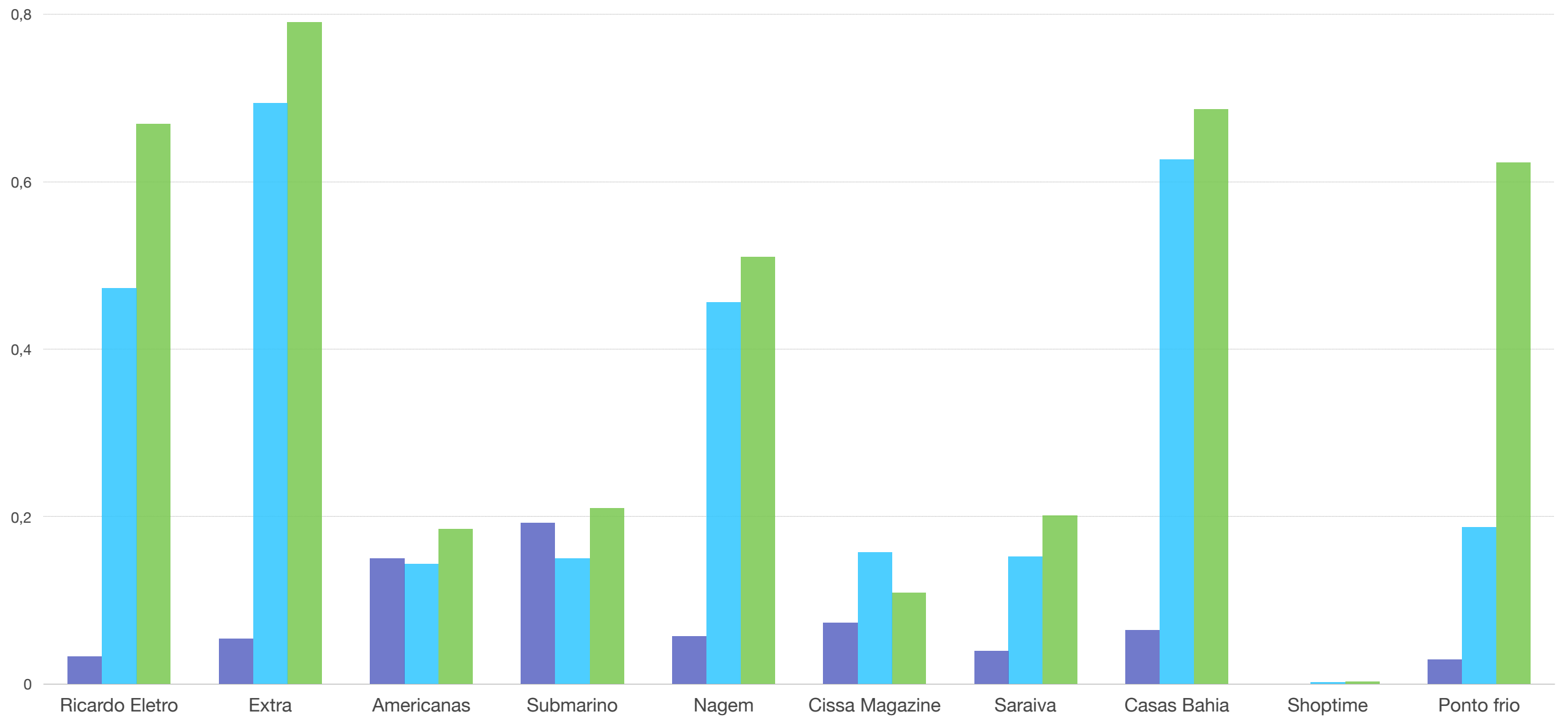
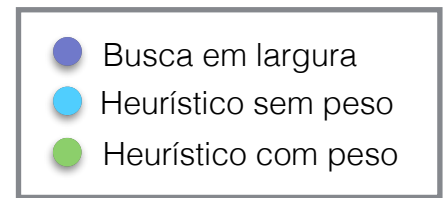
# Comparativo e conclusões

(Harvest ratio)

# Geral



# Por site



# Considerações

- O site shoptime teve um resultado muito ruim por usar URL's relativas em seus htmls. Como o crawler sempre checa se a URL encontrada está no mesmo domínio da semente, as URL's das páginas dos smartphones nunca eram visitadas.
- Alguns sites tiveram uma densidade de produtos muito maiores porque possuem muito mais aparelhos celulares e smartphones em relação a outros, que só tinham modelos das principais marcas (e.g. Motorola, Apple e Samsung)