

Efficient Crawling for Scalable Web Data Acquisition

*29th International Conference on Extending Database Technology
(EDBT 2026)*

Antoine Gauquier¹ Ioana Manolescu² Pierre Senellart¹

¹DI ENS, ENS, CNRS, PSL & Inria

²Inria & Institut Polytechnique de Paris

27 March 2026



PSL 



Inria



Ever increasing volume of content of many forms on the Web:

- Corpora of text and video → Training of ML or DL models
- Structured HTML (especially containing tables) → Knowledge Bases
- Open-access statistic dataset and studies → Analyzing how societies function and doing fact-checking.
- ...

First step for building applications based on Web-published data is to find those data files.

Context (2/3)

Data files come in a wide variety of **formats**:

- **Text** and **video** → .doc, .txt, .mp4, .flv...
- **Structured HTML** → .html
- **Statistic dataset and studies** → .pdf, .xls, .csv...

They are published by different **providers**:

- *National agencies*: INSEE (France), Labor/Census Bureaus (US), ABS (AU)
- *International bodies*: Eurostat, IMF, WB, UN Agencies → some provide **APIs or portals**
- *Smaller entities*: Trade unions, NGOs → publish **highly-specialized data**

Websites can be huge with only a small fraction of published data files

Exhaustively crawling those websites takes a lot of **space** and most importantly **time**:

- Taken by the crawl of a lot of “uninteresting” resources
- Crawling ethics requires to wait (typically ~ 1 sec.) between two HTTP requests \rightarrow For a website of **1M pages**, it would take roughly **11 days** of waiting. . .

Generic crawler for any given website (selected beforehand) \rightarrow we **cannot make any assumption about the website's structure/content**.

Problem Statement

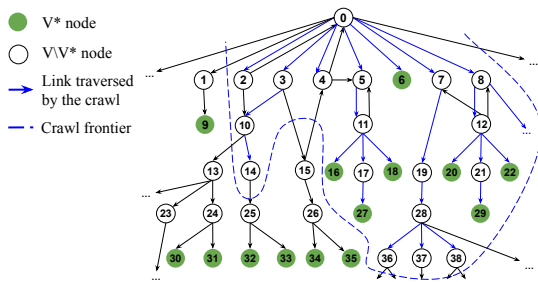
Problem: Given the starting URL of a website, we want to **retrieve** as many **targets** as possible, while **minimizing** the consumption of **resources**. We define:

- A **target**: file uploaded on the website, satisfying some constraints specified by the user (e.g., be a data file: CSV, Excel, PDF...).
- **Resources**: number of **HTTP requests** sent to the server, and **data volume** exchanged with the server.

We want to retrieve the maximum of targets while minimizing the effort to do so.

Formalization: Graph Crawling Problem (2/2)

Given a website graph $G = (V, E, r, \omega, \lambda)$ and a subset V^* of targets of V , the **graph crawling problem** is to find a crawl $T = (V', E')$ of G with $V^* \subseteq V'$, of minimal total cost.



This problem turns out to be **intractable**, as it is **NP-hard** (proof via reduction from **set cover** problem).

This is why we need a **heuristic** yet **effective** solution to solve our crawling problem: an **Efficient Crawler for Scalable Web Data Acquisition**.

Existing Web Crawlers (1/3)

What are **existing Web crawlers** from the literature? Among most simple and classic methods:

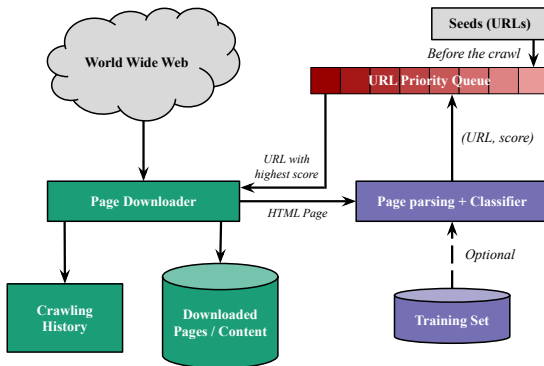
- **Random Crawler**: selects the next hyperlink to visit uniformly at random from the frontier. **Most naive technique**.
- **Breadth-First Search (BFS)**: frontier as a **First-In-First-Out** queue data structure. Efficient when targets are not too deep.
- **Depth-First Search (DFS)**: frontier in a **Last-In-First-Out** stack.

None of them is effective without prior knowledge of the **website's structure**, **which we don't have**

Existing Web Crawlers (2/3)

Focused Crawlers:

- Prioritize some pages during a crawl, following a **focus** (an objective).
- Usually, a **classifier** assesses the **likelihood** that an hyperlink leads to a page part of the focus.
- **Frontier**: a *priority-queue*.



Existing Web Crawlers (3/3)

Two kind of focused crawlers:

- **Topical**: **focus** is a predefined **topic** \rightarrow keywords or sample documents. Most common ones.
- **Non-Topical**: rarer; **focus** is something else, e.g.:
 - ANTHELION [MMB14]: semantic annotations in Web pages
 - ACEBOT [FS15]: as much diverse textual content as possible

In our setting: topical crawlers are **ill-suited**

- We want targets about **any possible subject** covered by the website
- Independently, defining a topic \implies choosing **one (or several) languages**

Best approach: **non-topical focused crawlers**

Our approach → two main research hypotheses:

- Hyperlinks similarly structured in the HTML page were they were found lead to similar content.
- It is possible to learn which hyperlinks are most likely to lead to targets, given their link structure.

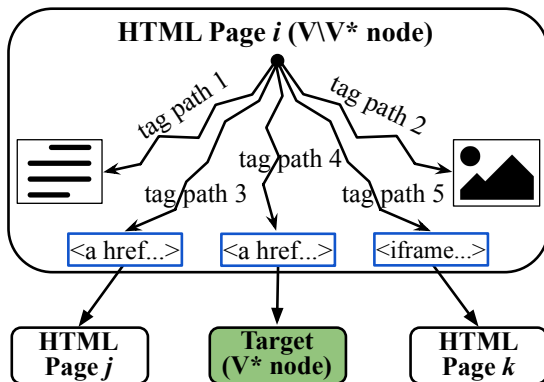
Idea: We can separate, for a given website, parts of it that are rich in targets, from ones that are not:

- without prior knowledge of the website's structure/content
- just assuming existence of some structure

Building Groups of Hyperlinks

How to do it: We represent each encountered hyperlink with its associated **tag path** (its **link structure**).

They are computed from the HTML page's **Document Object Model**:



Building Groups of Hyperlinks

We form **groups of hyperlinks** with a merging strategy based on **similarity** between the tag paths. . .

. . . to **separate the fruitful groups** from the less interesting ones, as we assume that **similar hyperlinks will lead to similar kind of content**.

```
https://cpm.justice.gouv.fr/content/notice-dépôt-d'une-candidature-au-concours-de-greffier-de-tribunal-de-commerce-annexe-1-0  
/html/body.html front not-logged-in one-sidebar sidebar-first page-node/div.jumbotron opm-content-container/div.container/div.row/div.col-sm-3 well-resp/div.block block-block last odd#block-block-7/p/a
```

```
https://www.justice.gouv.fr/lettre-direction-affaires-civiles-du-sceau  
/html/body.html front not-logged-in one-sidebar sidebar-first page-node/div.jumbotron opm-content-container/div.container/div.row/div.col-sm-3 well-resp/div.block block-block first odd#block-block-6/p/a
```

```
https://cpm.justice.gouv.fr/content/notice-de-renseignements-su-l'accès-à-la-profession-de-greffier-de-tribunal-de-commerce-et  
/html/body.html front not-logged-in one-sidebar sidebar-first page-node/div.jumbotron opm-content-container/div.container/div.row/div.col-sm-3 well-resp/div.block block-block last odd#block-block-7/p/a
```

```
https://www.justice.gouv.fr/grands-dossiers/justice-amiable  
/html/body.path-frontpage page-node-type-homepage/[...] /div.fr-col-md-6fr-col-12/article.fr-enlarge-link fr-card fr-card-detail-sm fr-card-horizontal/div.fr-card_body/div.fr-card_content/h3.fr-card_title fr-title/a.fr-card_link
```

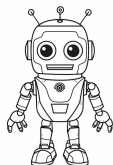
```
https://www.justice.gouv.fr/prise-charge-victimes-d'infractions-penales  
/html/body.path-frontpage page-node-type-homepage/[...] /article.fr-enlarge-link fr-card fr-card-detail-sm fr-card-horizontal/div.fr-card_body/div.fr-card_content/h3.fr-card_titlefr-title/a.fr-card_link
```

```
https://www.justice.gouv.fr/plan-daction-justice  
/html/body.path-frontpage page-node-type-homepage/[...] /div.fr-col-md-6 fr-col-12/article.fr-enlarge-link fr-card fr-card-detail-sm fr-card-horizontal/div.fr-card_body/div.fr-card_content/h3.fr-card_title fr-title/a.fr-card_link
```

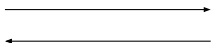
Learning Fruitful Groups: Reinforcement Learning

How can we dynamically separate fruitful groups from other ones?
Reinforcement Learning (RL).

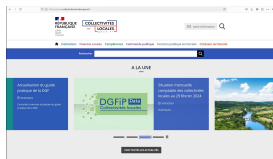
Leveraging **Multi-Armed Bandits** strategy [ACBF02], and especially **Sleeping Multi-Armed Bandits** [KNMS10].



Select a group
Follow a hyperlink from said group



Get a **reward**: a score describing how
“good” was this choice, to improve
future decisions ...



How many **targets** can
we reach **from this**
page?

Reward: Number of reachable, unobserved targets from the HTML page.

Exploration–Exploitation Trade-Off

At each crawling step, the choice of the group is a trade-off between:

- **Exploitation** of groups so far identified as fruitful
- **Exploration** of few times selected groups, that might hide fruitful hyperlinks.

Doing so **works well**... But requires to tell, **just given a hyperlink's URL**, if its a target or not:

- `https://www.assemblee-nationale.fr/14/pdf/ta/ta0686.pdf`
→ Easy
- `https://www.education.gouv.fr/media/90068/download`
→ Manageable
- `https://www.justice.gouv.fr/en/node/9961` → Requires learning...

Naive way: HTTP HEAD requests over the URL to get Content-Type from the header. But too costly.

Idea: Actively train a URL classifier, taking as input an URL, outputting one of the two following **classes**: **target** or **HTML**.

“on-the-fly” training allows the classifier to **adapt** to its environment and to eventual drastic changes in it.

It takes advantage of URLs that are **automatically labeled** when following a hyperlink during the crawl.

Dataset: Selected Websites and Characteristics (1/2)

| | Starting URL | Multilingual Fully Crawled | |
|-----------|---|----------------------------|---|
| <i>ab</i> | https://www.abs.gov.au/ | X | X |
| <i>as</i> | https://www.assemblee-nationale.fr/ | X | X |
| <i>be</i> | https://www.bea.gov/ | X | ✓ |
| <i>ce</i> | https://www.census.gov/ | X | ✓ |
| <i>cl</i> | https://www.collectivites-locales.gouv.fr/ | X | ✓ |
| <i>cn</i> | https://www.cnis.fr/ | X | ✓ |
| <i>ed</i> | https://www.education.gouv.fr/ | X | ✓ |
| <i>il</i> | https://www.ilo.org/ | ✓ | X |
| <i>in</i> | https://www.interieur.gouv.fr/ | X | ✓ |
| <i>is</i> | https://www.insee.fr/ | ✓ | ✓ |
| <i>jp</i> | https://www.soumu.go.jp/ | ✓ | X |
| <i>ju</i> | https://www.justice.gouv.fr/ | X | ✓ |
| <i>nc</i> | https://nces.ed.gov/ | X | ✓ |
| <i>oe</i> | https://www.oecd.org/ | ✓ | ✓ |
| <i>ok</i> | https://okfn.org/ | ✓ | ✓ |
| <i>qa</i> | https://www.psa.gov.qa/ar/Pages/default.aspx | ✓ | ✓ |
| <i>wh</i> | https://www.who.int/ | ✓ | X |
| <i>wo</i> | https://www.worldbank.org/ | ✓ | X |

Dataset: Selected Websites and Characteristics (2/2)

| | #Available (k) | #Target (k) | HTML with Tar. (%) | Target Depth |
|-----------|----------------|-------------|--------------------|----------------------|
| <i>ab</i> | 952.26 | 263.26 | 8.86 | 8.94 (\pm 2.56) |
| <i>as</i> | 994.36 | 93.56 | 4.34 | 4.47 (\pm 4.23) |
| <i>be</i> | 31.23 | 15.84 | 32.19 | 5.73 (\pm 3.21) |
| <i>ce</i> | 988.37 | 257.68 | 3.47 | 4.23 (\pm 0.48) |
| <i>cl</i> | 5.54 | 3.70 | 5.40 | 2.80 (\pm 0.82) |
| <i>cn</i> | 12.80 | 7.49 | 13.87 | 4.26 (\pm 1.59) |
| <i>ed</i> | 102.71 | 10.47 | 3.95 | 11.89 (\pm 13.22) |
| <i>il</i> | 993.92 | 26.56 | 2.53 | 9.63 (\pm 9.83) |
| <i>in</i> | 922.46 | 22.98 | 1.54 | 66.94 (\pm 39.43) |
| <i>is</i> | 285.55 | 168.88 | 41.34 | 5.20 (\pm 1.81) |
| <i>jp</i> | 998.26 | 256.34 | 6.30 | 21.37 (\pm 49.34) |
| <i>ju</i> | 56.61 | 14.85 | 4.85 | 86.91 (\pm 86.30) |
| <i>nc</i> | 309.97 | 84.94 | 18.87 | 3.63 (\pm 1.66) |
| <i>oe</i> | 222.58 | 45.04 | 15.61 | 6.28 (\pm 5.65) |
| <i>ok</i> | 423.12 | 12.95 | 0.74 | 2.64 (\pm 2.89) |
| <i>qa</i> | 4.36 | 2.45 | 4.15 | 3.03 (\pm 0.61) |
| <i>wh</i> | 351.86 | 55.59 | 14.19 | 4.43 (\pm 0.62) |
| <i>wo</i> | 223.67 | 23.10 | 2.38 | 4.52 (\pm 0.69) |

Crawlers (1/2)

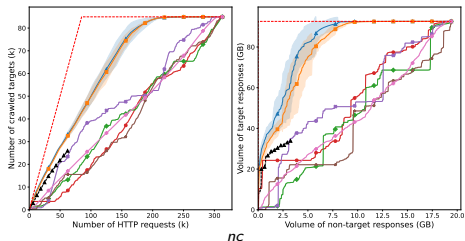
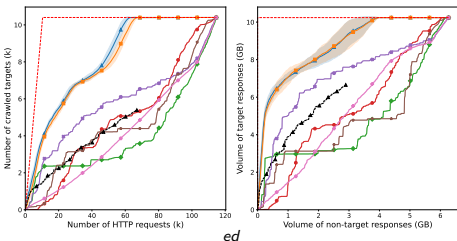
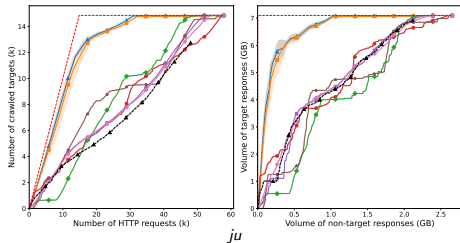
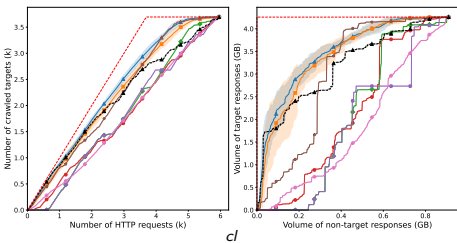
- 1 SB-CLASSIFIER: Our approach with URL classifier
 - 2 SB-ORACLE: Our approach with an **oracle** (perfect, costless URL classifier)
 - 3 RANDOM: Random crawler
 - 4 BFS: Breadth-First Search
 - 5 DFS: Depth-First Search
 - 6 FOCUSED: A standard focused crawler (focus = targets), **without topical features**:
 - approx. depth of the hyperlink
 - 2-gram representation of the hyperlink's URL
 - 2-gram representation of the hyperlink's anchor
- **ablated version** of our approach: neither tag path structure information, nor RL

Crawlers (2/2)

- 7 TP-OFF: Offline-Trained, Tag Path-Based Crawler; an adaptation of ACEBOT [FS15] to target retrieval
 - Crawls 3k first pages with BFS & learns tag paths groups on this subset → requires to compute benefit of each HTML page: we provide it at no cost with an oracle (unfair advantage)
 - Then matches hyperlinks to existing tag path groups (ordered by likelihood), or given a score of 0 if too distant→ ablated version of our approach: learning offline instead of online (closest to our approach)
- 8 TRES [KKP⁺21]: Topical, RL-Based Focused Crawler; using a Bi-LSTM classifier to assess topic relevance
 - requires keywords and sample pages to define the topic: we provide it with 74 terms about the topic “target retrieval” at no cost
 - requires positive examples of fruitful HTML pages: we provide it with 1k such pages, giving partial access to the solution at no cost
 - can only handle URLs to HTML pages in its frontier: we give access to a perfect URL classifier (oracle) at no cost

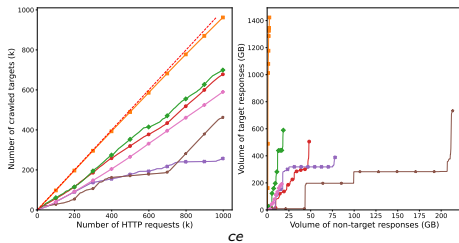
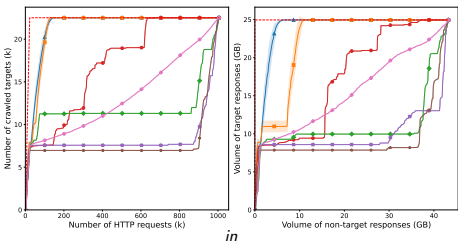
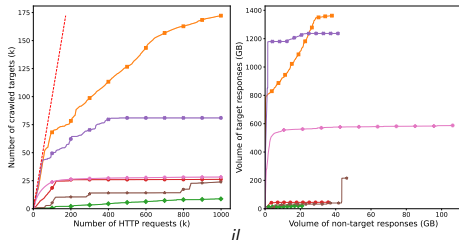
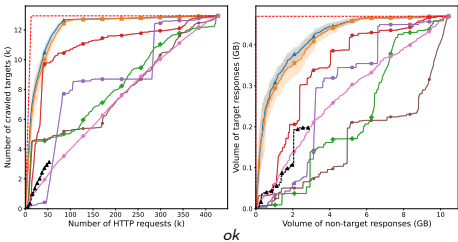
Results (Overview, 1/3)

SB-ORACLE SB-CLASSIFIER FOCUSED TP-OFF BFS DFS RANDOM TRES OMNISCIANT



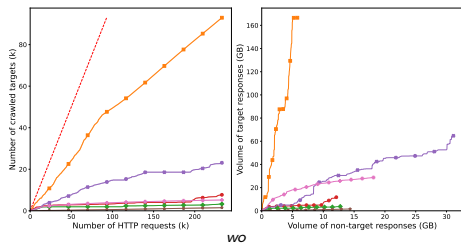
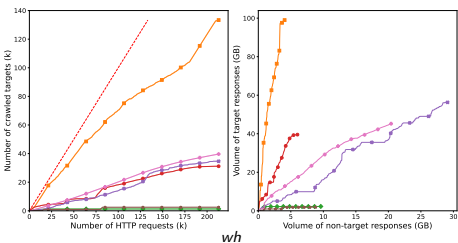
Results (Overview, 2/3)

SB-ORACLE SB-CLASSIFIER FOCUSED TP-OFF BFS DFS RANDOM TRES OMNISCIANT



Results (Overview, 3/3)

SB-ORACLE SB-CLASSIFIER FOCUSED TP-OFF BFS DFS RANDOM TRES OMNISCIENT



Takeaways

- Most applications are based on Web-published data: finding them is not easy. Exhaustive methods are too costly.
- Finding an optimal crawl of a website is intractable → need of a heuristic technique.
- We presented a novel crawling technique based on reinforcement learning to dynamically learn interesting parts of the Website, that are rich in targets.
- We equip our crawler with a URL classifier that is dynamically trained with automatically-generated data, giving the class (HTML or target) of an URL at no extra cost.
- The crawler outperforms standard baselines and three adapted from the literature, sometimes reaching more than 90% of targets within less than 10% of the webpages.

Thank you for your attention!
Any questions?

Code and dataset information

https://github.com/AntoineGauquier/efficient_crawler_for_scalable_web_data_acquisition/

Bibliography I

 Peter Auer, Nicolò Cesa-Bianchi, and Paul Fischer.

Finite-time Analysis of the Multiarmed Bandit Problem.

Machine Learning, 47(2):235–256, 2002.

 Muhammad Faheem and Pierre Senellart.

Adaptive web crawling through structure-based link classification.

In Robert B. Allen, Jane Hunter, and Marcia Lei Zeng, editors, *Digital Libraries: Providing Quality Information - 17th International Conference on Asia-Pacific Digital Libraries, ICADL 2015, Seoul, Korea, December 9-12, 2015, Proceedings*, volume 9469 of *Lecture Notes in Computer Science*, pages 39–51. Springer, 2015.

 Andreas Kontogiannis, Dimitrios Kelesis, Vasilis Pollatos, Georgios Paliouras, and George Giannakopoulos.

Tree-based focused web crawling with reinforcement learning.

CoRR, abs/2112.07620, 2021.



Robert Kleinberg, Alexandru Niculescu-Mizil, and Yogeshwer Sharma.
Regret Bounds for Sleeping Experts and Bandits.
Machine Learning, 80(2):245–272, 2010.



Robert Meusel, Peter Mika, and Roi Blanco.
Focused crawling for structured data.

In Jianzhong Li, Xiaoyang Sean Wang, Minos N. Garofalakis, Ian Soboroff, Torsten Suel, and Min Wang, editors, *Proceedings of the 23rd ACM International Conference on Conference on Information and Knowledge Management, CIKM 2014, Shanghai, China, November 3-7, 2014*, pages 1039–1048. ACM, 2014.

Why don't you just Google it?

Examples based on [Google](#) services (conducted early 2025):

- ① Classical [Google Search Engine](#) (and its [Programmable Search Engine](#)¹), by filtering with `site:` and `filetype:`
 - For `site:https://www.justice.gouv.fr/`
 - on `filetype:pdf` → 302 results, while > 9000 exist
 - on `filetype:tsv` → filetype is **not even recognized by GS**
 - For `https://www.ilo.org/`
 - on `filetype:pdf` → 641 results, while > 41k exist
 - on `filetype:zip` → 0 result, while > 2200 exist

¹<https://programmablesearchengine.google.com/>

- ② **Google Dataset Search²**: aims at **target search**, but...
 - For `site:https://www.justice.gouv.fr/` → 109 tabular data files, out of > 1100
 - For `https://www.ilo.org/` → 93 datasets, out of > 170k
 - For `https://www.census.gov/` → 312 datasets, out of > 800k
- ③ **Google Public Data Explorer**: aims at **public data exploration**, but...
 - build for human readers → only allows **manual exploration**
 - ranges over a **closed set of providers**, missing huge ones (US and UN agencies for instance)

²<https://datasetsearch.research.google.com/>

- Do not fully index any given website, and anyway limit available results to 1000.
- Don't offer any transparency or control over their process.

We devise a method that retrieves targeted data files within a budget, with hopefully fully explainable decisions.

Evaluation metrics

- **Intra-website:** % of requests needed to reach 90% of targets
- **Inter-website:** Misclassification Rate (MR)

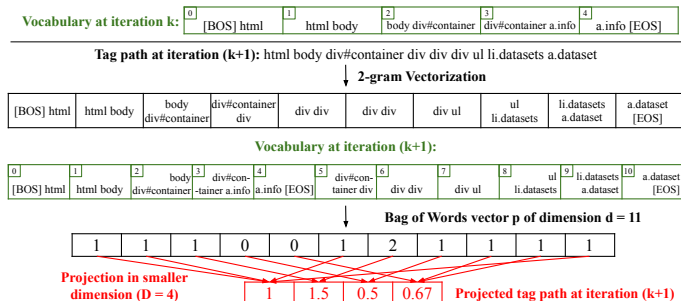
Feature sets

- **URL_ONLY:** URL features only
- **URL_CONT:** URL + anchor text + DOM path + surrounding text

Models LR: Linear Regression SVM: Support Vector Machine NB: Naive Bayes PA: Passive-Aggressive

| Variant | <i>be</i> | <i>cl</i> | <i>cn</i> | <i>ed</i> | <i>in</i> | <i>is</i> | <i>ju</i> | <i>nc</i> | <i>oe</i> | <i>ok</i> | <i>qa</i> | MR |
|--------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-------------|-------------|
| URL_ONLY-LR | 82.1 | 75.1 | 71.3 | 53.2 | 11.7 | 76.1 | 36.5 | 52.6 | 60.7 | 15.9 | 62.3 | 2.62 |
| URL_ONLY-SVM | 82.7 | 75.7 | 71.8 | 63.6 | 11.3 | 76.0 | 37.4 | 52.2 | 63.5 | 16.7 | 61.5 | 2.99 |
| URL_ONLY-NB | 82.9 | 75.2 | 72.1 | 53.7 | 11.4 | 76.3 | 35.8 | 52.7 | 59.7 | 18.0 | 63.1 | 2.92 |
| URL_ONLY-PA | 82.3 | 74.4 | 71.7 | 53.3 | 11.1 | 75.8 | 36.7 | 51.6 | 60.5 | 15.9 | 60.9 | 2.56 |
| URL_CONT-LR | 82.2 | 74.4 | 71.9 | 54.3 | 11.3 | 76.4 | 37.8 | 52.9 | 64.7 | 16.8 | 60.0 | 5.93 |
| URL_CONT-SVM | 82.6 | 75.0 | 71.8 | 52.8 | 11.6 | 76.4 | 38.8 | 53.1 | 61.1 | 18.7 | 60.1 | 6.36 |
| URL_CONT-NB | 84.1 | 74.7 | 71.9 | 53.6 | 11.4 | 75.7 | 35.5 | 52.3 | 59.9 | 19.1 | 60.4 | 7.15 |
| URL_CONT-PA | 82.5 | 75.1 | 71.9 | 53.6 | 11.6 | 76.2 | 38.4 | 52.1 | 62.6 | 16.1 | 60.6 | 4.12 |

Building Groups of Hyperlinks – Merging strategy



Merging strategy: all groups are represented with a **centroid**. For each tag path:

- 1 Compute its set of (approx.) k -NN, and an associated **similarity score**
- 2 If the top-score $> \theta$ (threshold): **merge** and update the **centroid**
- 3 Otherwise: create a **new group**

URL Classifier Results (Overview)

Confusion matrix of our URL classifier, on the fully-crawled websites (average over 15 runs)

| True \ Predicted | HTML (%) | Target (%) | Neither (%) |
|-------------------------|-----------------|-------------------|--------------------|
| HTML | 58.04 | 1.37 | 0.00 |
| Target | 0.75 | 32.19 | 0.00 |
| Neither | 5.34 | 2.41 | 0.00 |

SD Retrieval Across Sample Targets

Experimental setup

- Random sample of $7 \times 40 = 280$ targets
- 7 diverse websites

| | <i>be</i> | <i>ed</i> | <i>is</i> | <i>in</i> | <i>nc</i> | <i>oe</i> | <i>wh</i> |
|---------------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| SD Yield (%) | 82 | 35 | 93 | 40 | 83 | 60 | 40 |
| Mean # SDs / Target | 9.1 | 2.8 | 2.9 | 2.1 | 2.1 | 4.9 | 1.4 |

- SDs are present in a substantial fraction of targets across all websites
- Even non-statistical websites show notable SD presence
- SDs are concentrated: when present, typically > 2 SDs per target