

Do Not Crawl in the DUST: Different URLs with Similar Text

Ziv Bar-Yossef, Idit Keidar, Uri Schonfeld

WWW'07

CSE 450 Web Mining

Presented by Zaihan Yang

Introduction & Contribution

- Propose a novel algorithm DustBuster for uncovering DUST.
 - Discover DUST rules from a URL list
 - *Mainly focus on the substring substitution rules*
 - *Introduce 3 heuristic methods.*
 - *Eliminate redundant rules.*
 - *Validate DUST rules.*
 - Use DUST rules to transform URLs into canonical form
- Main feature: Mine DUST from crawl logs or web server logs instead of examining page content.

Problem Identification

- *What is DUST?* Different URLs with Similar Context.
E.g. <http://google.com/news> & <http://news.google.com>.
- *How generated?*
Aliases, redirection, dynamically generated pages, etc.
- *Features of DUST?*
Not casual: with certain rules.
Not universal: specific to web sites.
- *What advantage for uncovering them?*
Reduce overhead in crawling, indexing, and catching.
Increase accuracy of page metrics, like PageRank.

Problem Definition

- **URL:** strings over Σ starts with “^” and ends with “\$”.
- **DUST:** $(u1, u2)$ are DUST if $doc(u1)$ and $doc(u2)$ similar.
Shingling resemblance measure.
- **URL List:**
a URL http return code size of the document document sketch.
- **DUST rules:** instance of rules, support of rules
URL pair: $(u1, u2)$ Rule: $\alpha \rightarrow \beta$ if $u1 = p\alpha s$ and $u2 = p\beta s$.
 $u1 = ^{http://\underline{www.site.com/index.html}}\$$
 $u2 = ^{http://www.site.com}\$$
DUST rule: “index.html\$” \rightarrow ”\$”

Large Support Heuristic

- Main idea: The support of a valid DUST rule is large.
- How to compute support size
 - **Envelope** of string α
 - URL $u=p\alpha s$, (p, s) is the envelop of α
 - $u=\text{http://www.site.com/index.html}$ $\alpha=\text{“index”}$
 $p=\text{“^http://www.iste.com/”}$ $s=\text{“.html$”}$
 - $E_L(\alpha)$: a set of envelops of α in URLs, each of which appear in \mathcal{L} and has α as a substring.
 - Theorem:

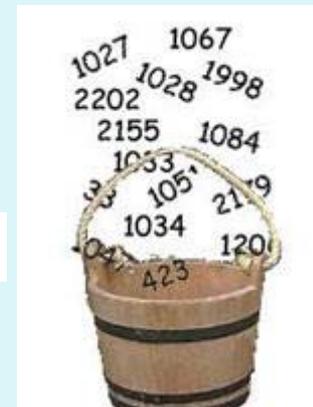
$$|\text{support}_L(\alpha \rightarrow \beta)| = |E_L(\alpha) \cap E_L(\beta)| \quad |\text{support}_L(\alpha \rightarrow \beta)| = |\text{support}_L(\beta \rightarrow \alpha)|$$

$\alpha \neq \beta$, α and β are non-empty and non-semiperiodic.

Small Bucket Heuristic

- Main idea: Much of the support of valid DUST rules belong to small buckets.
- **Bucket:**
 - For an envelope (p, s) , a bucket (p, s) is the set of all substrings α satisfying that $p\alpha s \in \mathcal{L}$.
 - Namely, if (p, s) belongs to many envelope set
 - $E_L(\alpha_1), E_L(\alpha_2), \dots, E_L(\alpha_k)$, then $\alpha_1, \alpha_2, \dots, \alpha_k$ constitutes the bucket of (p, s) .

[^http://www.domain.name/story_](http://www.domain.name/story_)



\$

Similarity Likelihood Heuristic

- Main Idea: The likely similar support of a valid DUST rule is large.
- Similar Content:
- Document sketch:
 - Obtained from previous crawl.
 - Shingling resemblance.
- Size match:
 - Obtained from web server logs.
 - For each URL, a min and max size, size interval.
 - $(u1, u2)$ is similar if interval overlaps.

Algorithm Framework

- Input: URL list
 - Detect likely DUST rules
 - Eliminate redundant rules
 - Validate DUST rules using samples:
 - Eliminate DUST rules that are “wrong”
 - Further eliminate duplicate DUST rules
- No Fetch Required
- 

Detecting likely DUST rules

- **Input:** a URL list L .
- **Output:** an ordered list of pairs, each representing two DUST rules whose support beyond a threshold MS .
- ST: Substring table
- IT: Instance table
- RT: Rule table

```
1: Function DetectLikelyRules(URLList  $\mathcal{L}$ )
2: create table ST (substring, prefix, suffix, size_range/doc_sketch)
3: create table IT (substring1, substring2)
4: create table RT (substring1, substring2, support_size)
5: for each record  $r \in \mathcal{L}$  do
6:   for  $\ell = 0$  to  $S$  do
7:     for each substring  $\alpha$  of  $r.url$  of length  $\ell$  do
8:        $p :=$  prefix of  $r.url$  preceding  $\alpha$ 
9:        $s :=$  suffix of  $r.url$  succeeding  $\alpha$ 
10:      add ( $\alpha$ ,  $p$ ,  $s$ ,  $r.size\_range/r.doc\_sketch$ ) to ST
11: group tuples in ST into buckets by (prefix,suffix)
12: for each bucket  $B$  do
13:   if ( $|B| = 1$  OR  $|B| > T$ ) continue
14:   for each pair of distinct tuples  $t_1, t_2 \in B$  do
15:     if (LikelySimilar( $t_1, t_2$ ))
16:       add ( $t_1.substring, t_2.substring$ ) to IT
17: group tuples in IT into rule_supports by (substring1,substring2)
18: for each rule_support  $R$  do
19:    $t :=$  first tuple in  $R$ 
20:   add tuple ( $t.substring1, t.substring2, |R|$ ) to RT
21: sort RT by support_size
22: return all rules in RT whose support size is  $\geq MS$ 
```

Eliminating Redundant Rules

- (“co.il/story?id=”, “.co.il/story_”) & (“story?id=”, “story_”)
- **Refinement of rule pairs**
 - A rule A refines a rule B if $\text{support}(A) \subseteq \text{support}(B)$
 - Rule $\alpha' \rightarrow \beta'$ refines rule $\alpha \rightarrow \beta$, if there is envelope (γ, δ) satisfying $\alpha' = \gamma\alpha\delta$ and $\beta' = \gamma\beta\delta$.
- **Eliminate method**
 - If A refines B and $\text{support}(A) \subset \text{support}(B)$, then keep the broader one.
 - If A refines B and $|\text{support}(A)| = |\text{support}(B)|$, then remove the refined one.

```
1: Function EliminateRedundancies(pairs_list  $\mathcal{R}$ )
2: for  $i = 1$  to  $|\mathcal{R}|$  do
3:   if (already eliminated  $\mathcal{R}[i]$ ) continue
4:   for  $j = 1$  to  $\min(\text{MW}, |\mathcal{R}| - i)$  do
5:     if ( $\mathcal{R}[i].\text{size} - \mathcal{R}[i+j].\text{size} >$ 
            $\max(\text{MRD} \cdot \mathcal{R}[i].\text{size}, \text{MAD}))$  break
6:     if ( $\mathcal{R}[i]$  refines  $\mathcal{R}[i+j]$ )
7:       eliminate  $\mathcal{R}[i+j]$ 
8:     else if ( $\mathcal{R}[i+j]$  refines  $\mathcal{R}[i]$ ) then
9:       eliminate  $\mathcal{R}[i]$ 
10:    break
11: return  $\mathcal{R}$ 
```

Validating DUST rules

- Validate rules via fetching a small number of pages.

```
1:Function ValidateRule( $\mathcal{R}$ ,  $\mathcal{L}$ )
2: positive := 0
3: negative := 0
4: while (positive <  $(1 - \epsilon)N$  AND negative <  $\epsilon N$ ) do
5:   u : $\leftarrow$  a random URL from  $\mathcal{L}$  on which applying  $\mathcal{R}$  results
      in a different URL
6:   v := outcome of application of  $\mathcal{R}$  to u
7:   fetch u and v
8:   if (fetch u failed) continue
9:   if (fetch v failed OR DocSketch(u)  $\neq$  DocSketch(v))
10:     negative := negative + 1
11:   else
12:     positive := positive + 1
13:   if (negative  $\geq \epsilon N$ )
14:     return FALSE
15: return TRUE
```

Figure 3: Validating a single likely rule.

```
1:Function Validate(rules_list  $\mathcal{R}$ , test_URLList  $\mathcal{L}$ )
2 create an empty list of rules LR
3: for i = 1 to  $|\mathcal{R}|$  do
4:   for j = 1 to i - 1 do
5:     if ( $\mathcal{R}[j]$  was not eliminated AND  $\mathcal{R}[i]$  refines  $\mathcal{R}[j]$ )
6:       eliminate  $\mathcal{R}[i]$  from the list
7:       break
8:   if ( $\mathcal{R}[i]$  was eliminated)
9:     continue
10:  if (ValidateRule( $\mathcal{R}[i].\alpha \rightarrow \mathcal{R}[i].\beta$ ,  $\mathcal{L}$ ))
11:    add  $\mathcal{R}[i].\alpha \rightarrow \mathcal{R}[i].\beta$  to LR
12:  else if (ValidateRule( $\mathcal{R}[i].\beta \rightarrow \mathcal{R}[i].\alpha$ ,  $\mathcal{L}$ ))
13:    add  $\mathcal{R}[i].\beta \rightarrow \mathcal{R}[i].\alpha$  to LR
14:  else
15:    eliminate  $\mathcal{R}[i]$  from the list
16: return LR
```

Figure 4: Validating likely rules.

- Validation count: N
- Refutation threshold ϵ
- α shrinks β if $|\alpha| \geq |\beta|$
- Prefer to be considered.

URL canonization

- Canonization problem
 - A canonical URLs subset $Cu_s \subseteq U_s$
 - A canonization: Mapping $C: U_s \rightarrow Cu_s$
- Algorithm
 - Input: a URL u and a list of valid DUST rule R .
 - Repeatedly apply u to all the rules in R . until
 - u is unchanged, or
 - A predetermined max iteration is reached.

Experiment Results

- Dataset

➤ 4 websites:

Web Site	Log Size	Unique URLs
Forum Site	38,816	15,608
Academic Site	344,266	17,742
Large News Site	11,883	11,883
Small News Site	9,456	9,456

- Parameter Settings:

Substring Length	Bucket Size	MRD	MAD	MW	MS	Validation Count	Refutation Threshold	Max Iterations
35	(6, 11)	5%	1	1100	3	100	5%-10%	10

- Metrics:

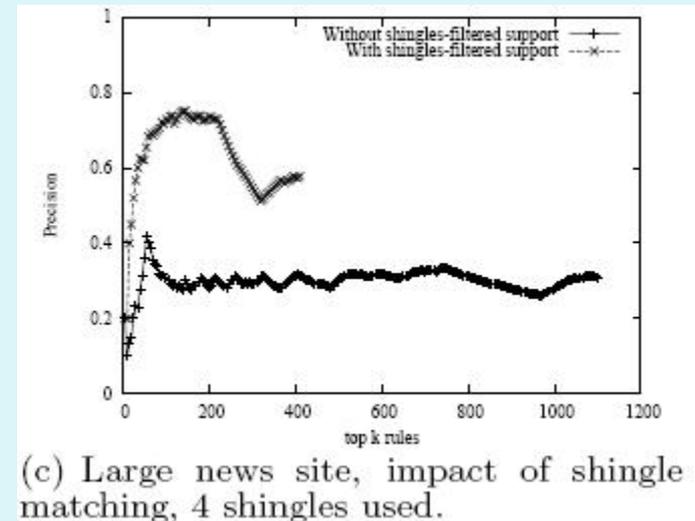
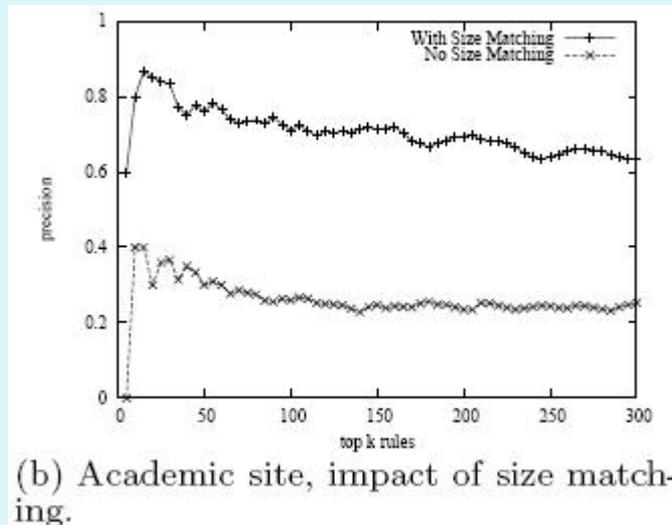
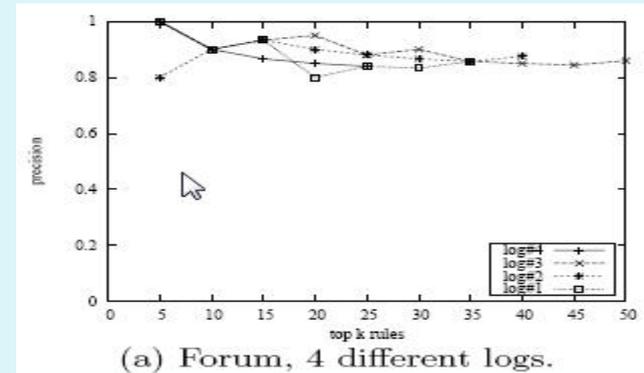
Precision Discovered Redundancy Coverage

Results

- Detecting and Eliminating

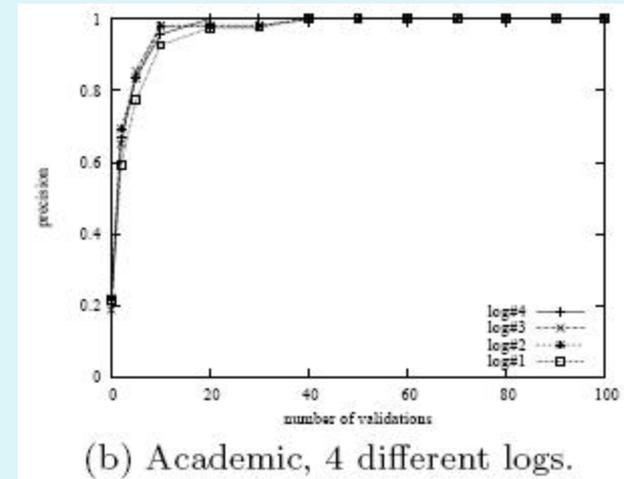
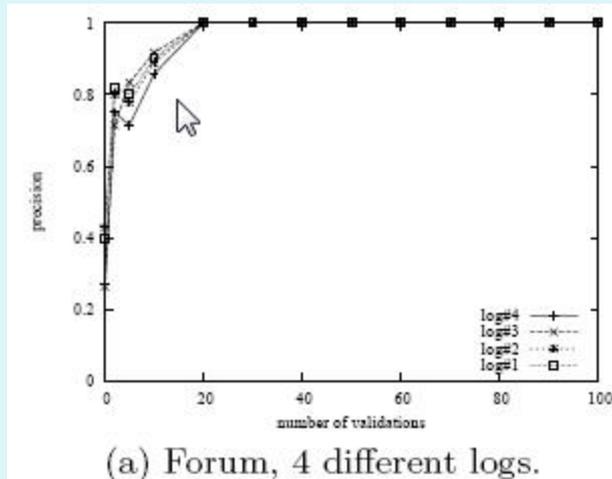
Web Site	Rules Detected	Rules Remaining after 2nd Phase
Forum Site	402	37 (9.2%)
Academic Site	26,899	2,041 (7.6%)
Large News Site	12,144	1,243 (9.76%)
Small News Site	4,220	96 (2.3%)

Table 2: Rule elimination in second phase.



Results

- Validation



Web Site	Valid Rules Detected
Forum Site	7
Academic Site	52
Large News Site	62
Small News Site	5

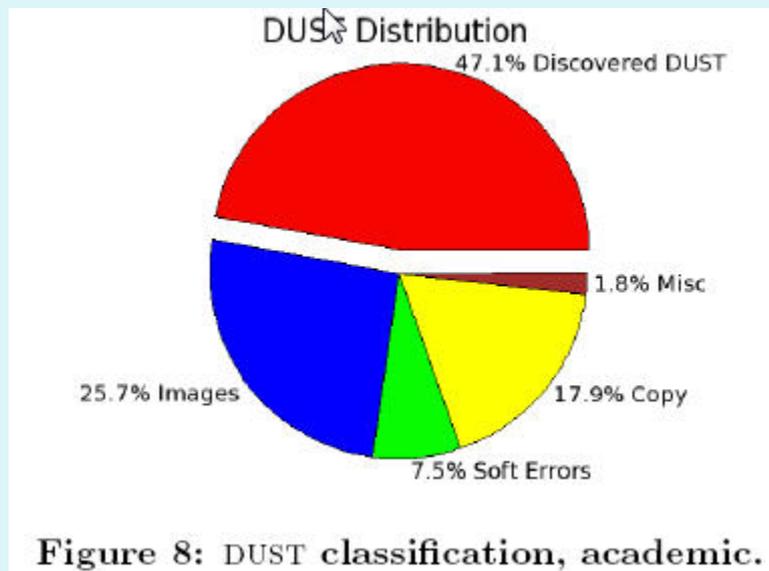
Table 3: The number of rules found to be valid.

- 1 “.co.il/story_” → “.co.il/story?id=”
- 2 “\&LastView=\&Close=” → “”
- 3 “.php3?” → “?”
- 4 “.il/story_” → “.il/story.php3?id=”
- 5 “\&NewOnly=1\&tvqz=2” → “\&NewOnly=1”
- 6 “.co.il/thread_” → “.co.il/thread?rep=”
- 7 “http://www.../story_” → “http://www.../story?id=”

Figure 7: The valid rules detected in the forum site.

Results

Coverage



Savings in crawl size

Web Site	Reduction Achieved
Academic Site	18%
Small News Site	26%
Large News Site	6%
Forum Site(using logs)	4.7%

Table 4: Reductions in crawl size.

Conclusions

- Introduced a problem of mining site-specific DUST rules.
- Proposed the DustBuster algorithm for mining DUST from a URL list.
- Mining DUST rules can
 - Reduce crawling overhead by up to 26%.
 - Reduce indexing overhead
 - Benefit canonizing URL names, and increase accuracy of page metrics.