

Smart Crawler - things to look for

Sites have

- hot deals
- daily sale
- Friday Sale
- bargain bin
- weekly ~~sp~~ special
- daily ad
- specials

- coupons

go to these sites

OCR the gif
or

read the ~~url~~ url's

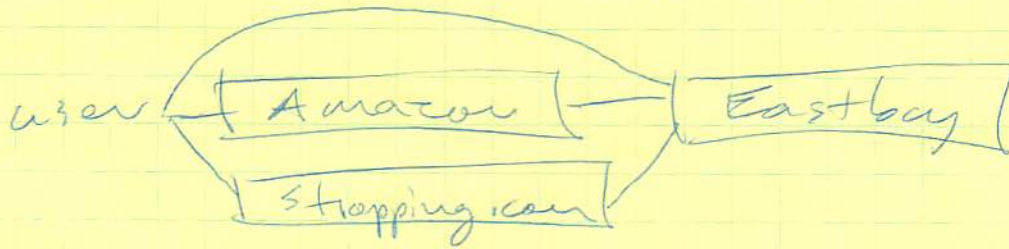
update ^{each item in} the db

create one page w/ all of the deals

~~Duplicating~~ ~~Agent Site~~

Example Amazon has storefronts

Each storefront is a different store



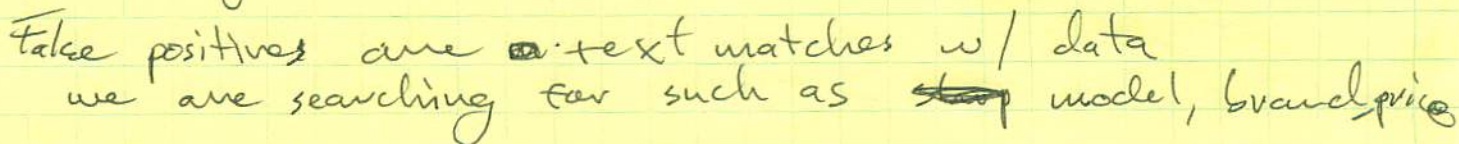
HTML "Anchor"

an " " is an HTML code fragment ^{report} ^{db template} which appears on each page on a site. ^{On a shopping site} The anchor ~~is~~ can be one of the following:

- "click to buy"
- "add to shopping cart"
- "add to " bag "
- "view shopping cart"
- "update cart"
- drop down menu w/ colors/size

The anchor can appear in the following formats in the HTML page

- a gif (need to OCR the gif to recognize the letters in the image)
- a script (need to interpret the script to determine the ~~words~~ words associated w/ the script)



Since the χ^2 test is a test of the fit of the observed frequencies to the expected frequencies, it is a test of the fit of the observed frequencies to the expected frequencies.

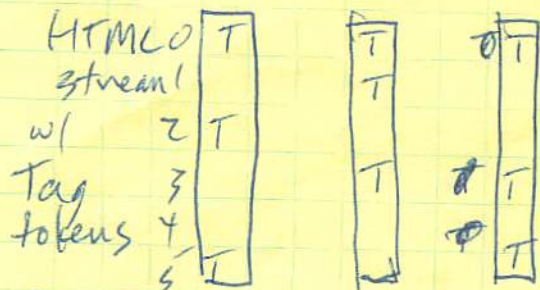


The χ^2 test is a test of the fit of the observed frequencies to the expected frequencies. It is a test of the fit of the observed frequencies to the expected frequencies.

The χ^2 test is a test of the fit of the observed frequencies to the expected frequencies. It is a test of the fit of the observed frequencies to the expected frequencies.

The χ^2 test is a test of the fit of the observed frequencies to the expected frequencies. It is a test of the fit of the observed frequencies to the expected frequencies.

① ~~linearize~~ linearize HTML

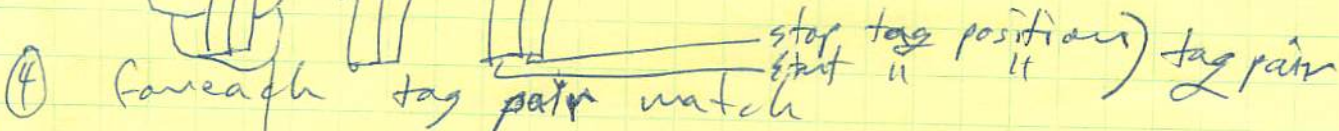


② find tag positions in HTML



tag type
tag position in HTML matrix

③ match tags in lists



positions in tag position vector

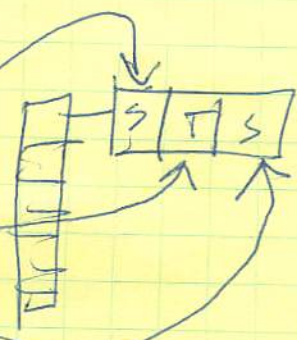
match the ~~look at~~ examine the contents of the tag pair if the tag pair contains

- ① currency
- ② price
- ③ description
- ④ string that matches model name

record ~~int~~ positions & type

this is the signature

vector start tag position
type (s)
stop tag position



linearize the HTML

```

T      → T
...
/T     → /T

```

Match the tags

create a new list of indices of tags into the HTML list

T = tag
 x = text
 s = string
 a = attribute

| | |
|---|----|
| 0 | T |
| 1 | x |
| 2 | /T |
| 3 | T |
| 4 | a |
| 5 | a |
| 6 | /T |
| 7 | T |
| 8 | s |
| 9 | / |

Tag list

| | | | | | |
|---|----|---|----|---|----|
| 0 | 2 | 3 | 6 | 7 | 9 |
| T | /T | T | /T | T | /T |

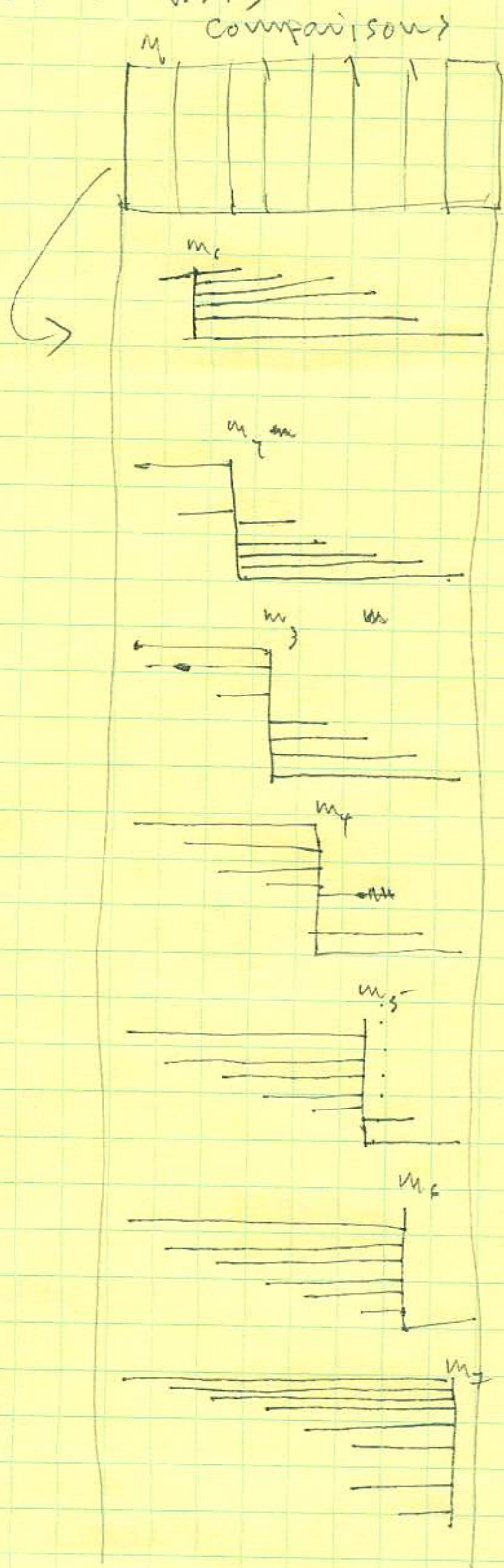
repeat for each candidate page

master is compared against
all other ~~sets~~ lists

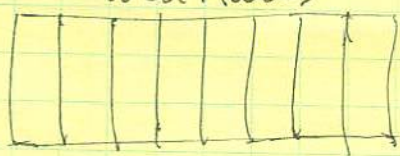
contains 04/13

contains HTML

This wants for inter
page compare
or intra page compare
(tables)



common tag
locations



Signature

list of
tags

now
break down
the contents
of the strings



now
compare
the
extracted HTML
signs
do they
match

compare
common tag
locations
merge the results 2
find differences?

extraction of db fieldnames and/or
db values from (e.g. English) sentences / fragments /
natural language phrases

~~from~~ sentences / frags / phrases ~~is~~ contains

The ↓ must be separated in the ↓

The separation process:

- ① And boundaries between ↓
 to ← can be
- punctuation (#) or parenthesis or comma or colon, ...
 - db field names
 - known db field values (company name)
 - alphanumeric text which is a model #

②

| | | |
|---------|------------------------|---------|
| ge | amazing microwave oven | geam101 |
| company | model name | model # |

To analyze the above string

- ① search for alphanumeric text
 if found search in master model # dictionary
 this is a shortcut compare the comp. name
 and the model name in the db to the
 text string fragments to other strings on the page
- ② if # fails search for company name.

links

Time/Date identification
 there are many formats
 search for ~~the~~ " strings in the page

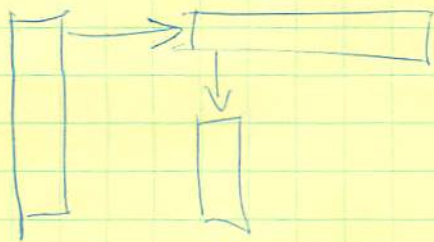
currency identification
 there are many formats
 search for the " strings in the page

Institutions ident (schools, univ's)
 look up strings in institution list

Linear search of HTML

read each token in the HTML page and
look it up in a ~~the~~ taxonomy list

list

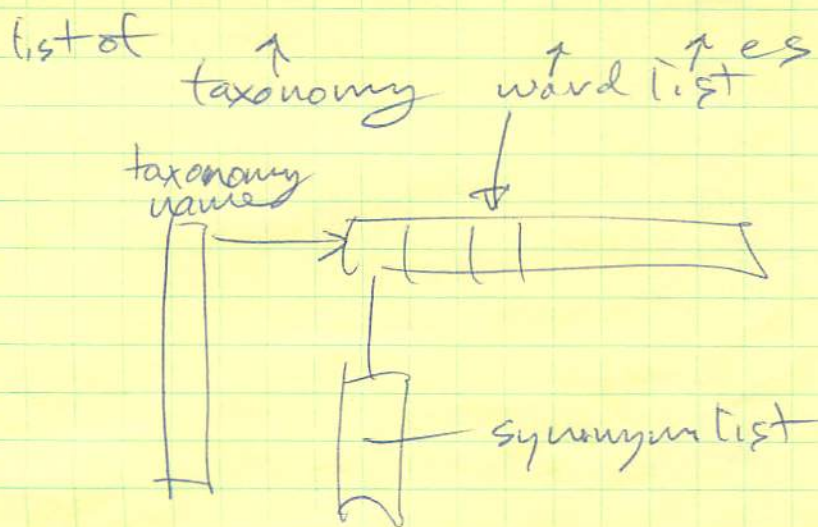


taxonomy list
group of keywords by
topic

synonym in list

~~the~~

Patricia Tree for fast look up?



There are several different taxonomy lists

list of data base template synonym lists
 (each field in the record has a field name and a " value the field name can be many different names which is stored in a synonym list

example

| field name | price | \$100.00 |
|------------|-------|----------|
| the field | | |

 ← either a formatted field
 (currency
 date or time format
 etc.,
 or a word from a word list
 (e.g. product names
 model numbers).