Stitching

- The state of the s
- Paste & copy of fragments of other web pages, uniting them into a single page
- → automatic way to create «interesting» content to populate a site (various search engine also use *global bonuses* to measure how much *information* does a site offer)

Broadening



- Insert not only the keywords we selected, but also opportune synonyms or related keywords/phrases
- This is anyway useful to better cover user queries, but not only:

Broadening (cont.)



It is also helpful because many search engines also use measures of similarity among keywords in order to give added bonuses

Example

If we look for Disney and within a web page there is Winnie the Pooh, the score of our page can have a little bonus



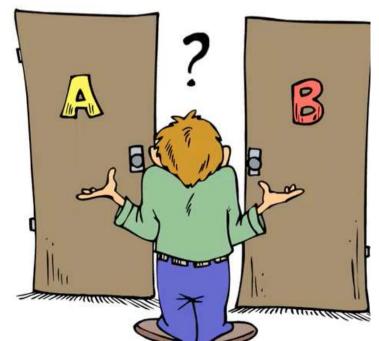
What is missing...?

- Beyond the techniques of textual spamdex, there is also another problem even more fundamental:
- What keywords to choose!



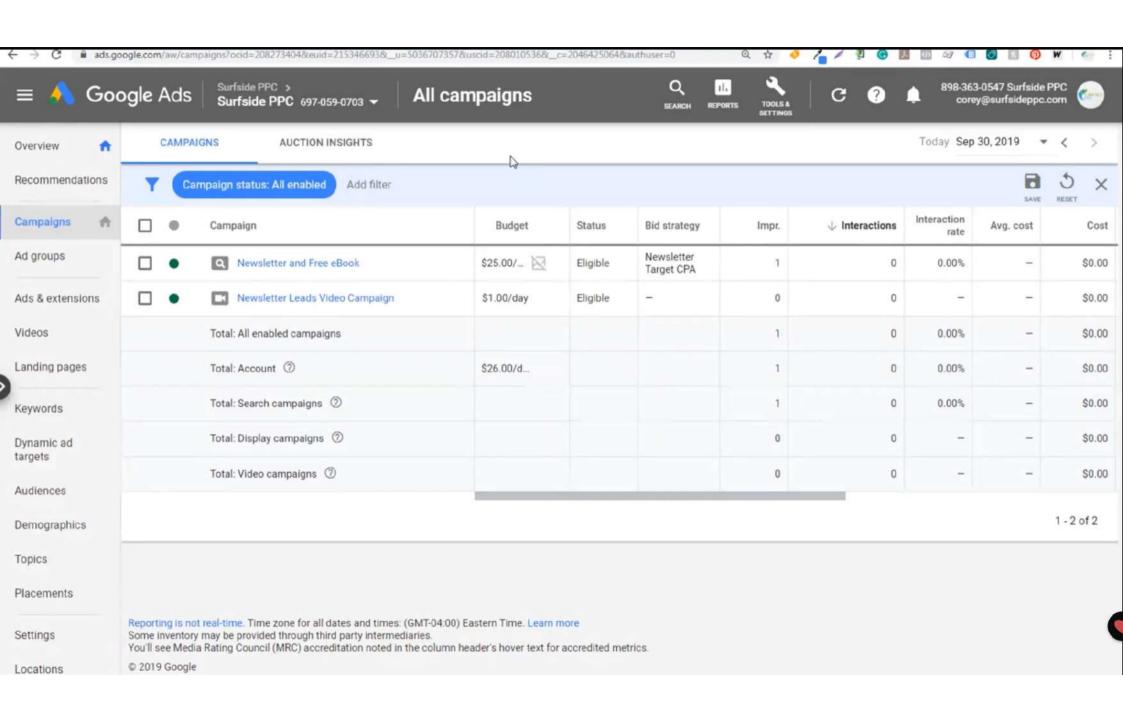
How can we do it?

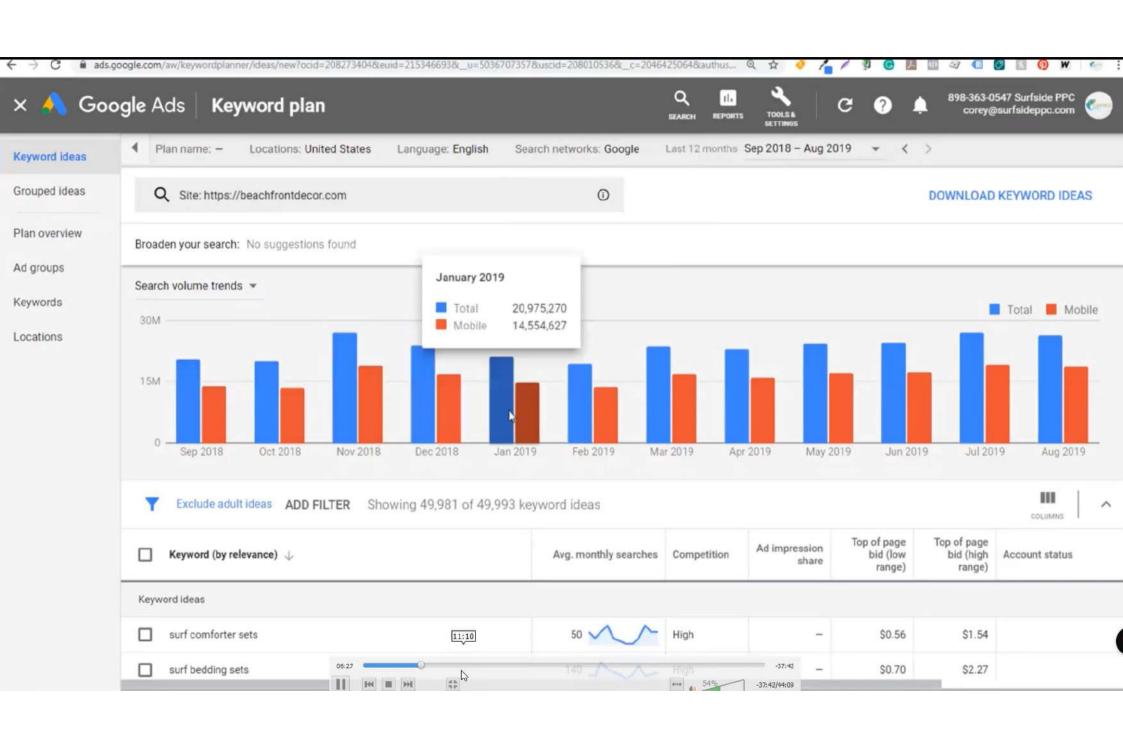
- We need to know what users want
- But how??
- There are various methods, let's see a couple...

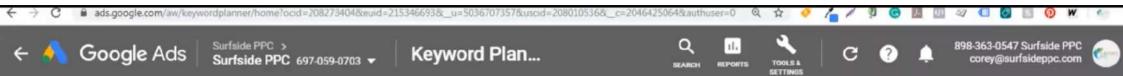


Smart methogs

Google Ads Keyword Planner (a.k.a. Keyword Tool *gk (\$\$...)

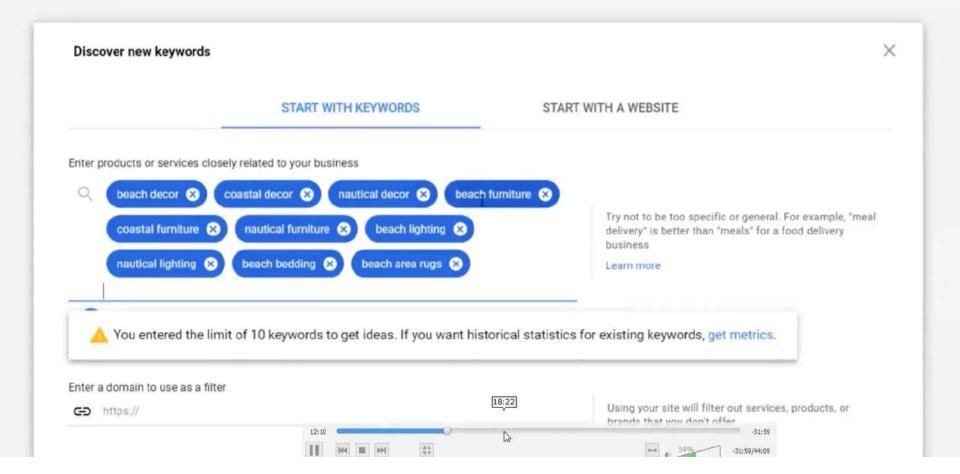


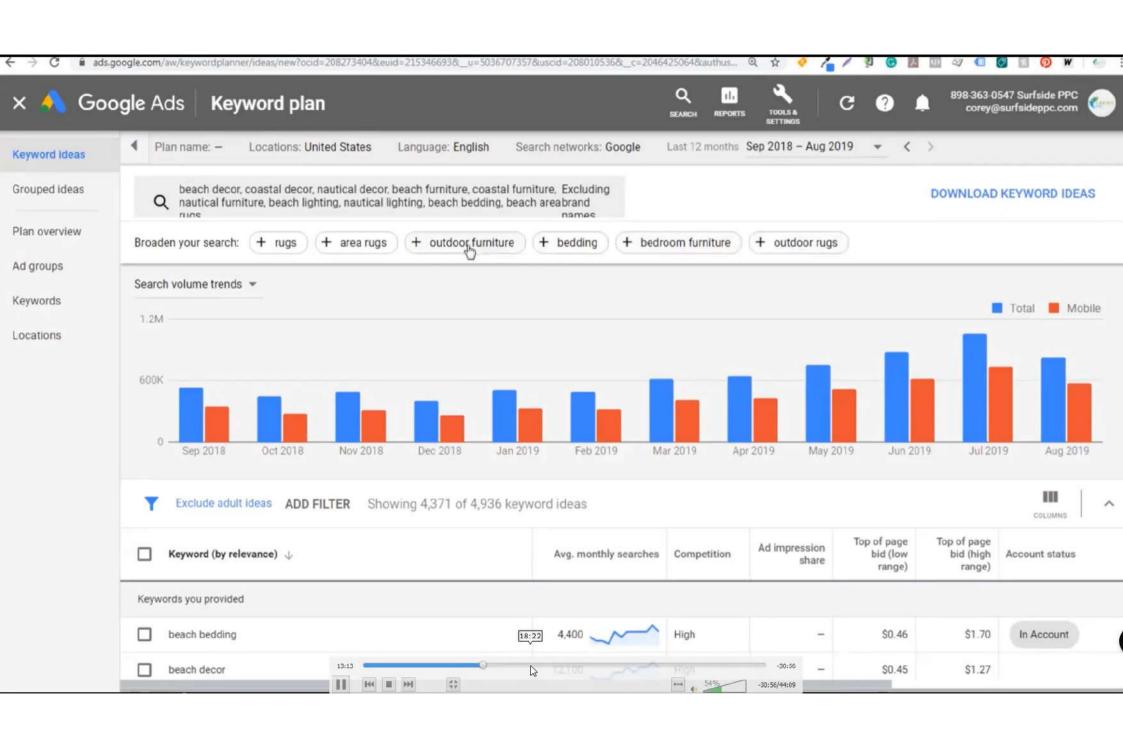




Select an active account

Surfside PPC 🧨





Smart methogs

- Google Ads Keyword Planner (a.k.a. Keyword Tool *gk (\$\$...)
- (likewise in Bing etc...)
- ♦ Or.... <u>alternatives</u>... (!)





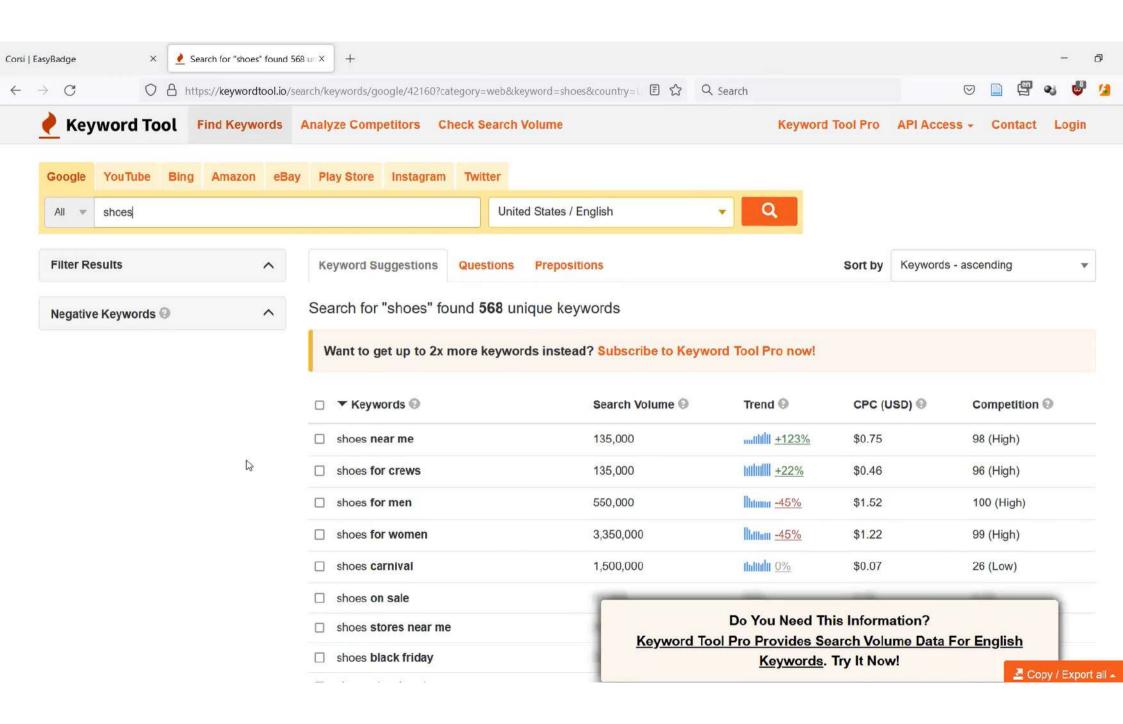
Find Great Keywords Using Google Autocomplete



B

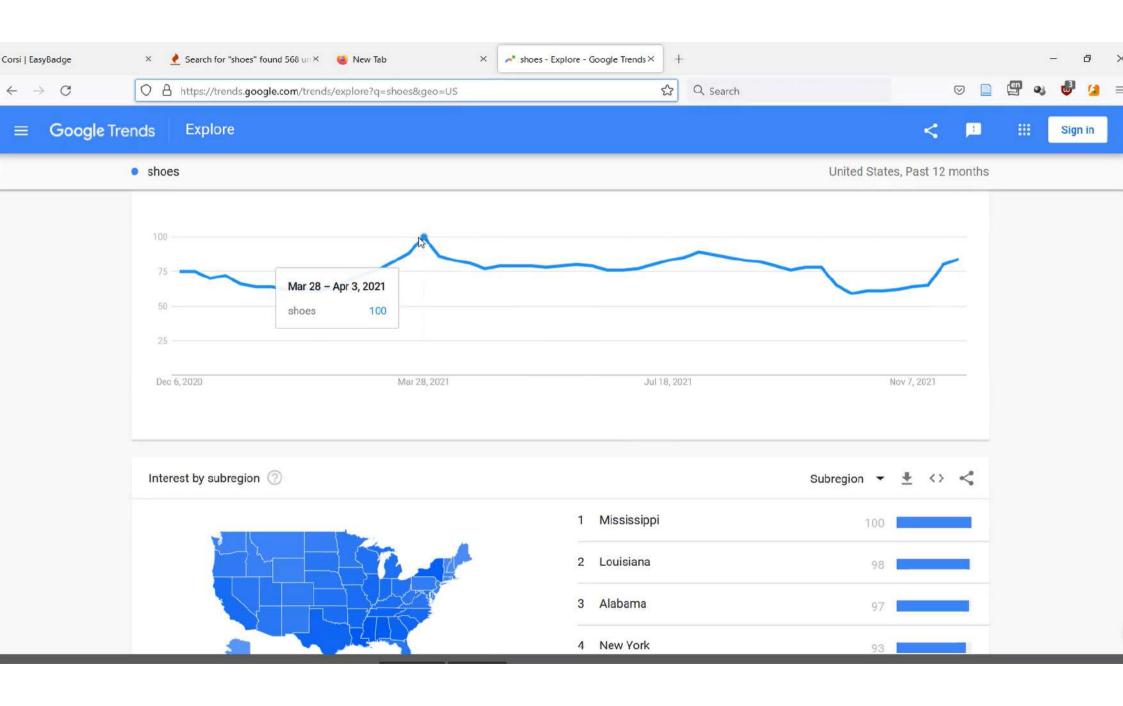
Keyword Tool Is The Best Alternative To Google Keyword Planner And Other Keyword Research Tools

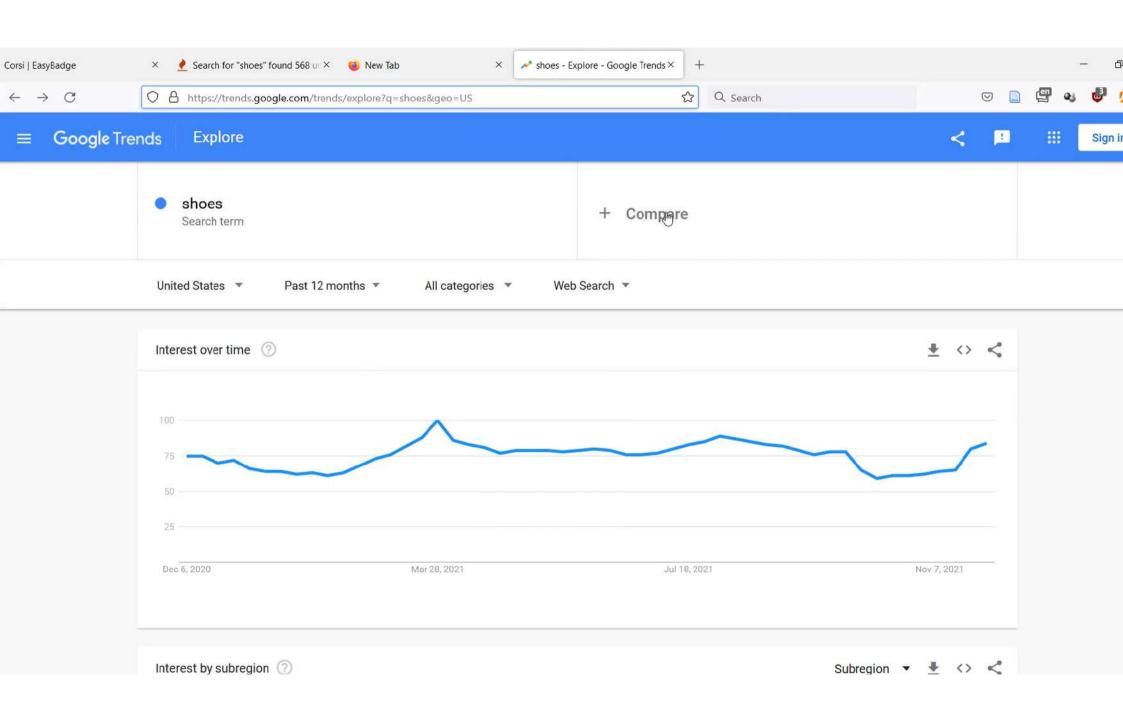
Here are a few reasons why:



Smart methogs

- Google Ads Keyword Planner (a.k.a. Keyword Tool *gk (\$\$...)
- (likewise in Bing etc...)
- ◆Or.... alternatives... (!)
- Together with some related analysis...





Problems of term spamming

- We said it at the very beginning when talking about body spam:
- Text spamming generally changes the content of the pages, so users are affected (likely in a bad way!!)

Think...

• We «power up» a page by inserting keywords, piece of other pages etc etc



◆The page gets higher in search engines
→ Users click and land on the modified
page → content corrupted or not so
relevant → (usability) he gets angry!!



Hiding

• We can use a series of techniques called *hiding*, which hide the «trash» that we inserted with spamming



Content hiding (examples)

```
<body background="white">
 <font color="white">
 text to disappear...
 </font>
 </body>
```



Content hiding (examples)

```
<a href="\pippo.html">
  <img src="\webbug.gif">
  </a>
```



Redirection...



Redirection

- Also called "302 technique
- Easy way:
- <meta http-equiv="refresh" content="0;url=pippo.html">
- Actually not so effective (countermeasures!)



Redirection (cont.)



- More effective: use javascript (!)
- Because search engines have a hard time to understand Javascript (and in many cases they just ignore it!)
- Example:

```
<script language="javascript"> <!- -
location.replace("pippo.html")</pre>
```

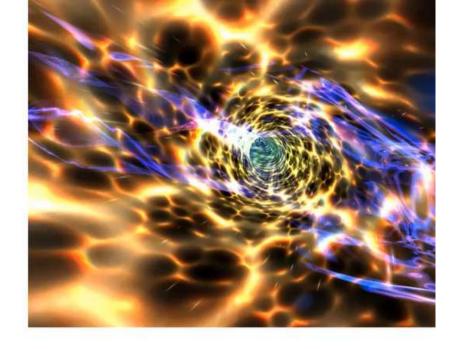
```
- -></script>
```

Cloaking...



Let's pass now...

... to the other component, the hypertextual one

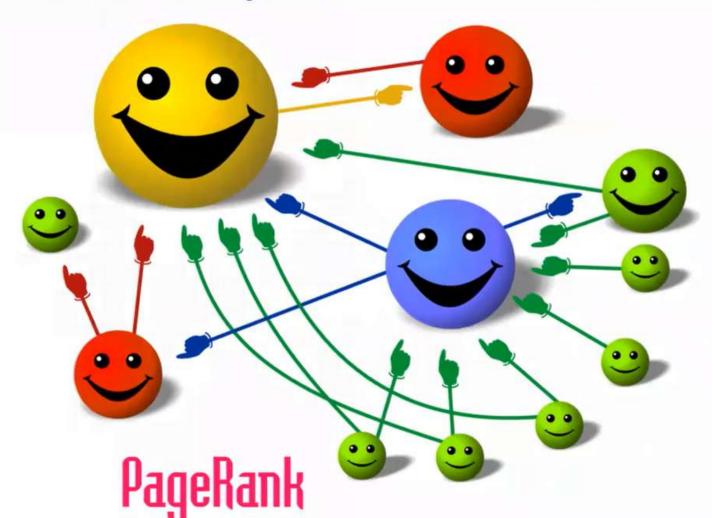


Complementary to the textual one, and contributes with a good deal of points derived from the *network topology*



The super-famous Pagerank

More or less you heard how it works...

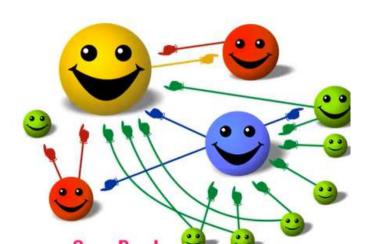


Pagerank in more detail...

(d_w = number of outgoing links)

$$\pi_v = \sum_{(w,v)\in E} \frac{\pi_w}{d_w}$$

$$\sum_{v} \pi_{v} = 1$$



Intuition

Like water reservoirs...

$$\pi_v = \sum_{(w,v)\in E} \frac{\pi_w}{d_w}$$

$$\sum_{v} \pi_{v} = 1$$



How to calculate it?

For instance, iteratively...



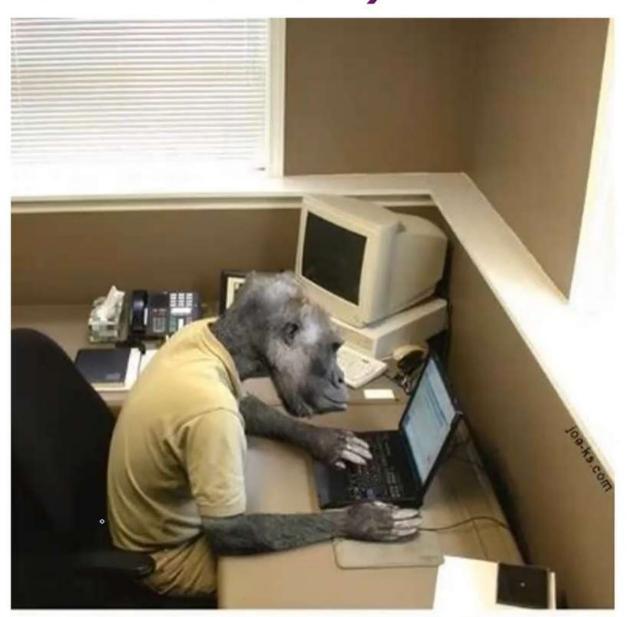
PageRank(p)
$$\leftarrow \Sigma_{q \in B(p)} (PageRank(q) / N(q))$$

Reformulation...

Markov chains and random walks (the «drunkard's walk»)

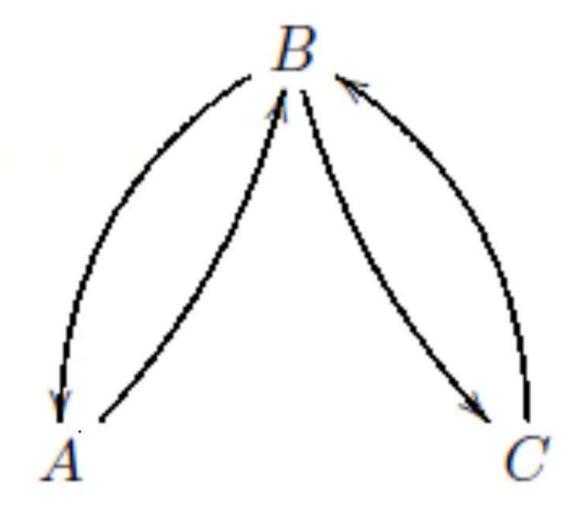


In the web drunkard → monkey ("random surfer")



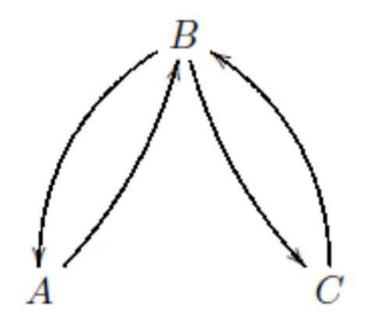
Example

What Pagerank?



Example (cont.)

Either we calculate, or we reason by symmetry:



- A=C, B has double flow than A and C, so (if the total liquid is 1):
- The pagerank of A and C is 1/4, those of B is 1/2.

All easy?

Let's see some other case...



Descent to Hell...

Instance of a more general problem in the web structure: spider traps

when a web spider gets «trapped» into an infinite navigation



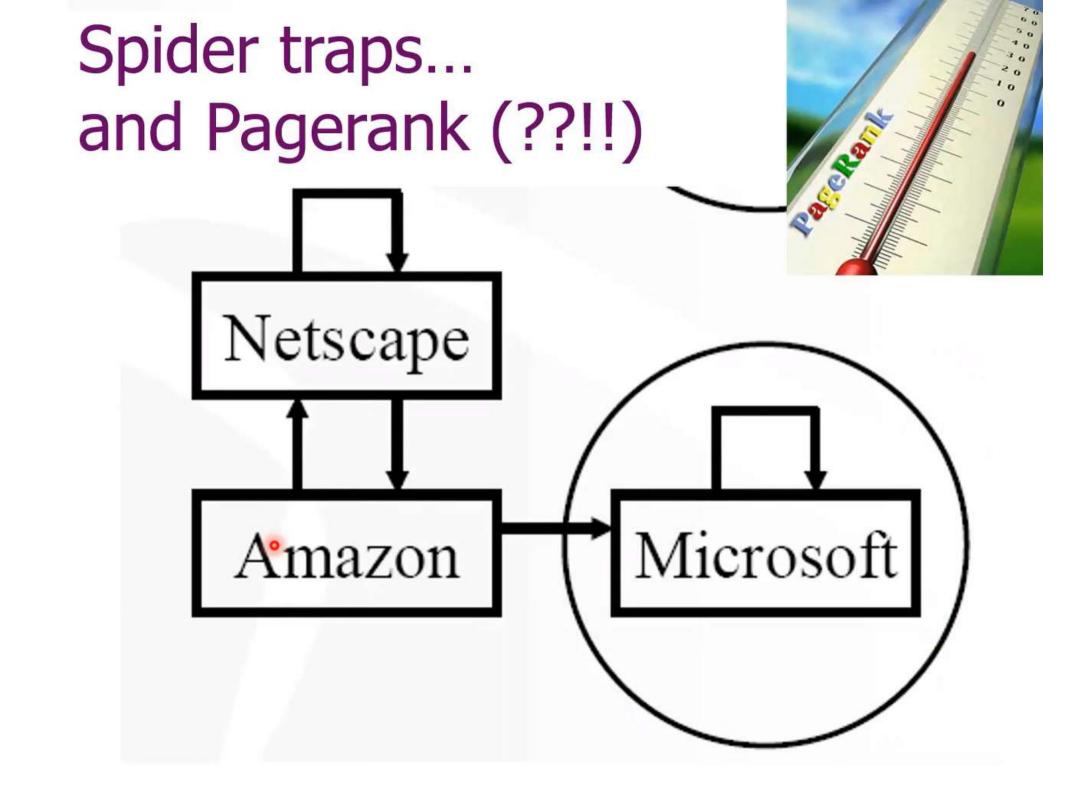


Spider traps...

- Unfortunately are very frequent
- For instance just think of an event calendar...

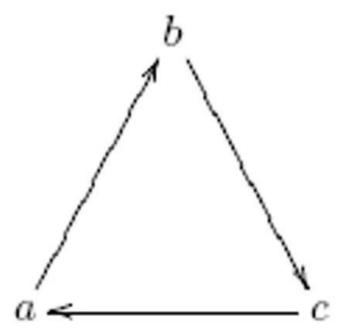


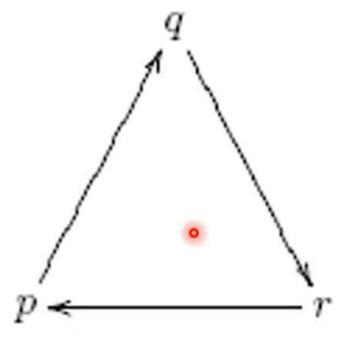




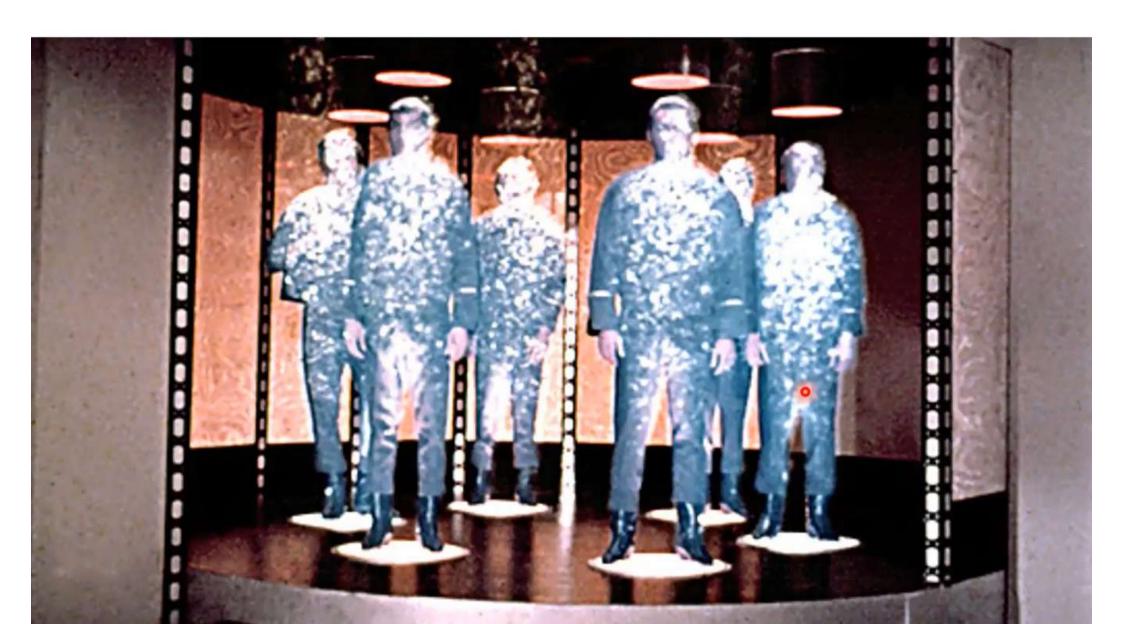
Other scenario: Islands

What pagerank?





Solutions....??



Let's change...!

The formula has been patched (teleportation)!

$$\pi_v = (1-\epsilon) \left(\sum_{(w,v) \in E} \frac{\pi_w}{d_w} \right) + \frac{\epsilon}{N}$$

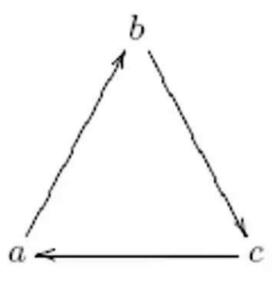


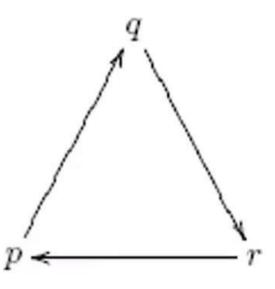


So...



◆After a few steps we get out of a spider trap, and even from an island ☺







Note...

Range of possible choices: teleport 0 = original pagerank, teleport 1 = score equal for everybody



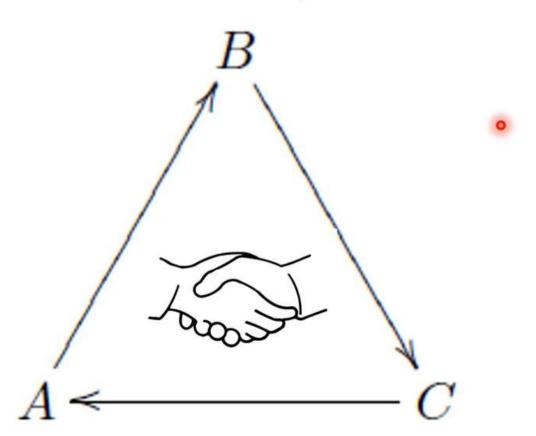
Note on elegance... Totalrank

Computational cost: almost the same as Pagerank

$$\mathbf{T} = \int_0^1 \mathbf{r}(\alpha) d\alpha$$

Ok, back to start

- Pagerank in this case?
- ◆Easy: A=B=C have 1/3 ©



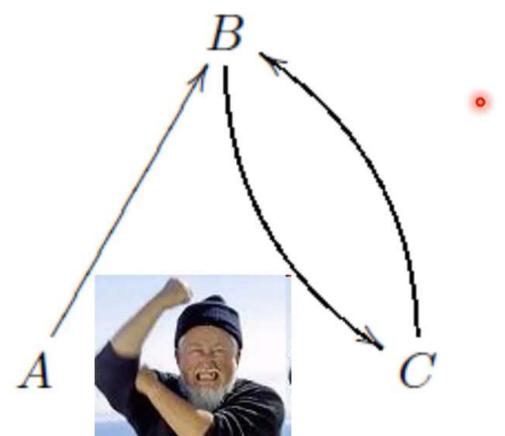


But, wait...!!

Here comes C, which changes a link!



• Pagerank of B and C = $\frac{1}{2}$, A = 0 (!!!)

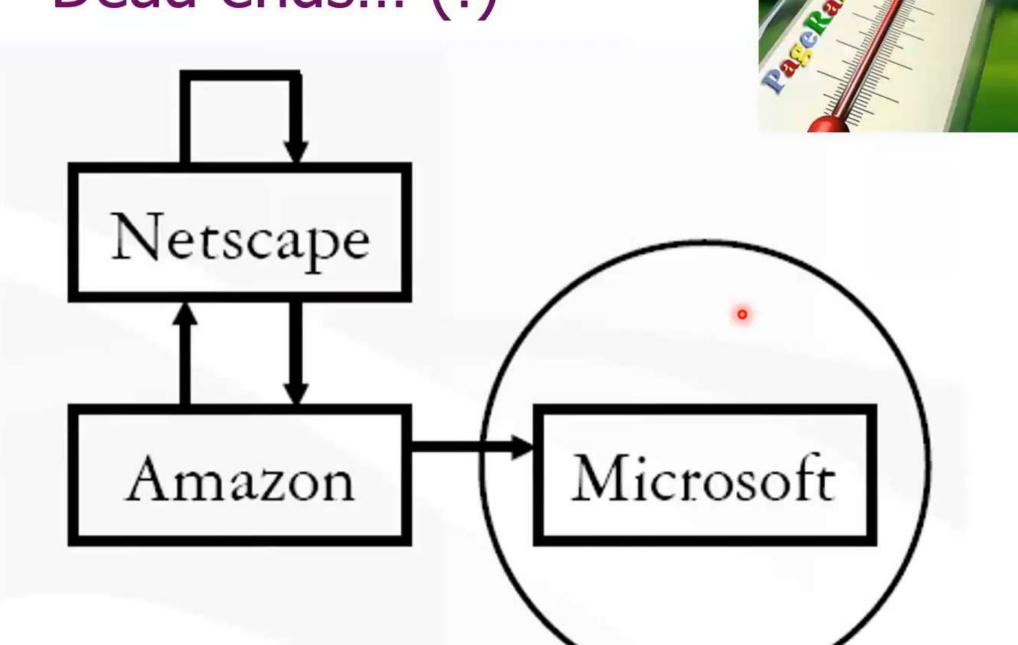


Symptom of a bigger problem...





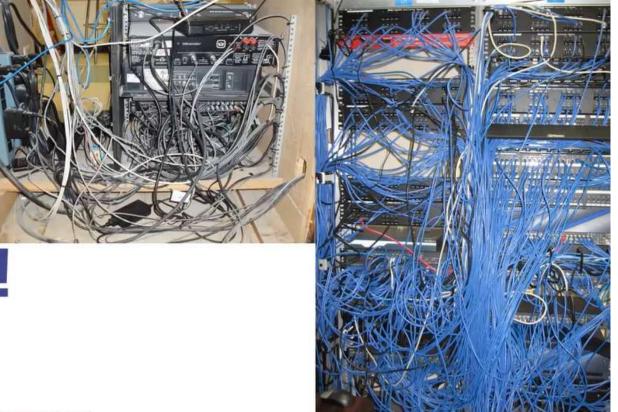
Dead ends... (!)

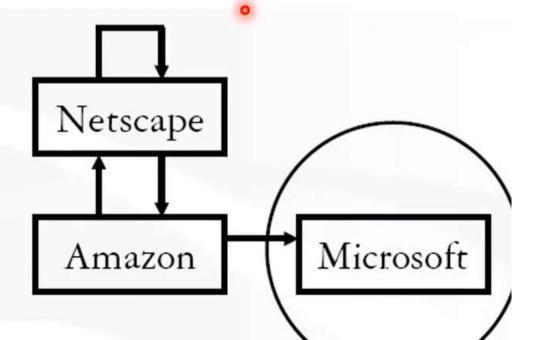


Soluzione?

A MESS!!







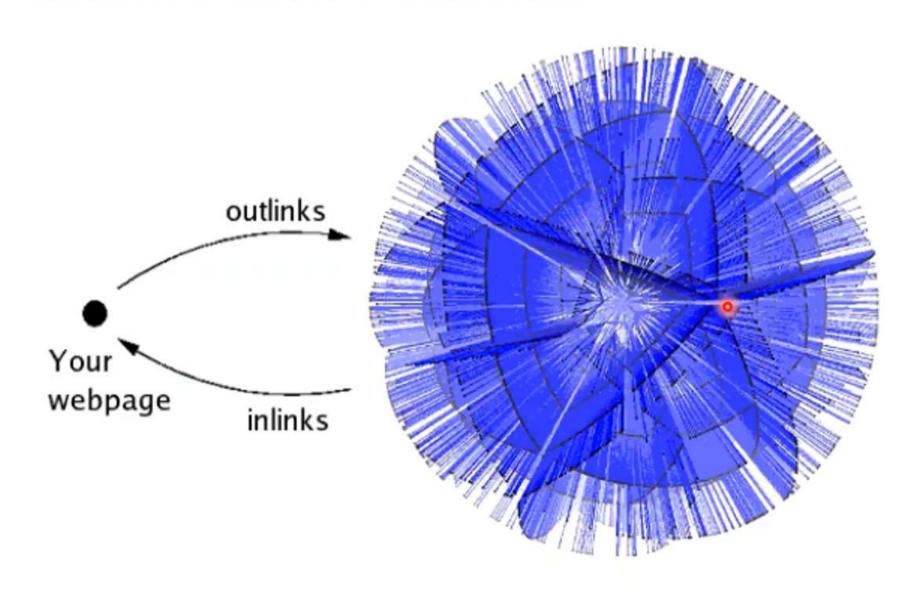
So...

• We have seen how Pagerank workds (...), now we can see how to use it for spamdex...

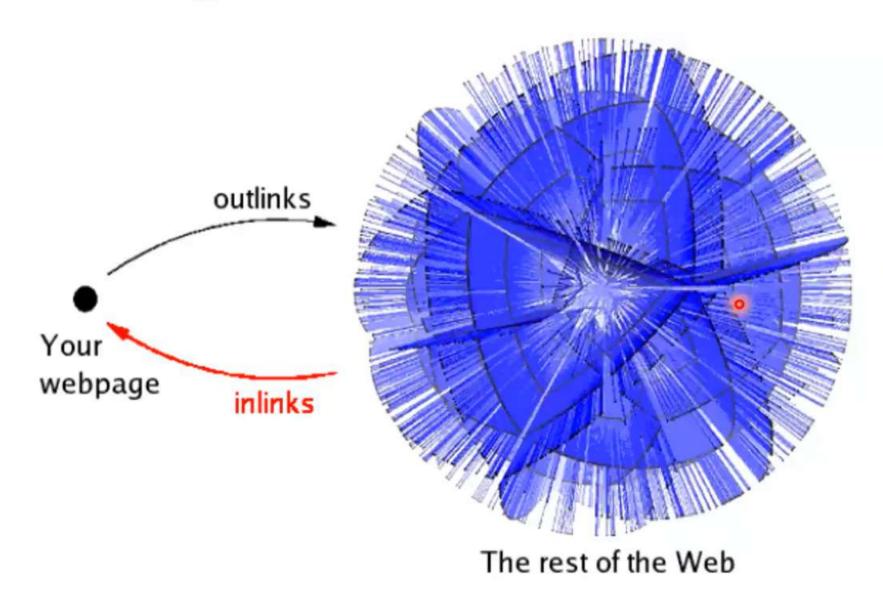




Two fundamental factors: Inlinks and outlinks



Strategic move Adding inlinks



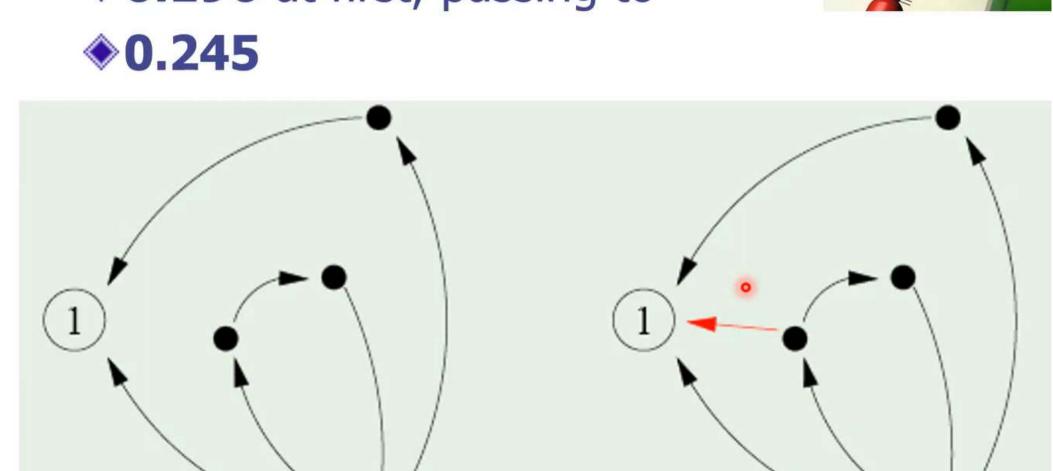
Result

Pagerank always gets higher



Example

◆0.196 at first, passing to



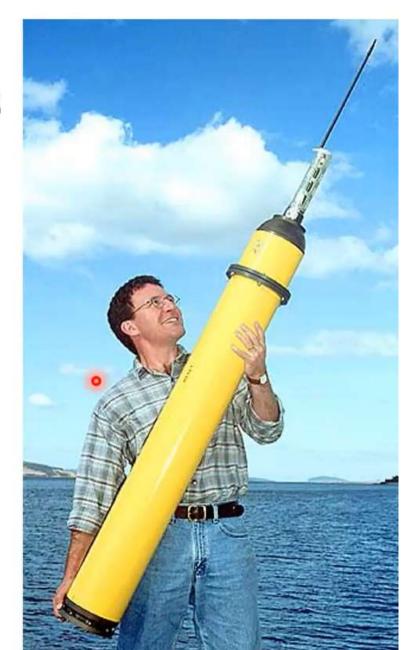


Techniques

- Infiltration
- Honey pot
- Link exchange
- Resurrection

Infiltration

- We «infilitrate» in various sites and try to insert links to our site
- For instance inside directories, blogs, wikis, comment sections etc



Honey pot

Yummy!!



Honey pot (cont.)

- Create "yummy" content, and then naturally receive incoming links
- We can do it by smart paste & copy from content of other sites (!)



Link Exchange

Let's join forces with someone else (related to another technique we will soon see)



Resurrection

When a domain «dies»... catch it!!

