

TrustRank: Combating the Web Spam

Mining of Massive Datasets
Leskovec, Rajaraman, and Ullman
Stanford University



Combating Spam

- **Combating term spam**
 - Analyze text using statistical methods
 - Similar to email spam filtering
 - Also useful: Detecting approximate duplicate pages
- **Combating link spam**
 - **Detection and blacklisting of structures that look like spam farms**
 - Leads to another war – hiding and detecting spam farms
 - **TrustRank** = topic-specific PageRank with a teleport set of **trusted pages**
 - **Example:** .edu domains, similar domains for non-US schools

TrustRank: Idea

- **Basic principle: Approximate isolation**
 - It is rare for a “good” page to point to a “bad” (spam) page
- Sample a set of **seed pages** from the web
- Have an **oracle (human)** to identify the good pages and the spam pages in the seed set
 - **Expensive task**, so we must make seed set as small as possible

Trust Propagation

- Call the subset of seed pages that are identified as **good** the **trusted pages**
- Perform a topic-sensitive PageRank with **teleport set = trusted pages**
 - **Propagate trust through links:**
 - Each page gets a trust value between **0** and **1**
- **Solution 1: Use a threshold value and mark all pages below the trust threshold as spam**

Why is it a good idea?

- **Trust attenuation:**

- The degree of trust conferred by a trusted page decreases with the distance in the graph

- **Trust splitting:**

- The larger the number of out-links from a page, the less scrutiny the page author gives each out-link
- Trust is **split** across out-links

Picking the Seed Set

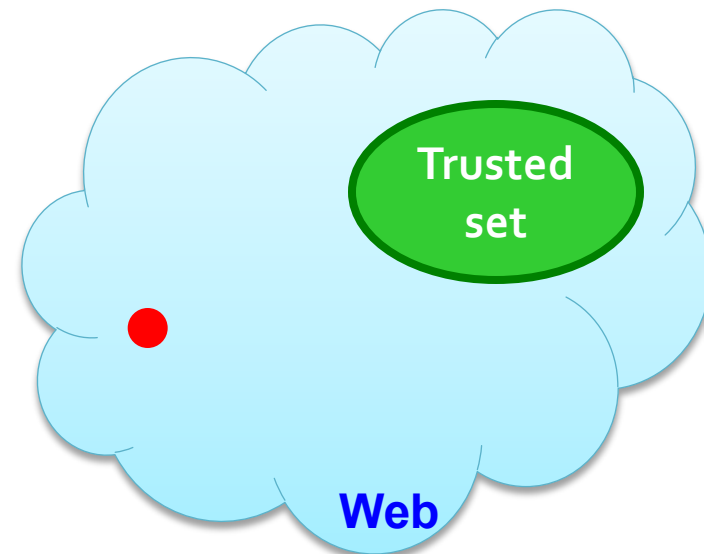
- **Two conflicting considerations:**
 - Human has to inspect each seed page, so **seed set must be as small as possible**
 - Must ensure every **good page** gets adequate trust rank, so need make **all good pages reachable from seed set by short paths**

Approaches to Picking Seed Set

- Suppose we want to pick a seed set of k pages
- **How to do that?**
- **(1) PageRank:**
 - Pick the top k pages by PageRank
 - The idea/hope is that you can't get a bad page's rank really really high
- **(2) Use trusted domains** whose membership is controlled, like .edu, .mil, .gov

Spam Mass

- In the **TrustRank** model, we start with good pages and propagate trust
- **Complementary view:**
What fraction of a page's PageRank comes from **spam** pages?
- In practice, we don't know all the spam pages, so we need to estimate



Spam Mass Estimation

Solution 2:

- r_p = PageRank of page p
- r_p^+ = PageRank of p with teleport into **trusted** pages only
- **Then:** What fraction of a page's PageRank comes from **spam** pages?

$$r_p^- = r_p - r_p^+$$

- **Spam mass of p** = $\frac{r_p^-}{r_p}$

