

Social Networks

Luigi Giugliano¹, Steven Rosario Sirchia¹

¹Università degli studi di Salerno

July 5, 2016

Introduction

The purpose of our work is to test different algorithms in the three fundamental areas for assembling any search engine offering a Sponsored Search system:

- **Ranking** of web documents
- **Matching** of words inside documents
- **Auctions** for acquiring advertisement slots.

We will briefly talk about the proposed algorithms, and then compare running times and results obtained from their execution more in detail, suggesting what combination of algorithms seems to be the best for realizing a new search engine.

Creating the Dataset

Our experiments ran on a set of approximately 30000 pages created this way:

- we choose a web-page for each of the 15 categories listed in <https://www.dmoz.org/>
- for every of these web pages we crawled 2000 pages by using the Wibbi online crawler
- Moreover, from each pair of sets of 2000 pages, we choose at random 10 pairs of vertices (u, v) with u being a page in the first set and v being a page in the second set and added a link from u to v (if this link was absent)

Creating the Dataset

Here is the complete list of the websites chosen.

Category	Website	Description
Arts	www.imdb.com	The Internet Movie Database
Business	www.moodyys.com	Corporate finance, banking
Computers	www.ibm.com	International Business Machines Corporation.
Games	www.ign.com	Videogame news
Health	www.who.int	World Health Organization
Home	www.cooks.com	Recipe search
Kids	www.cartoonnetwork.com	The home of cartoons online

Creating the Dataset

Category	Website	Description
News	www.foxnews.com	Breaking News
Recreation	www.lego.com	Producer of bilding blocks.
Reference	www.britannica.com	Encyclopaedia Britannica Online.
Science	www.nasa.gov	Comprehensive, world-class center for aeronautics
Shopping	www.amazon.com	Most know shopping website
Society	www.un.org	Daily United Nations news, documents and publications
Sports	www.nba.com	The official site of the National Basketball Association
Regional	www.lonelyplanet.com	Offers travel advice, detailed maps, travel news

Overview

1 Ranking

- Page Rank
 - Results
- HITS
 - Results
- Comparing the Results

2 Matching

- Best Match
- Improved Best Match
- Results

3 Search Engine

- Results

4 Auction

- First Price Auction
- Generalized Second Price Auction
- Results
 - Bots
 - Selecting Bots
 - Experiment on “real case”

5 Summing Up

Overview

- 1 Ranking
 - Page Rank
 - Results
 - HITS
 - Results
 - Comparing the Results

2 Matching

3 Search Engine

4 Auction

5 Summing UP

Page Rank

Page Rank

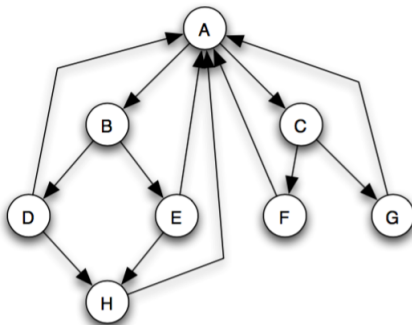
The intuition behind *Page Rank* is:

“a page is important if it is cited by other important pages”.

This intuition rises from the usual endorsement mode, for example, among academic or governmental pages, among bloggers, or among personal pages more generally. It is also the dominant mode in the scientific literature.

Algorithm

We can think of PageRank as a kind of “fluid” that circulates through the network, passing from node to node across edges, and pooling at the nodes that are the most important.



Algorithm Steps

- In a network with n nodes, we assign all nodes the same initial PageRank, set to be $1/n$.

Algorithm Steps

- In a network with n nodes, we assign all nodes the same initial PageRank, set to be $1/n$.
- We choose a number of steps k .

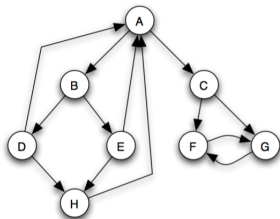
Algorithm Steps

- In a network with n nodes, we assign all nodes the same initial PageRank, set to be $1/n$.
- We choose a number of steps k .
- We then perform a sequence of k updates to the PageRank values, using the following rule for each update:

Basic PageRank Update Rule: Each page divides its current PageRank equally across its out-going links, and passes these equal shares to the pages it points to. (If a page has no out-going links, it passes all its current PageRank to itself.) Each page updates its new PageRank to be the sum of the shares it receives.

The “Wrong” nodes

There is a difficulty with the basic definition of PageRank, however: in many networks, the “wrong” nodes can end up with all the PageRank.



The Wrong nodes are a small sets of nodes that can be reached from the rest of the graph, but have no paths back.

Scaled PageRank

We can use the mechanism of fluid presented before: there is a **counter-balancing process** preventing that all the water stands only on downhill places on the earth.

Scaled PageRank Update Rule

First apply the Basic PageRank Update Rule.

Then scale down all PageRank values by a factor of s , shrinking the total from 1 to s .

We divide the residual $1 - s$ units of PageRank equally over all nodes, giving $(1 - s)/n$ to each.

Results

We are going to present the result of the experiment, comparing the **execution time** of PageRank on following inputs:

- Graph of 1000 nodes
- Graph of 2000 nodes
- Graph of 5000 nodes
- Graph of 10000 nodes
- Graph of 20000 nodes
- Full Graph (30000 nodes)

All the graphs are generated by chunking the Full Graph.

Graph of Times



Rank Values

We now present the result of the PageRank algorithm, considering the score of the pages:

- **Min:** $5e - 06$
- **Max:** $1.8e - 03$
- **Mean:** $1.1e - 05$
- **Std:** $5.4e - 05$

HITS

This hubs-and-authorities algorithm, sometimes called HITS (*hyperlink induced topic search*), was originally intended not as a preprocessing step before handling search queries, as PageRank is, but as a step to be done along with the processing of a search query, to rank only the responses to that query.

This kind of approach is used by the Ask search engine.

The Intuition Behind HITS

HITS views important pages as having two different types of importance.

- Certain pages are valuable because they provide information about a topic. These pages are called **authorities**.
- Other pages are valuable not because they provide information about any topic, but because they tell you where to go to find out about that topic. These pages are called **hubs**.

HITS Algorithm

For calculating the HITS values for the pages, we shall assign two scores to each Web page. One score represents the *hubbiness* of a page, that is the degree to which it is a good hub, and the second score represents the degree to which the page is a good authority.

These values are then calculated as:

- **Hubbiness**: the sum of the Authority value of the outgoing nodes.

HITS Algorithm

For calculating the HITS values for the pages, we shall assign two scores to each Web page. One score represents the *hubbiness* of a page, that is the degree to which it is a good hub, and the second score represents the degree to which the page is a good authority.

These values are then calculated as:

- **Hubbiness**: the sum of the Authority value of the outgoing nodes.
- **Authority**: the sum of Hubbiness value of the incoming nodes.

HITS Algorithm

For calculating the HITS values for the pages, we shall assign two scores to each Web page. One score represents the *hubbiness* of a page, that is the degree to which it is a good hub, and the second score represents the degree to which the page is a good authority.

These values are then calculated as:

- **Hubbiness**: the sum of the Authority value of the outgoing nodes.
- **Authority**: the sum of Hubbiness value of the incoming nodes.

These values are normalized so that the largest value is 1.

Results

We are going to present the result of the experiment, comparing the **execution time** of HITS on following inputs:

- Graph of 1000 nodes
- Graph of 2000 nodes
- Graph of 5000 nodes
- Graph of 10000 nodes
- Graph of 20000 nodes
- Full Graph (30000 nodes)

All the graphs are generated by chunking the Full Graph.

Tuning for improving performance

On the first attempts of running the algorithm on the full graph we observed that one iteration took about 30 minutes, due to the nature of the algorithm.

In each iteration we explore all the graph and calculate the incoming nodes for the current node ...

Considering that the graph never changes, we precomputed all the incoming nodes for each node so we can obtain the incoming nodes in $O(1)$.

Tuning for improving performance

In the HITS algorithm there are two stop rules:

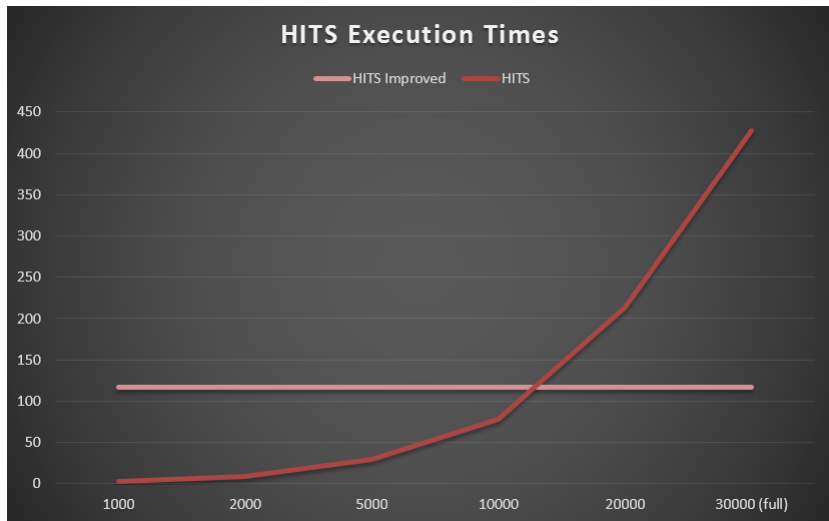
- Max number of iteration
- Min confidence on the errors reached

We noticed that in some cases the Algorithm runs until the maximum number of steps is reached, and when this happens the relative error trends to stabilize on a fixed level. We decided then to add a new stop rule:

- If two successive relative errors are distant at maximum a small ϵ , then stop the algorithm.

We tried this improvement on the full graph, we plotted the time as a pink line ... the difference is embarrassing

Graph of Times



Parallel Version

We also implemented the parallel version, but the performance resulted to be worst compared with the non-parallel version: this is due to overhead of coping all the structures needed for the calculation of scores.

We now present times of 10 iterations of both version:

- **Non-Parallel:** 4,15"
- **Parallel:** 53,41"

Hits Authority Values

We now present the result of the HITS algorithm, considering the authority score of the pages:

- **Min:** 0
- **Max:** 1
- **Mean:** 0.0017
- **Std:** 0.0408

Hits Hubbiness Values

We now present the result of the HITS algorithm, considering the hubbiness score of the pages:

- **Min:** 0
- **Max:** 1
- **Mean:** 0.0585
- **Std:** 0.2328

Comparing the Results

The results are incomparable, both considering execution time and values:

HITS time is up to 40 times the PageRank time.

The values are different for the nature of the algorithms.

But we can combine the observations we made earlier to infer the structure of the network, in fact both PageRank and HITS authority have the majority of documents with a low score, and also the variance proves this.

For the HITS hubbiness we notice that the variance is higher, as well as the mean, this means that the values are more widespread along the network.

Network structure



Network structure



Overview

- 1 Ranking
- 2 Matching
 - Best Match
 - Improved Best Match
 - Results
- 3 Search Engine
- 4 Auction
- 5 Summing UP

The idea of Best Match

Given a query q , containing n query words, and a set of documents S , we define Best Match as a method that finds a subset of document S' such that:

- each document s_i of S' has a "reasonable" number of query words in it

According to this definition the basic Best Match consists of:

- counting how many query words documents have
 - we call this value "score" of a document
 - it is at maximum n
- ordering in decreasing order of score the documents (optional)
- return all documents whose score is "reasonable"
 - we use a threshold to define what is "reasonable"

Refining the Best Match

Two are the basic refinements to have a more efficient Best Match:

- 1 using an inverted index
 - in the form (word \rightarrow list of documents containing the word)
 - the keys of the dataset are the query words
 - we can have in $O(1)$ all the documents with a determined word
- 2 using the frequency instead of assigning score 1 to each query word found
 - defined as number of occurrences in document d / $\text{length}(d)$
 - requires precalculation of occurrences for all words and all S
 - it represents the **relevance** of documents to a particular word or query

Improved Best Match

We implemented also the following improved version of Best Match:

- 1 Sort documents in each inverted index in order of frequency of the term at which the inverted index refers
- 2 For every query term define its possible impact on the score as the frequency of the most frequent document in its index
- 3 Sort the query terms in decreasing order of impact
- 4 Consider the first 20 documents in the index of the first query term (if the first query term has an index with less than 20 documents, then complete with the first documents in the index of the next query term)

Improved Best Match

- 5 Compute the score for each of these documents
- 6 Consider the first term in which there are documents that have not been scored
- 7 Consider the first non-scored document in the index of this term
- 8 If the frequency of the current term in the current document plus the sum of the impact of next terms is larger than the score of the 20-th scored document, then score this document and repeat from 7, otherwise consider the next

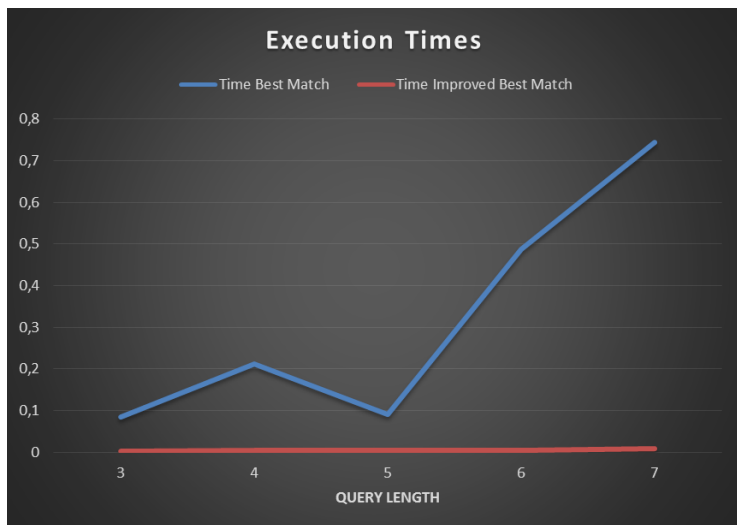
Experiment configuration

For obtaining the time comparison between BestMatch and ImprovedBestMatch we run both algorithms on 25 random queries of different length :

- 5 for each length from 3 to 7

Then we mean the results and plotted them on a graph

Time comparison BestMatch vs ImprovedBestMatch

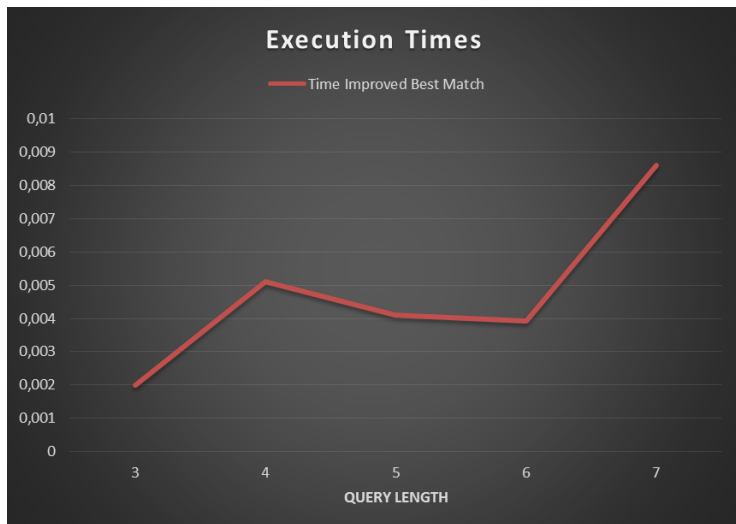


Time comparison BestMatch vs ImprovedBestMatch

First of all Improved Best Match “**wins**” on all query lengths, moreover it seems to be constant. This is due to the fact that there is a huge difference on the times of the algorithm at least of 1 order of magnitude.

We can notice that times of query of length 5 in best match are smaller of the one on length 4; this can be attributable to the random generation of queries, in fact, if we take 5 very uncommon words there are **few** documents that have one of these word so the algorithm ends quickly.

Improved Best Match time



Overview

- 1 Ranking
- 2 Matching
- 3 Search Engine**
 - Results
- 4 Auction
- 5 Summing UP

Search Engine

The implemented Search Engine combine the algorithms seen before.

Given a query:

- 1 Find the documents that match the query using:
 - Best Match
 - Improved Best Match
- 2 Order these documents by:
 - Page Rank Score
 - HITS authority

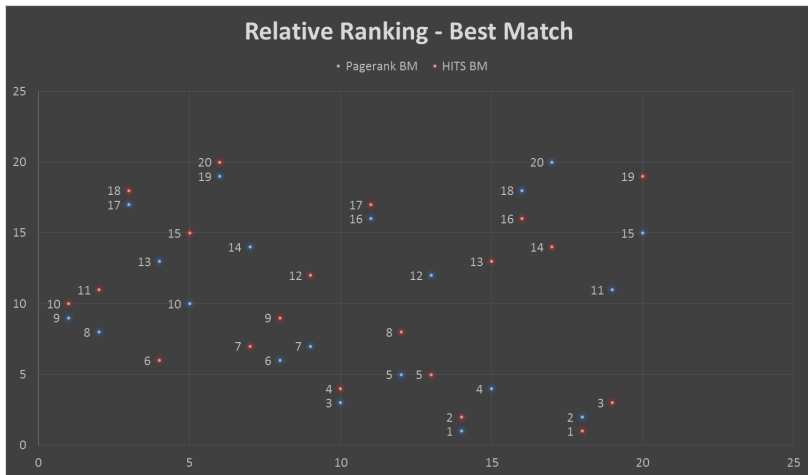
Running configurations

- We ran the experiment for a total of 25 queries
- 5 for each length between 3 and 7
 - In each group of 5 queries 1 of them is manually constructed by us for having a high likelihood with a real query
 - The other 4 are randomly generated among the dataset
- The manual generated query are used not only for calculate the execution time, but also for fully analysing the output of search engine

Ranking on Improved Best Match

	MATCHING DOCUMENTS	SCORE	BM	IBM	Pagerank BM	HITS BM	Pagerank IBM	HITS IBM
1	http://www.nasa.gov/audience/forstudents/k-4/stories/what-is-the-iss-k4.html	0,162393162	Y	Y	9	10	9	11
2	http://www.nasa.gov/audience/forstudents/5-8/features/what-is-the-iss-58.html	0,132013201	Y	Y	8	11	8	12
3	http://www.nasa.gov/audience/foreducators/expeditions/stem/stem-math-index.html	0,131578947	Y	Y	17	18	17	19
4	http://www.lego.com/en-us/juniors/games/firetruck	0,125	Y	Y	13	6	13	7
5	http://www.un.org/en/documents/contact.asp	0,125	Y	Y	10	15	10	16
6	http://www.nasa.gov/audience/foreducators/kidsclub/text/Phase_6_Putting_it_All_Together.htm	0,124645892	Y	Y	19	20	19	20
7	http://www.lego.com/en-us/juniors/games/gas-station	0,12244898	Y	Y	14	7	14	8
8	http://www.lego.com/en-us/juniors/games/gas-station?icmp=COUSGamesGJuniorsGasSta	0,12244898	Y	Y	6	9	6	10
9	http://www.nasa.gov/audience/forstudents/k-4/stories/what-is-orion-k4.html	0,109195402	Y	Y	7	12	7	13
10	http://www.nasa.gov/audience/foreducators/teachingfromspace/dayinthelife/index.html	0,101694915	Y	Y	3	4	3	5
11	http://www.nasa.gov/audience/foreducators/expeditions/stem/stem-tech-index.html	0,101694915	Y	Y	16	17	16	18
12	http://www.lego.com/en-us/juniors/games/pony?icmp=COUSGamesGJuniorsPony	0,1	Y	Y	5	8	5	9
13	http://www.lego.com/en-us/juniors/games/pony	0,1	Y	Y	12	5	12	6
14	http://www.lego.com/en-us/city/codepage	0,1	Y	Y	1	2	1	2
15	http://www.nasa.gov/audience/forstudents/k-4/stories/what-is-juno-k4.html	0,098765432	Y	Y	4	13	4	14
16	http://www.nasa.gov/mission_pages/station/research/index.html	0,097560976	Y	Y	18	16	18	17
17	http://www.nasa.gov/about/highlights/AN_Structure_OtherAgencies.html	0,095238095	Y	Y	20	14	20	15
18	http://www.ign.com/wikis/halo-master-chief-collection	0,091922006	Y	Y	2	1	2	1
19	http://www.ign.com/wikis/halo-master-chief-collection?save=successful	0,091922006	Y	Y	11	3	11	4
20	http://www.nasa.gov/audience/foreducators/spacelife/explorationdesign/overview/	0,090909091	Y	N	15	19	//////////	//////////
21	http://www.ign.com/wikis/halo-master-chief-collection/Halo_2	0,090909091	N	Y	//////////	//////////	15	3

Ranking on Best Match



Ranking on Improved Best Match



Considerations

First of all we can observe that in both types of matching we found the document we used for generating the query, as expected.

Another thing we notice is that the two results differ only by one document, but these documents have the same matching score, so this difference is due to the order of visiting documents.

We notice that there is an **high** correlation within the first and the last position of the documents ordered by PageRank or by HITS, this means that even if the values of the ranking algorithms are not comparable absolutely, tent to have a **relative** comparability.

Overview

- 1 Ranking
- 2 Matching
- 3 Search Engine
- 4 Auction**
 - First Price Auction
 - Generalized Second Price Auction
 - Results
 - Bots
 - Selecting Bots
 - Experiment on “real case”
- 5 Summing UP

First Price Auction

In this kind of auction, bidders submit simultaneous “sealed bids” to the seller. The terminology comes from the original format for such auctions, in which bids were written down and provided in sealed envelopes to the seller, who would then open them all together.

**The highest bidder wins the object
and pays the value of her bid.**

Non - truthfulness of FPA

In a sealed-bid first-price auction, the value of your bid not only affects whether you win but also how much you pay.

Bidding your true value is not a dominant strategy. By bidding your true value, you would get a payoff of 0 if you lose (as usual), and you would also get a payoff of 0 if you win, since you'd pay exactly what it was worth to you.

As a result, the optimal way to bid in a first-price auction is to “shade” your bid slightly downward, so that if you win you will get a positive payoff. Determining how much to shade your bid involves balancing a trade-off between two opposing forces.

Generalized Second Price Auction

Also called Vickrey auctions. Bidders submit simultaneous sealed bids to the sellers;

**The highest bidder wins the object
and pays the value of the second-highest bid.**

These auctions are called Vickrey auctions in honor of William Vickrey, who wrote the first game-theoretic analysis of auctions. Vickrey won the Nobel Memorial Prize in Economics in 1996 for this body of work.

Truthfulness of GSP

Truthful bidding is a dominant strategy in a sealed-bid second-price auction. The heart of the argument is the fact noted at the outset: in a second-price auction, your bid determines whether you win or lose, but not how much you pay in the event that you win.

So in a second-price auction, it makes sense to bid your true value even if other bidders are overbidding, underbidding, colluding, or behaving in other unpredictable ways.

Description of Bots

We used the following bots:

- **Best_response** with balanced tie-breaking rules
- **Best_response_competitive** submit the highest possible bid that gives the preferred_slot
- **Best_response_altruistic** submit the lowest possible bid that gives the preferred_slot
- **Competitor** always submit a bit grater than the highest bid
- **Budget_saving** always submit the minimum between last-non winning bid and advertiser value

Description of Bots

We used the following bots:

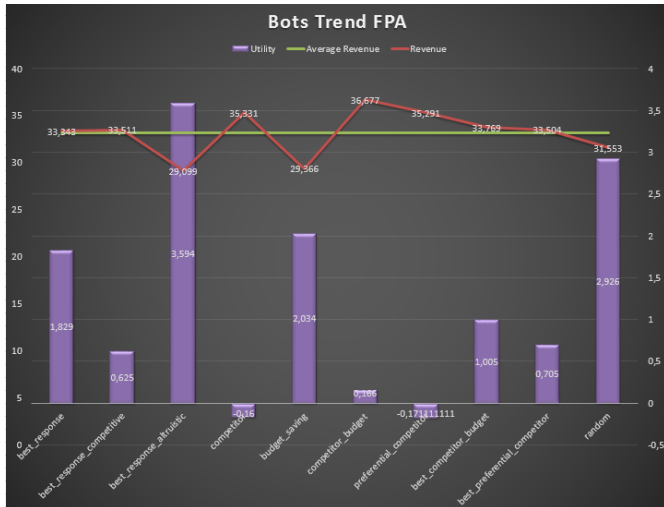
- **Competitor_budget** Use competitor when budget is more than its half, best_response otherwise
- **Preferential_competitor** Use competitor when value is more than a threshold, budget_saving otherwise
- **Best_competitor_budget** Use best_response_competitive when budget is more than its half, best_response otherwise
- **Best_preferential_competitor** Use best_response_competitive when value is more than a threshold, budget_saving otherwise
- **Random** bids randomly

Selection of interesting Bots

The configuration of the experiment for finding the most interesting bots is the following:

- **Number of Query Words** 1
- **Number of Auction** 10
- **Slot 1 click-through** 0.4
- **Slot 2 click-through** 0.15
- **Number of Advertiser** 3
 - 1 bot to test and 2 enemies (of the same kind)
- **Values** 7
- **Budgets** 25
- **Number of runs** 10000

FPA Utility



Considerations on FPA advertiser-side

Observing FPA graph it seems that *best_response_altruistic* is the best bot, yielding the **highest** utility.

This is due to the fact the bots that try to save their budgets, offer the lowest bid possible. Since these bids are “far” from the advertiser’s values and FPA is not truthful, they have a high chance to have a good slot at low price.

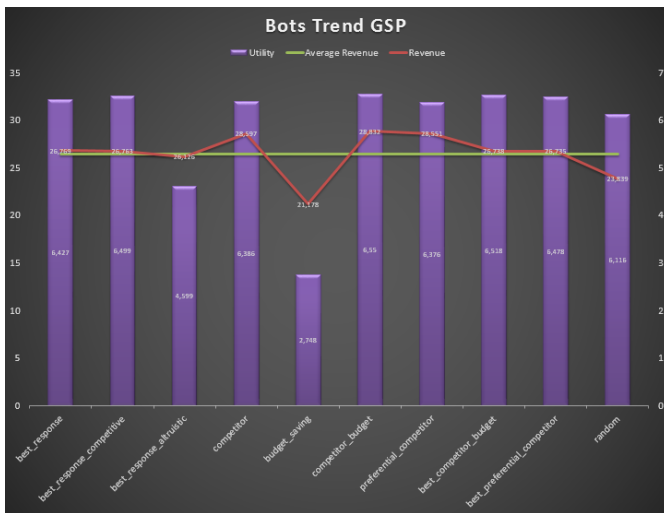
Considerations on FPA seller-side

The bots who try to save their budget on seller-side are the **worst**, because they lower not only their own bids, but also the bids of the whole auction.

It follows that the other kinds of bots which raise their bids are more **lucrative** for us.

For the reasons explained before we need a balance between high revenue and high utility, and the bot that seems to accomplish this is **best_response**.

GPS Utility



Considerations on GSP advertiser-side

We notice that, since GPS is truthful, the bots that offer bids near to their “real” value perform **better**, obtaining a higher utility.

For this reason all the bots that tent to lower the bids, straying from the “real” value, perform a lot **worst**, obtaining a lower utility.

Considerations on GSP seller-side

We notice that there are a lot of bots with almost the same utility, around $\sim 6,5$.

Since we know the bots have different behaviours, we would like to understand how this difference affect the outcome of the bots.

For this reason we decided to calculate an heuristic about the quantity and quality of slots obtained.

Our heuristic

The heuristic for a bidder is calculated as follow:

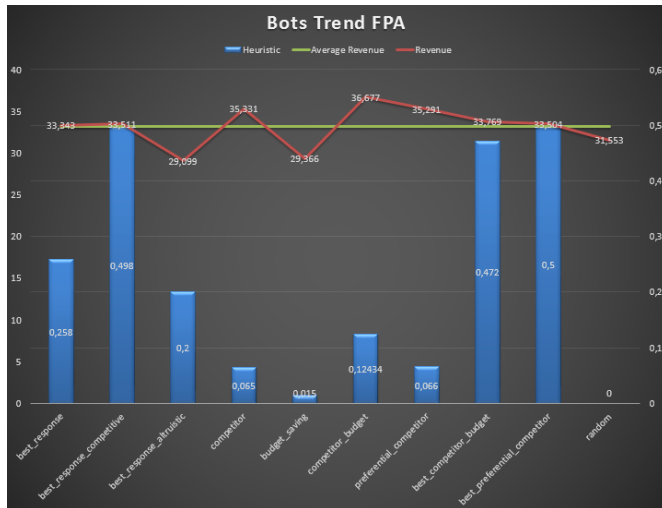
On each step of an auction, given the **preferred_slot** and the result of that step.

- Add 1 if the bidder obtains the **preferred_slot**, or he wants nothing and obtain nothing
- Add 2 if the bidder obtains a better slot, since he bids for a worse slot.
- Subtract 1 if the bidder obtain a worse slot, since he bids for a better slot.

We mediate this score on the auction step.

The heuristic straddles in the range $[-1; 2]$

Heuristic on FPA



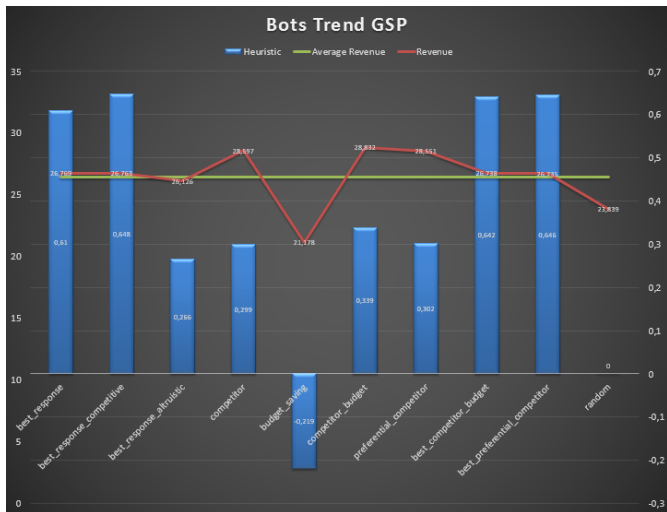
Considerations on FPA heuristic

Considering the bots who have the higher utility (best_response_altruistic, budget_saving) we notice that they have a very **low** heuristic score, this is due to the fact they tent the obtain the slots that the don't "want" but a very **low price** resulting in a high utility.

On the other hand aggressive bidders, tent to have a better heuristic score but a worse utility, since they pay a lot for the slot they obtain.

For the FPA the revenue is not dependant on the heuristic score, but on the "nature" of the bot.

Heuristic on GSP



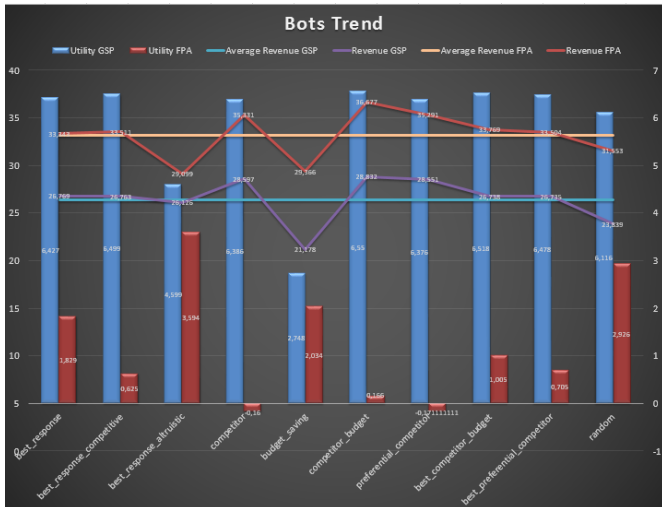
Considerations on GSP heuristic

Thanks to the heuristic score we can divide the bots in several groups, the “blind” competitors, the “wise” competitors and the “stingy”.

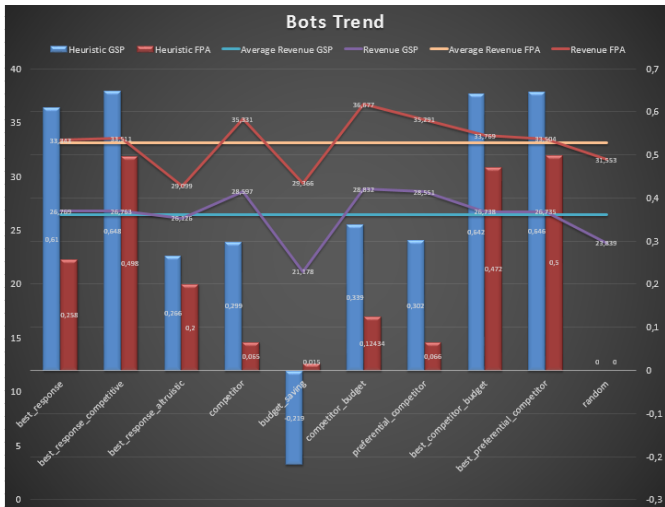
	Utility	H_score	Revenue
"Blind" competitor	High	Low	High
"Wise" competitor	High	High	Average
"Stingy"	Low	Low	Low

Table : Resume

Comparing FPA and GPS results



Comparing FPA and GPS results



Comparing FPA and GPS results

Comparing the auction types, we notice that on a “fixed” auction they have the following behaviour:

- FPA auction produce the **highest** revenue for the seller, this follow from the fact that you pay exactly what you bid. As for clients satisfaction, we notice that they have bot at **low** heuristic and utility compared to GPS.
- GPS instead produce **lowest** revenue and **high** utility and heuristic score for the client, this is due the nature of GPS:
the winner of a slot can offer a very high bid without having too much impact on what he pays, because he doesn't pay his bid but the bid immediately lower than his.

The choice depends of what we want to “favour”

Revenue equivalence Theorem

As we notice the revenue are different, this is due to the fact that we stopped the auction at 10 steps and moreover the budget of the advertisers are different so using balance algorithms, sometime they end up with not enough budget for bidding.

Bug setting the max step to: 100, and giving an higher budget equal to all the advertiser we have this result:

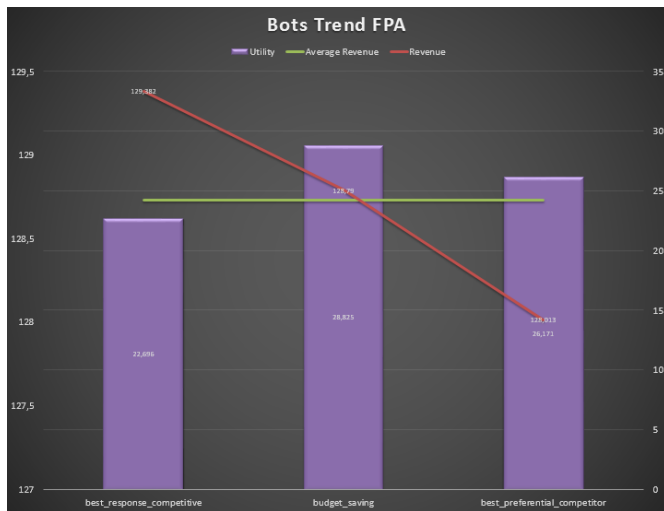
```
55
{'prova': {'adv_slots': {'id2': 'z', 'id1': 'y'}, 'adv_bids': {'prova': {'y': 7.0, 'x': 7.0, 'z': 7.0}}, 'adv_pays': {'y': 7.0, 'z': 7.0}}}
Totale Val:      0.672727272727
Totale Uti:      9.05714285714
Totale Rev:      756.4
AUCTION: fpa l: best_response ENEMIES: best_response
55
{'prova': {'adv_slots': {'id2': 'z', 'id1': 'y'}, 'adv_bids': {'prova': {'y': 7.0, 'x': 7.0, 'z': 7.0}}, 'adv_pays': {'y': 7.0, 'z': 7.0}}}
Totale Val:      0.672727272727
Totale Uti:      9.05714285714
Totale Rev:      756.4
```

Running Configuration

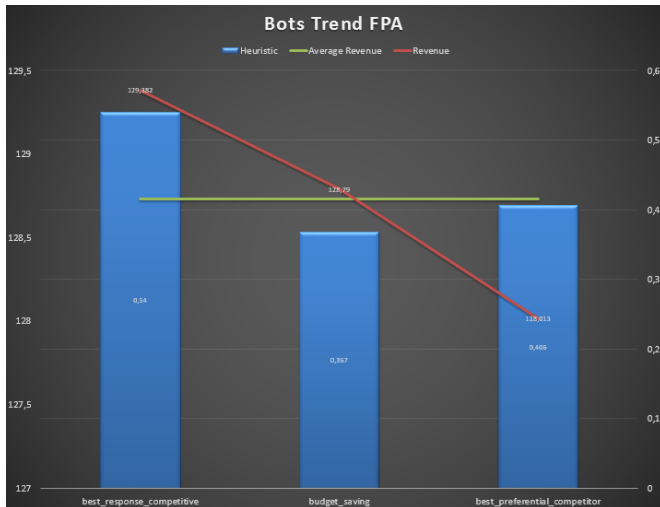
- Number of Query Words 10
- Number of Auction 20
- Number of Slot for each word [2, 4]
- Slot i click-through [0,1]
- Number of Advertiser 6
 - 1 bot to test and 5 enemies (of the same kind)
- Values [0, 20]
- Budgets [10, 50]
- Minimum Interest Threshold [5, 15]
- Number of runs 500

We choose two bots that have the best result among all the auction, and the worst one.

Utility FPA



Heuristic FPA

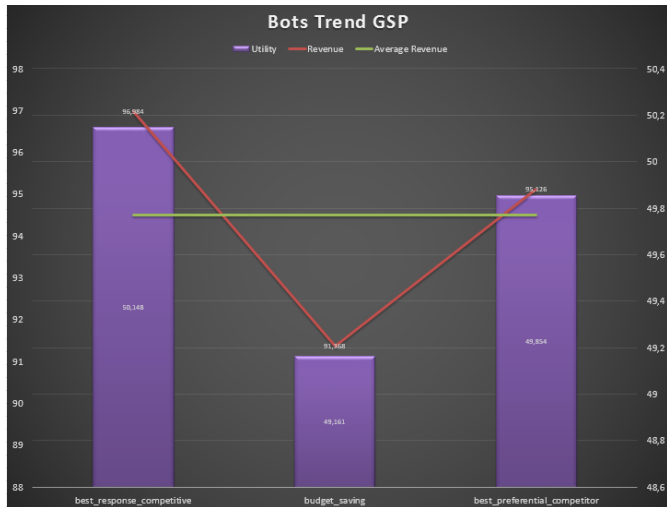


Consideration on FPA

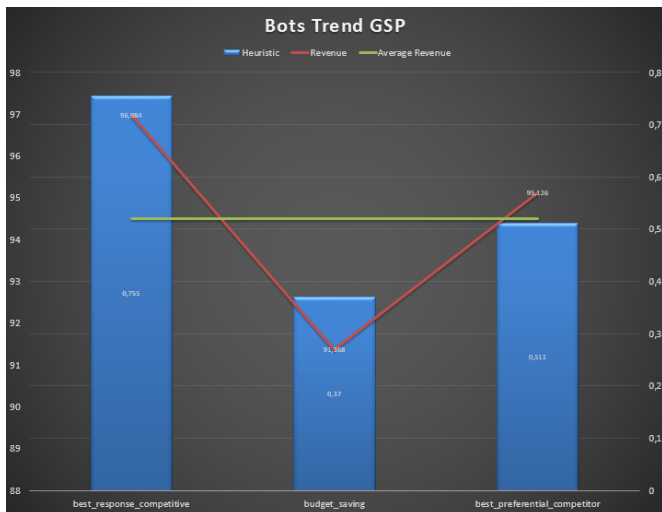
As expected the utility shows the **same** behaviour as in the fixed auction and difference on the revenue are almost irrelevant.

We notice **worsening** on the best_preferential_competitor due the fact that his behaviour is based on the threshold, that can change on different runs, this leading to more variability.

Utility GSP



Heuristic GSP



Consideration on GSP

In regards to utility and revenue, we notice that GPS shows the **same** behaviour, as said in FPA best_preferential_competitor shows a **worsening** for the same reason.

Instead budget_saving **improves** his behaviour, still remaining the worst, because in each run the budget is not fixed.

FPA vs GSP

As seen in the fixed example, we can observe that:

- FPA lead to a **higher** revenue compared to GPS
- GPS hold a **higher** utility and heuristic in relation to FPA

Overview

- 1 Ranking
- 2 Matching
- 3 Search Engine
- 4 Auction
- 5 Summing UP**

Summing UP

Based on all the experiments taken we can assume that the best configuration for a search engine is:

- **Ranking:** indifferent since choice is based on what we are more interested in
- **Matching:** Improved Best Match, since is quicker and gives the same result
- **Auction:**
 - For the advertiser GSP
 - For the seller FPA

Thank you for the attention.