

REPORT K-SHELL DECOMPOSITION

Pankaj Verma, 2015CSB1022
Soumyadeep Roy, 2015CSB1035

December 1, 2018

1 INTRODUCTION

K-Shell decomposition helps in finding out the shell numbers of various nodes in a particular graph. This is particularly helpful in finding out the core of any graph. With the help of this knowledge we are working on ways to reach the core of the graph from any random point with only local knowledge to rely upon.

2 Literature Review

2.1 The H-index of a network node and its relation to degree and coreness

This paper discusses how degree and coreness are related with each other with the help of a property called H-index. H-index is the property of a node that is defined by the relationship given below :

Let $G(V, E)$ denote a graph G with V vertices and E edges making it up. Now choosing any particular node let the degree of the node i be denoted by k_i . Then the zero-order H-index of node i is $h_i^{(0)} = k_i$, i.e the degree of the node. Let the $(n - 1)$ -order H-index of neighbours of i be denoted as $h_{j_1}^{(n-1)}, h_{j_2}^{(n-1)}, h_{j_3}^{(n-1)}, \dots, h_{j_{k_i}}^{(n-1)}$. Now the n -order H-index can be calculated as :

$$h_i^{(n)} = \mathcal{H}(h_{j_1}^{(n-1)}, h_{j_2}^{(n-1)}, h_{j_3}^{(n-1)}, \dots, h_{j_{k_i}}^{(n-1)}).$$

NOTE : Here $y = \mathcal{H}(x_1, x_2, x_3, \dots, x_n)$ signifies that y is the maximum integer such that there exist atleast y elements in $(x_1, x_2, x_3, \dots, x_n)$ each of which is no less than y .

This paper helps in establishing the fact that as the order of H-index value increases to infinity, the value of the H-index converges to the coreness value of the node.

The paper also provides a method for doing asynchronous updating in a dynamic graph, as such graphs are subject to rapid changes in the network. By traditional approach the addition of a single link in such network will need recalculation of H-index sequence for entire network. However Asynchronous updating can still guarantee a convergence to coreness.

2.1.1 Theorem : Convergence of H-index to coreness

For every node $i \forall v \in V$ of an undirected simple network $G(V, E)$, its H -index sequence $h_i^{(0)}, h_i^{(1)}, h_i^{(2)}, h_i^{(3)}, \dots$ will converge to the coreness of node i ,

$$c_i = \lim_{n \rightarrow \infty} h_i^{(n)}$$

Proof: We have $h_i^{(n)} \geq 0$ and $h_i^{(0)} \geq h_i^{(1)}$. After applying mathematical induction we get $h_i^{(n)} \geq h_i^{(n+1)}$. Since $h_i^{(0)}, h_i^{(1)}, h_i^{(2)}, \dots$ is a continuously decreasing sequence and all elements are nonnegative, therefore this sequence has got a lower bound. Let us consider $G'(V', E') \subset G(V, E)$, then for any node i' and any integer $n \geq 0$, $h_{i', G}^{(n)} \geq h_{i', G'}^{(n)}$. Second we set k_{min} as the minimal degree of G . Then for any node i and any integer $n \geq 0$, $h_i^{(n)} \geq k_{min}$ stands. It is clearly valid for $n = 0$. So after applied mathematical induction and it is hence proved. If we set G' the c_i core of G . here we can see $G' \subset G$ and in G' , $k_{min} \geq c_i$, therefore $h_i^\infty \geq k_{min} \geq c_i$.

we denote $G''(V'', E'')$ the induced subgraph containing nodes j s.t. $h_j^\infty \geq h_i^\infty$. node i itself belong to G'' . In G' for any node $j \in V''$ $h_j^\infty > h_i^\infty$. there at least h_i^∞ neighbours of j node with h^∞ value no less than h_i^∞ . Hence G'' 's degree of node are no less than h_i^∞ , so G'' is a subgraph of G 's h_i^∞ -core and it gives us $c_i \geq h_i^\infty$ where c_i is the coreness of i . In the end, from above two inequality we conclude that theorem is proved.

Using the above theorem and definition of H-index, the paper implemented synchronous updating in eight representative real network. Among these two are social network (Facebook and Sex), two collaboration network (Jazz and NS), one communication network (Email), one information network (PB), one transportation network (USAir) and one technological

network(Router).

The paper concludes that the sequence of H-indices quickly converges to the coreness. Apart from degree, H-index and coreness, all intermediate states $h^{(2)}, h^{(3)}, h^{(4)}, \dots$, can also be considered as centrality measures.

2.1.2 Theorem : Asynchronous updating

Given an undirected simple network $G(V, E)$ for every node $j \in V$, we define $g_j = k_j$. In each iteration of the asynchronous updating process, a node i is randomly selected and if g value updated, that is,

$$\mathcal{H}(g_{j_1}, g_{j_2}, g_{j_3}, \dots, g_{j_{k_i}}) - > g_i$$

Where $j_1, j_2, j_3, \dots, j_{k_i}$ are the neighbouring nodes of i . if $|V|$ is finite, this updating process will reach a steady state $(g_1^\infty, g_2^\infty, g_3^\infty, \dots, g_{|V|}^\infty)$ after a finite number of iterations such that the updating at any node will not changes its g value, namely,

$$\forall i \in V, g_i^\infty = \mathcal{H}(g_{j_1}^\infty, g_{j_2}^\infty, g_{j_3}^\infty, \dots, g_{j_{k_i}}^\infty).$$

In the steady state, for every node i we have $g_i^\infty = c_i$.

Proof: Initially we have time $t=0$ and for every node j , $g_j^{(0)} = k_j$. At each time step randomly select a node and perform the \mathcal{H} operator on it. When $t > 0$ node i selected and $g_i^{(t)} = \mathcal{H}(g_{j_1}, g_{j_2}, g_{j_3}, \dots, g_{j_{k_i}})$. First prove that if any node $j \in V$ selected at t_1 and t_2 time such that $t_2 > t_1 \geq 0$ then $g_j^{(t_1)} \geq g_j^{(t_2)}$. At $t = 1$, $g_j^{(t_1)} \geq g_j^{(t_2)}$. We apply mathematical induction. above inequality holds when $n \geq t_1$ and $n \geq t_2$ and we next prove this hold for $n + 1 \geq t_1$ and $n + 1 \geq t_2$. if i selected at time step $t=n+1$ and take t' an arbitrary earlier updating time step of node i such that $n \geq t' \geq 0$.

$$g_i^{t'} = \mathcal{H}(g_{j_1}^{\phi_1}, g_{j_2}^{\phi_2}, \dots, g_{j_{k_i}}^{\phi_{k_i}})$$

$$g_i^{n+1} = \mathcal{H}(g_{j_1}^{\varphi_1}, g_{j_2}^{\varphi_2}, \dots, g_{j_{k_i}}^{\varphi_{k_i}})$$

Here for any $m(k_i \geq m \geq 1)$, $n \geq \varphi_m \geq \phi_m$ and $g_{j_m}^{\phi_m} \geq g_{j_m}^{\varphi_m}$ therefore $g_i^{t'} \geq g_i^{n+1}$

for $0 \leq t_0 \leq t_1 \leq t_2 \leq \dots$ we have $g_i^{t_0}, g_i^{t_1}, g_i^{t_2}, g_i^{t_3} \dots$ is a monotonously non increasing sequence and each element also nonnegative so g_i^∞ has limitation.

we first prove that for any node $j \in V$, $g_j^\infty \geq c_j$. it is proved by contradiction and for second part analogous to proof of theorem 1, after convergence any node $i \in V$ all nodes j such that $h_j^\infty \geq h_i^\infty$. it induced subgraph of G' is g_i^∞ core which is c_i so we conclude $c_i \geq g_i^\infty$.

so above two inequality we proved theorem.

2.2 Estimating Shell-Index in a Graph with Local Information

This main idea of this paper is to find the shell index of a node using local information. This will help to avoid the major drawback of K-Shell decomposition i.e. requirement of the entire graph, which is not feasible for large scale dynamic networks.

2.2.1 Theorem

The Shell index of a node u can be computed as $k_s(u) = \mathcal{H}(k_s(v) | \in ngh(u))$, where $ngh(u)$ is the set of the neighbours of node u .

Proof : In the i^{th} iteration, nodes which have i or less connections with other nodes are removed. As it is the i^{th} iteration, therefore these removed nodes are connected with nodes which have shell index i or greater than i .

2.2.2 Hill Climbing Based approach to identify top rank nodes

In this algorithm inputs are Graph G , Initial node u , Repeat count number k (maximum no of times a crawler is allowed to reach a local maxima) and $maxindex$ (maximum $H_2 - Index$ in the graph G). Starting from the node u the crawler traverses to one of its non visited neighbours which also has the highest $H_2 - Index$. Upon continuing this process the crawler at one point reaches a maximum. If the maximum is a local maxima then the counter (initially 0) is increased by 1 and the flow of travel is passed on to one of its non visited neighbours randomly. This process continues until the crawler reaches the $maxindex$ or the counter reaches the value of k , at which point traversal is terminated.

Some modifications in the given algorithm includes traversal to highest degree neighbour instead of randomly chosen node in case of local maxima.

3 Work Done

Various methods were applied to reach from periphery to core based on only a limited number of node traversal. The following methods are discussed below :

3.1 Random walk

Here starting from any node in the graph, the core of the graph was to be reached. Movement from one node to the next neighbouring node was done randomly.

3.2 Hill climbing

Here starting from any node in the graph, the core of the graph was to be reached based on the value of $H_2 - Index$ of the neighbouring node. Two variations of the hill climbing method was seen.

3.2.1 Travelling highest degree node which is least travelled

Traversal to subsequent nodes is decided based on the highest $H_2 - Index$ neighbours of a node and among those the least travelled node. This allows traversal to be towards higher $H_2 - Index$ nodes mostly.

3.2.2 Travelling highest degree node which is least travelled but with a chance of random walk in every step

Traversal for this process is similar to the previous method but with a slight change. In every step of traversal there is a probability that the next node will be decided randomly from its neighbours rather than the usual process of highest $H_2 - Index$.

4 Dataset

- Facebook Combined
- Ca-CondMat
- Ca-GrQc
- Com-dblp.ungraph

- Email-EU-core
- RoadNet-CA
- P2p-Gnutella31
- Loc-gowalla_edges
- Gemsec_deezer_dataset RO_edges
- Gemsec_deezer_dataset HR_edges
- Gemsec_deezer_dataset HU_edges

4.1 Dataset Details

Dataset List					
Sl.no	Dataset	Directed or Undirected	Number of Nodes	Number of Edges	Maximum Shell number
1	Facebook Combined	Undirected	4039	88234	96
2	Ca-CondMat	Directed	23133	186936	21
3	Ca-GrQc	Directed	5242	28980	26
4	Com-dblp.ungraph	Undirected	317080	1049866	47
5	Email-EU-core	Undirected	1005	16706	35
6	RoadNet-CA	Directed	1965206	-	-
7	P2p-Gnutella31	Directed	62586	147892	6
8	Loc-gowalla_edges	Directed	196591	1900654	51
9	Gemsec_deezer_dataset RO_edges	Undirected	41773	125826	7
10	Gemsec_deezer_dataset HR_edges	Undirected	54573	498202	21
11	Gemsec_deezer_dataset HU_edges	Undirected	47538	222887	11

5 Results

5.1 Random Walk

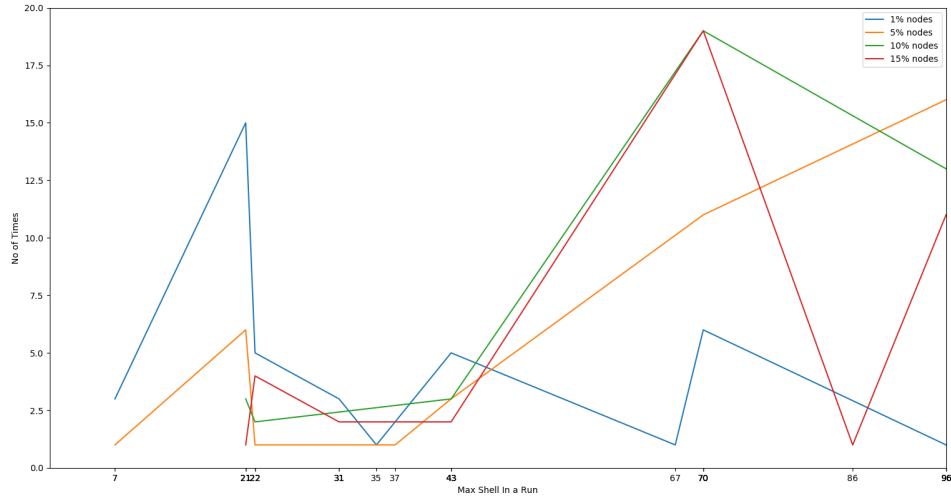


Fig 1.1 Graph for Random Walk for Dataset 1.

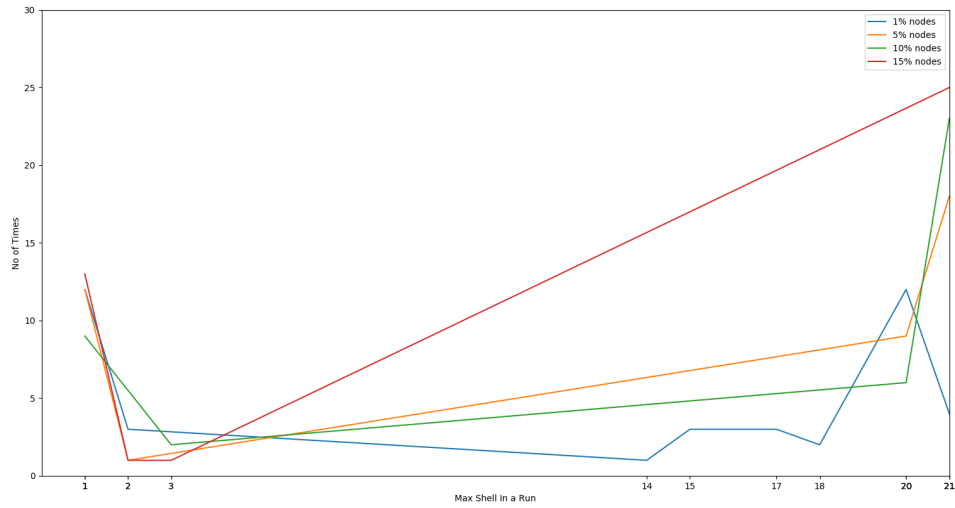


Fig 1.2 Graph for Random Walk for Dataset 2.

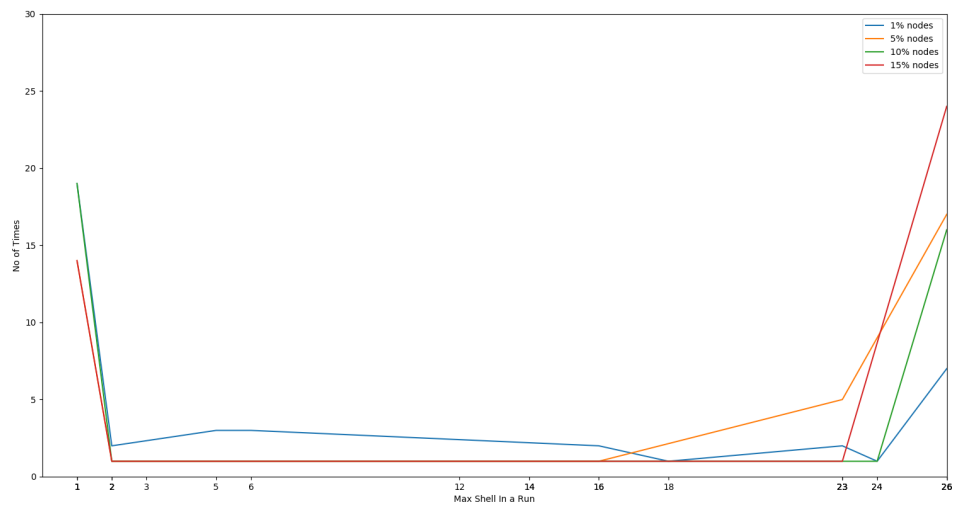


Fig 1.3 Graph for Random Walk for Dataset 3.

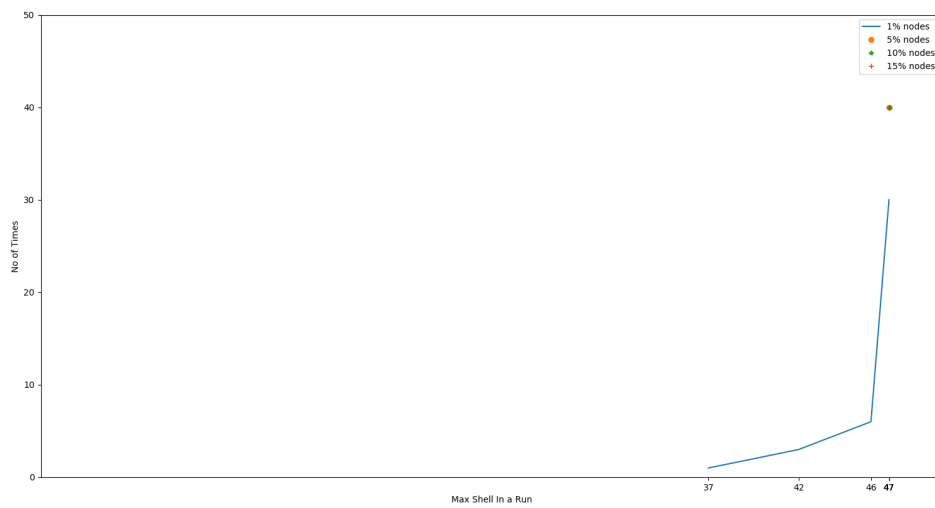


Fig 1.4 Graph for Random Walk for Dataset 4.

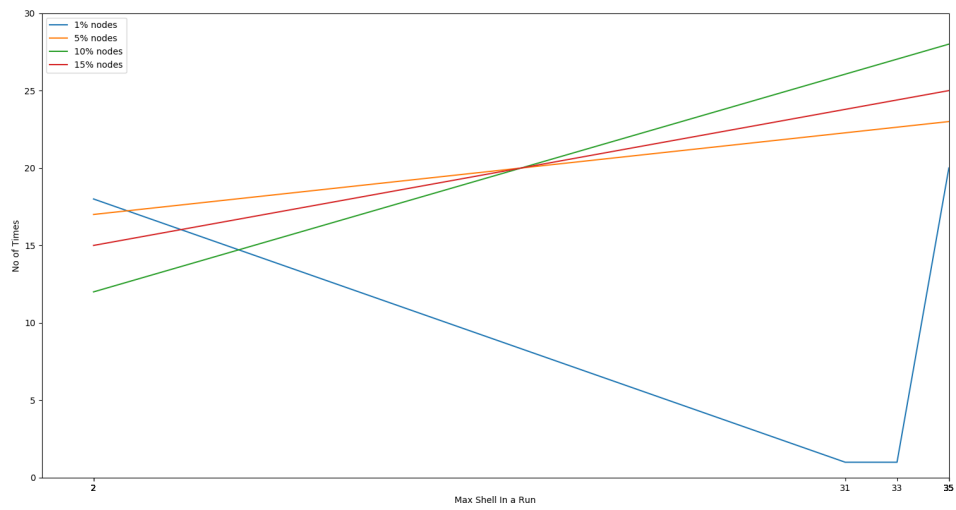


Fig 1.5 Graph for Random Walk for Dataset 5.



Fig 1.6 Graph for Random Walk for Dataset 7.



Fig 1.7 Graph for Random Walk for Dataset 8.

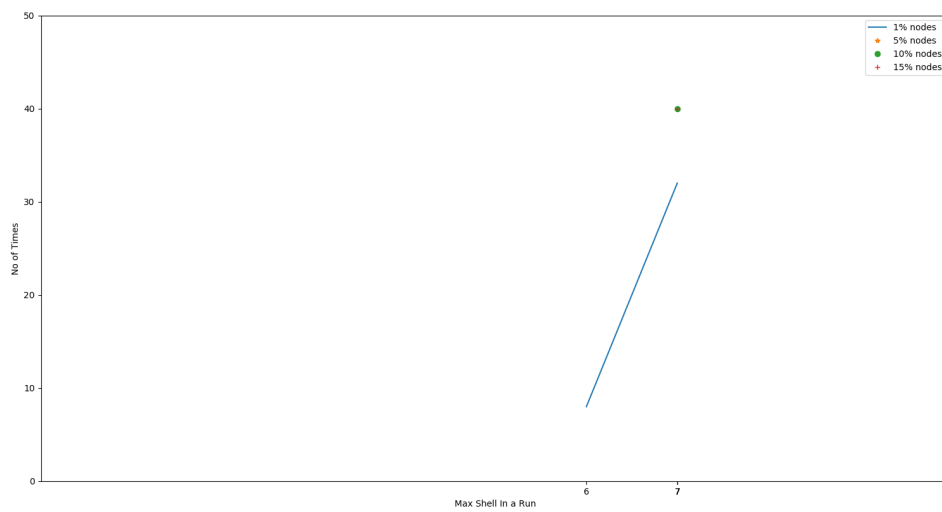


Fig 1.8 Graph for Random Walk for Dataset 9.



Fig 1.9 Graph for Random Walk for Dataset 10.

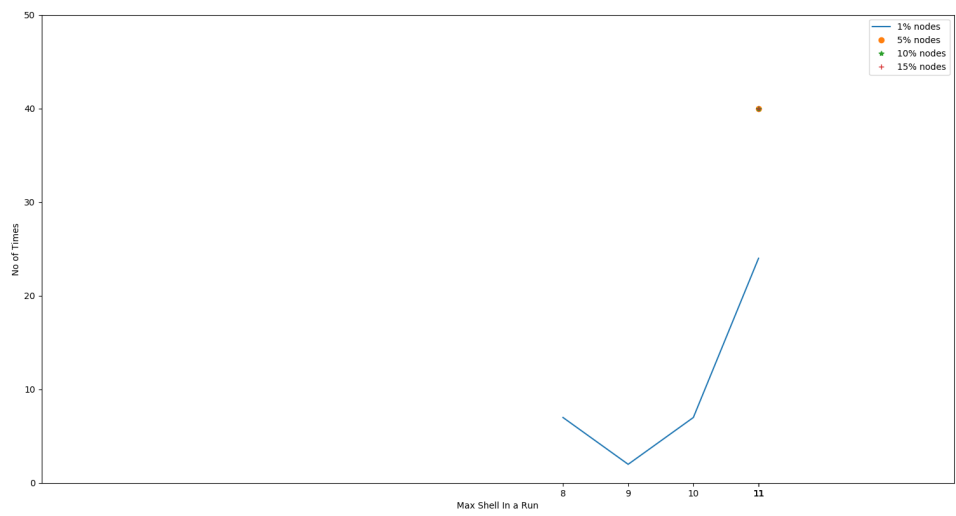


Fig 1.10 Graph for Random Walk for Dataset 11.

After plotting the graph of random walk for all the datasets it is seen that in all datasets, random walk leads to reaching the maximum shell number. Also in most graphs, increasing the percentage of nodes for random walk almost always increases the number of time the maximum shell number is reached.

For datasets 4, 7, 8, 9, 10 and 11 it is seen that for some percentage of nodes all the runs reach the maximum shell. This may signify that for these datasets either the number of shells are too less compared to the number of nodes or that the lower numbered shells are very well connected to the core shells.

In contrast to above datasets, dataset 1, 2, 3 and 5 is seen to contain an intermediate shell, as maximum shell reached in considerable no of runs in most percentage of nodes. For dataset 1 these intermediate shells are notably 70 and 21. Similarly dataset 2 has shell 1 and 20, dataset 3 has shell 1 and dataset 5 has shell 2. This goes to show how well these intermediate shells are interconnected among themselves. This can be verified by checking the Traversal pattern for runs in these graphs, where once these shells are entered while traversing, it is hard to get out.

5.2 Hill Climbing

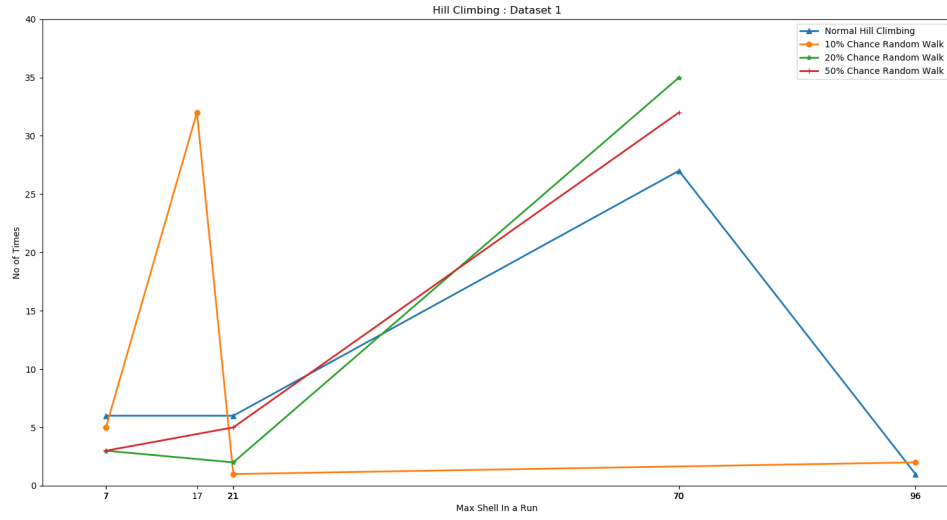


Fig 2.1 Graph for Hill Climbing for Dataset 1.

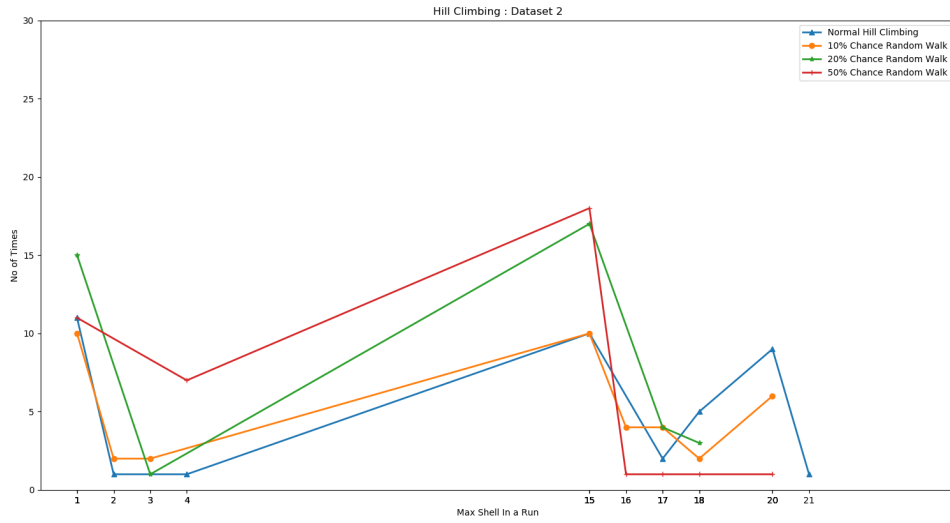


Fig 2.2 Graph for Hill Climbing for Dataset 2.

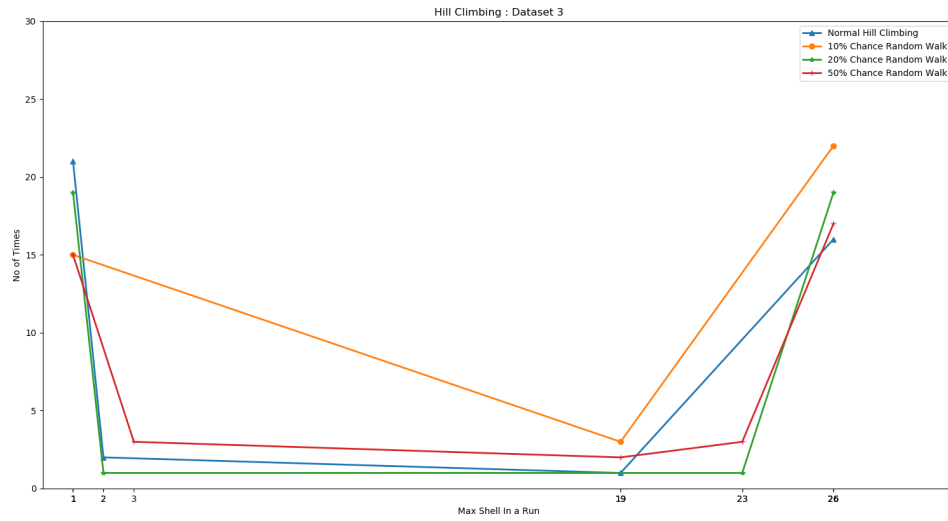


Fig 2.3 Graph for Hill Climbingk for Dataset 3.

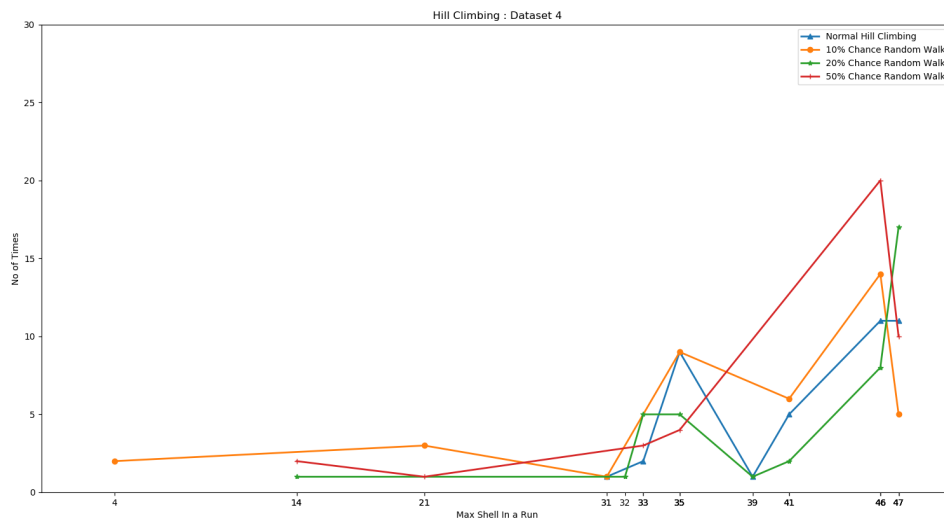


Fig 2.4 Graph for Hill Climbing for Dataset 4.

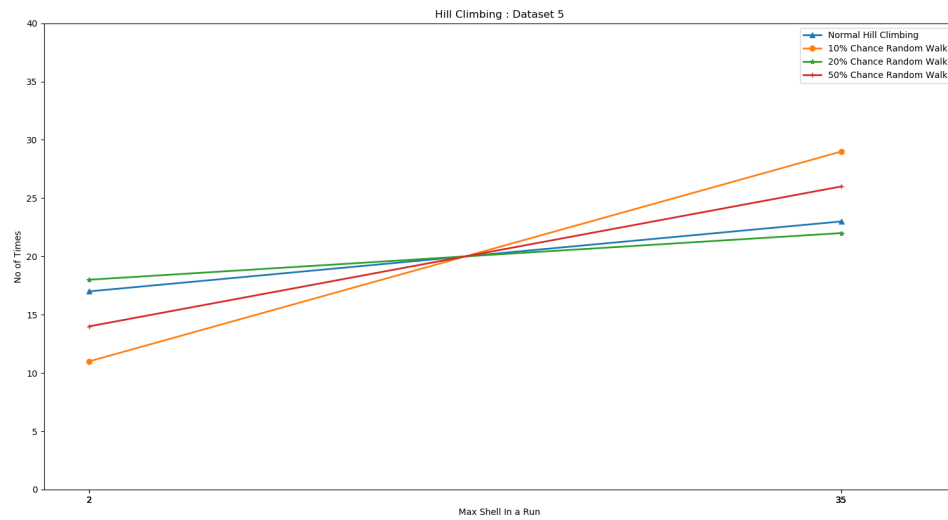


Fig 2.5 Graph for Hill Climbing for Dataset 5.

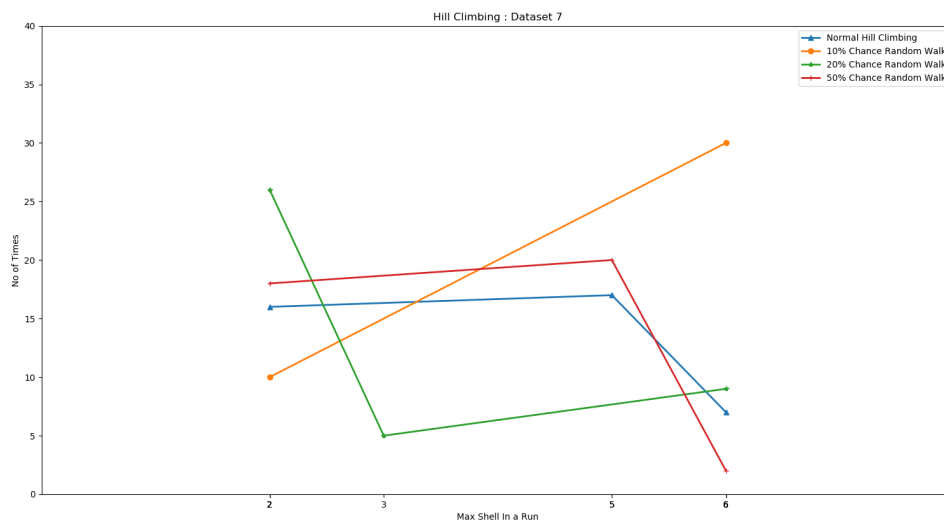


Fig 2.6 Graph for Hill Climbing for Dataset 7.

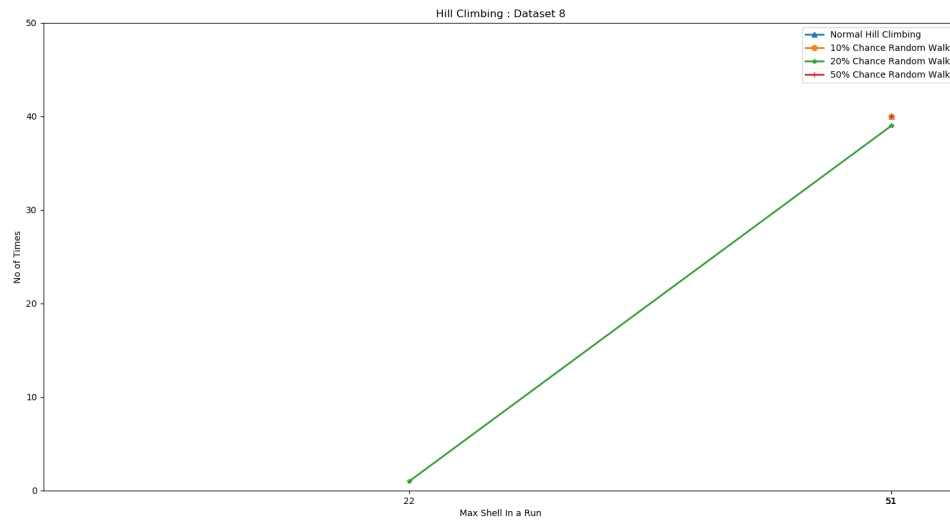


Fig 2.7 Graph for Hill Climbing for Dataset 8.

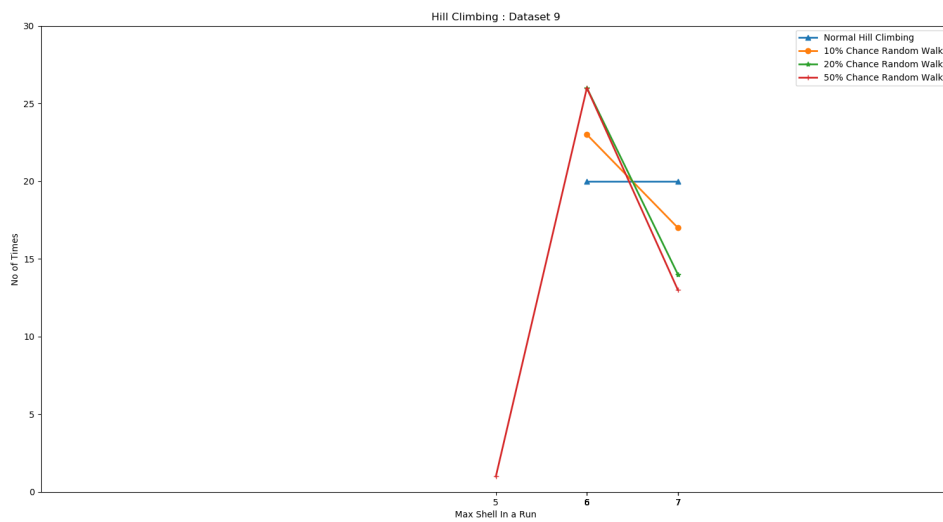


Fig 2.8 Graph for Hill Climbing for Dataset 9.

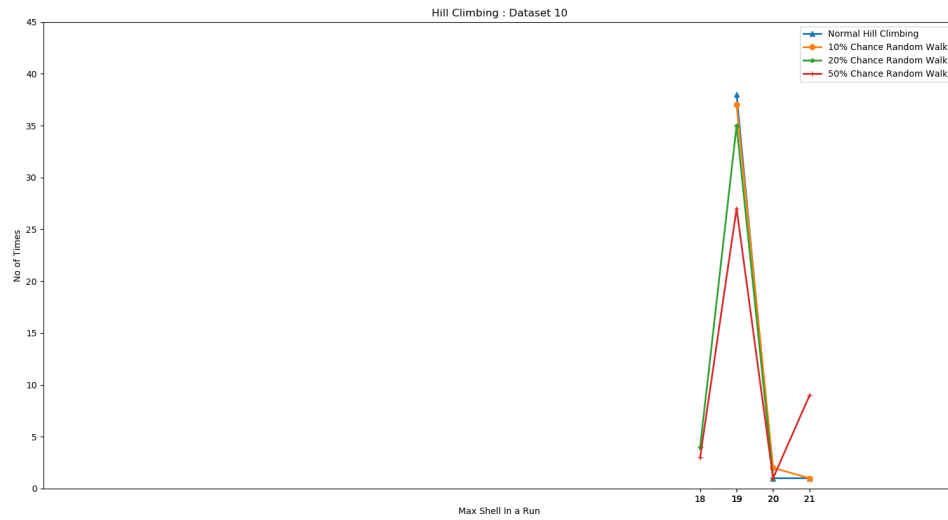


Fig 2.9 Graph for Hill Climbing for Dataset 10.

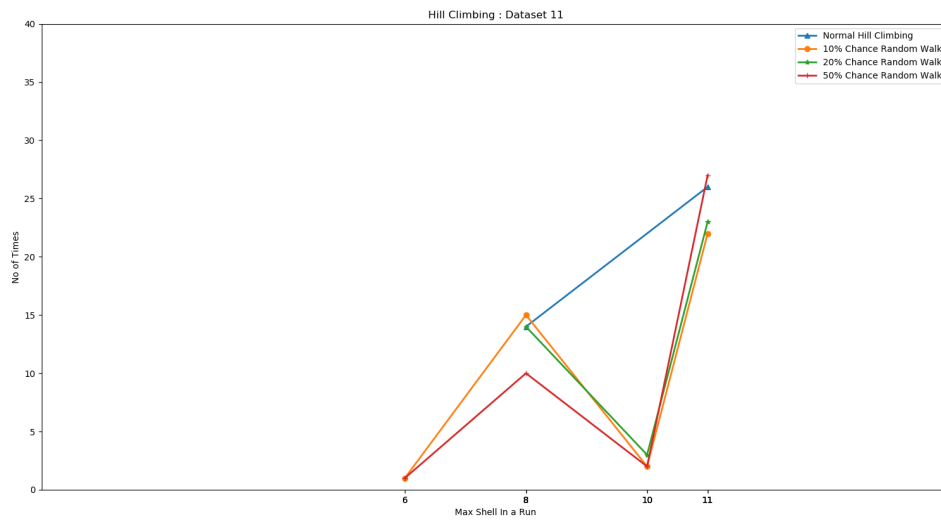


Fig 2.10 Graph for Hill Climbing for Dataset 11.

Unlike the graphs for Random walks, the graphs for Hill climbing are much more diverse. However the datasets have some very consistent similarity in between the two.

Datasets 1, 2, 3 and 5 show the same intermediate states in hill climbing similar to the intermediate states in random walk. Dataset 1 shows 70, 21 and 7, dataset 2 shows 1, 15 and 20, dataset 3 shows 1 and dataset 5 shows 2. New intermediate states are shown in Dataset 4, 7, 9, 10 and 11. For dataset 4 there is 46 and 35, dataset 7 has 2, dataset 9 has 6, dataset 10 has 19 and dataset 11 has 8. Unlike in random walk, here the intermediate state may refer to the location of local maxima.

For Dataset 2 and 3, we find the intermediate value to be 1. Now since 1 is the lowest shell possible, therefore there is no way that local maxima exists in shell 1. On checking the traversal route of these two datasets, the existence of traversal loops is found. For example a route like 86060 – 82784 – 86060 – 872023 – 86060 – 50238 – 86060 – 82784 ... in Dataset 2. One of the major reasons of such loops is caused by Dataset 2 and 3 being directed graphs. Directed graphs cause such problem as in the previous example nodes like 50238 are directed towards 86060 but the reverse is not true. Therefore if at any point of traversal a loop like above is entered, then the hill climbing process is halted. A unique version of the traversal loop is found in a run of Dataset 3. In this run the traversal is 9471 – 19586 – 7013 – 4364 – 18952 – – 18952 - The interesting fact here is that the loop consisted of 21 nodes all belonging to shell no 19. Also all the nodes in the loop had more than 1 neighbour. However the control was passed to the highest highest $H - 2index$ and among all the neighbours, the node which created the loop had the highest $H - 2index$. The introduction of percentage chance of random walk in each step helps avoid such problems. This can be seen in graph for Dataset 3, where addition of randomness has increased the chances of reaching the maximum shell and also decreasing the chances of staying in very low intermediate shell. Similar phenomenon is partially visible in Dataset 2.

In Dataset 10, a substantial number of runs get stuck in intermediate shell 19 especially through reaching the node 25352, which causes a loop. Similarly there is a slight improvement by increasing the randomness. However in many datasets it is seen that involvement of randomness has led to the decrease in result. This is visible in dataset 9. The reason for this is unknown but it may be due to fact of traversing to degrading route rather

than the correct route. In dataset 11, all the runs that reached the maximum shell is seen to reach the maximum shell is very short number of steps. Whereas the runs that are stop in the intermediate step usually take a huge number of steps. This is clearly due to the fact of proper connectivity in some parts of the graphs. Similar is the case of Dataset 1 where the very few times that the run reached the maximum shell, it took 2 to 3 steps to reach there.

Hill climbing on other datasets causes at least some portions of the run to reach the maximum shell. The success of reaching the maximum shell in case of hill climbing may be hindered by the strict restriction on its traversal pattern, something that did not bother Random Walk.