

Analyzing Python Package Dependency

Avadhoot Agasti
aagasti@indiana.edu

Abhishek Gupta
abhigupt@iu.edu

ABSTRACT

The Python packages have complex nature of inter-dependencies. Creating network structure of the package dependency can help in understanding several aspects of these dependencies. For example, we can easily understand the most basic packages and the packages which have complex dependencies. We can also understand the nature of this network (scalefree, random graph etc) and then utilize the concepts from these network structure. Further, this analysis can be extended to any other programming language to understand the package network. The scope of this project is to focus on python based packages and navigate to all dependent packages to build a network graph of these packages and then analyze this network.

KEYWORDS

scale-free, random networks, small-world, packages

1 INTRODUCTION

When we try to install a python library, or rather when we use a python library first thing that we need to do is to *pip install* it. Most of the times, a python library is dependent on another library. Our idea is to build a network or graph of the python library dependency. In this network, every library will be a node of the network and the dependency will form an edge. Such a graph can be queried to answer several questions like

- Which are the most core packages which are widely, directly or indirectly, used in large number of other packages?
- In which subject-area the new development is happening. We can pivot this solution around the small number of packages which are included in large number of packages. For example, if networkX is being used by lot of new packages then we can say that there is lot of development happening in network science
- What all packages will be impacted due to changes in a base package (e.g. if we find a severe bug in networkX, what are the other packages which can be potentially impacted due to the bug fix)

2 MOTIVATION

There is not much work done around this topic and after reading the concepts of network science it looks very interesting to find out how different module relate and how these relations can be interpreted. This paper *Power Laws in Software* [?] does analysis of power law distribution in a software application at class level and function level. It did analysis of java, perl, c/c++ etc applications and did establish a pattern that these applications do follow power law distribution.

It depicts the size of datasets studied for this this analysis. However, we couldn't find a paper or research which analyzes these

Table 2: Review of other evidence.

Dataset	size	k_{in}	k_{out}
Java, C/C++ [Valverde and Solé 2003]	27-5,285	1.94-2.54	2.41-3.39
C/C++ [Myers 2003]	187-5,420	1.9-2.5	2.4-3.3
Java [Wheeldon and Counsell 2003]	NA	0.71-3.66	
Smalltalk [Marchesi et al. 2004]	1,797-3,022	2.07-2.39	2.3-2.73
Object graphs [Potanin et al. 2005]	15,064-1,259,668	2.5	3

Figure 1: Review of Languages

patterns on all modules present in a given language. Its very interesting problem to solve and find out how various python packages are dependent on each other and does it make sense to bundle some of the very common packages used along with base python distribution. It is interesting to understand the dependency structure but also to understand the community structure of these modules and what kind of network they form. Any other phenomenon exist when we analyze this graph. This analysis can be used to optimize the imports, remove cyclic dependencies between modules. Similar approach can be extended to other programming languages like java, ruby, perl etc.

3 REQUIREMENTS

The goal of this analysis is

- to create a network graph of these dependencies
- understand the dependencies and important nodes
- identify if its a scale-free network
- identify any other properties depicted by this graph
- update the graph and build the graph as and when needed

At the end of this exercise we should be able to run this analysis on available python [1] packages. Also, we need to understand and compute the other properties of the graph like average degree, average clustering co-efficient, average path length etc Also verify if the graph is *scale free* or just a *random* graph.

4 TECHNICAL SOLUTION

4.1 Data Sources

Our approach is to identify the dependent packages by looking at the dependencies using *pip show* output. For example

```
pip show networkx
```

Name: networkx

Version: 1.11

Summary: Python package for creating and manipulating graphs and networks

Home-page: <http://networkx.github.io/>

Requires: decorator

We perform most of our coding in *Python* and *networkx* itself. We used *Gephi* for visualizing the network.

During this exercise we ran into few technical challenges for example to perform this analysis we have to install each python

package locally which could be really resource intensive specially in terms of disk space. Also, while doing the analysis and loading the graph for 128k packages could be memory and cpu intensive. To get around this problem we had to reduce the analysis to less number of python packages and keep adding more packages as we make progress. Further enhancements can be made by running the crawler on pypi.org [?].

5 GRAPH ANALYSIS

We created a Jupyter notebook to analyze the graph generated by pydgraph script. We executed this notebook on the graph generated by one of our installations which had about 197 packages installed. Below section explains the results of this analysis.

5.1 Graph encoding indirect dependencies

The graph created by pydgraph script captures the direct relationship between the nodes. This essentially means that if package A depends on package B then there will be edge between node A and node B. However, if package C depends on package B, then indirectly, it also depends on package A. In the graph created by pydgraph, the edge between A and C is not created. We created another variation of the graph (called GNew) which adds a link for such indirect dependencies. So, in GNew, an edge between A and C exists.

5.2 Degree Distribution

The fig 4 shows the degree distribution of original graph while the fig 5 shows the degree distribution of the new graph which has indirect dependencies encoded. After observing these graphs, we can say that

- The original graph (G) does not seem to follow degree distribution pattern of scale free graph. Especially on the log-log scale, it is not a straight line like in scale free graph.
- However, the updated graph (GNew) where indirect dependencies are directly added to each node, the graph looks like scale free graph with power law degree distribution

5.3 In and Out Degree

The analysis of in-degree and out-degree of GNew can explain us the packages which has large number of dependencies and packages which are used by other packages a lot. The higher value of in-degree means the package depends on several other packages. Similarly, higher value of out-degree means the package is used by several other packages. After observing fig 6 we can say that nodes (packages) like "uniq" or "zope.publisher" have larger dependency on other packages. Fig 7 explains that packages like "setuptools", "six" and "zope.interface" are used by several other packages.

5.4 Friendship Paradox

For original graph, only 21 percent nodes follow friendship paradox. However, for indirect-dependency-graph, the friendship paradox holds true for about 83 percent nodes. This is very useful since using indirect dependency graph, we can quickly identify the hubs in this network. Essentially, we can identify the popular packages very easily without even knowing degree distribution of entire network.

5.5 Communities

Fig 8 explains that the network has large number of communities. In this specific installation, it can be observed that:

- There are 10 identified communities (used resolution factor of 5). The bigger communities can be easily explained. For example, community around testing packages like pytest, community around notebook packages (notebook, nbconvert, ipykernel), community around django and community around zope.
- This installation has bunch of completely separate packages like 1) pip and pipdeptree 2)flake and 3)networkx and decorator
- The node size indicates betweenness centrality. It looks like pytest is the one which is used mostly across all communities.

5.6 Change Propagation

This network can also be utilized for several other purposes. For example, if package A is undergoing a major release, we can identify all the dependent packages to understand the impact of the change. The paths between two nodes can explain the different ways changes in Package A impacts Package B. Fig 9 and Fig 10 explains this.

6 RELATED WORK

- The project Pipedtree - [2] is command line utility which allows user to see the installed packages in the form of dependency tree. However, this utility is focused more on solving dependency conflicts than the kind of analysis we propose to perform.
- Analysis of 30K Github projects [3]. The project was aimed to analyze the 30k different Java, Ruby and Javascript projects on Github to understand the top libraries being used. While their analysis was similar to what we plan to do, the approach was not network based.
- This paper *Power Laws in Software* [?] does analysis of power law distribution in a software application at class level and function level. It did analysis of java, perl, c/c++ etc applications and did establish a pattern that these applications do follow power law distribution

7 FURTHER ENHANCEMENT

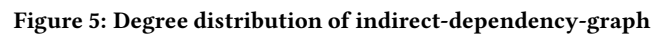
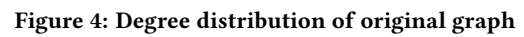
Further enhancements can be done to crawler stage where instead of relying on modules installed locally it can crawl on module available on pypi.org. This will reduce an additional step of installing the modules locally and save the disk space as well as more flexible to crawl on web. Another enhancement can be made to make it more generic to do the analysis on any given language for example java maven dependencies or perl packages or ruby version manager and packages etc

8 ACKNOWLEDGEMENTS

The authors thank Prof. YY Ahn for his technical guidance. The authors would also like to thank TAs of Network Science class for their valued support.

All project and report document can be found at [github project](#).

- [1] Python.org. Python, packages. Web Page. Accessed: 2017-02-02.
- [2] Python.org. pipdeptree, command line utility to show dependency tree of packages. Web Page. Accessed: 2017-02-02.
- [3] Tal Weiss. OverOps, blog by tal weiss. Web Page. Accessed: 2017-02-02.



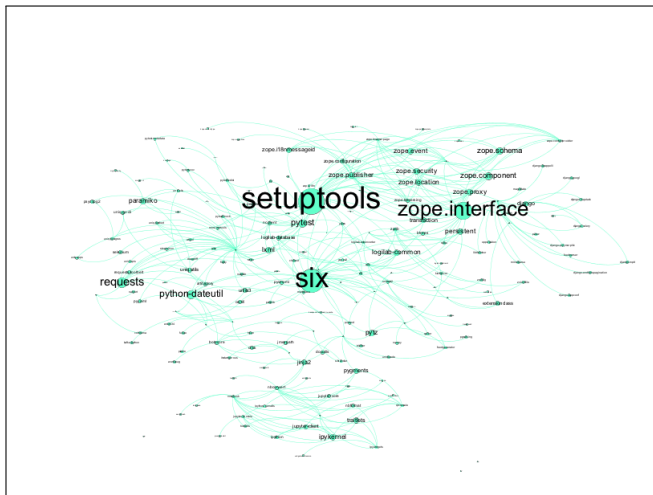


Figure 7: Out-Degree of indirect-dependency-graph

```
In [116]: package_under_change = "jinja2"
my_package = "notebook"
paths = nx.all_simple_paths(G, package_under_change, my_package, cutoff=None)

for p in paths:
    print(p)

['jinja2', 'nbconvert', 'notebook']
['jinja2', 'notebook']
```

Figure 10: Impact of Change

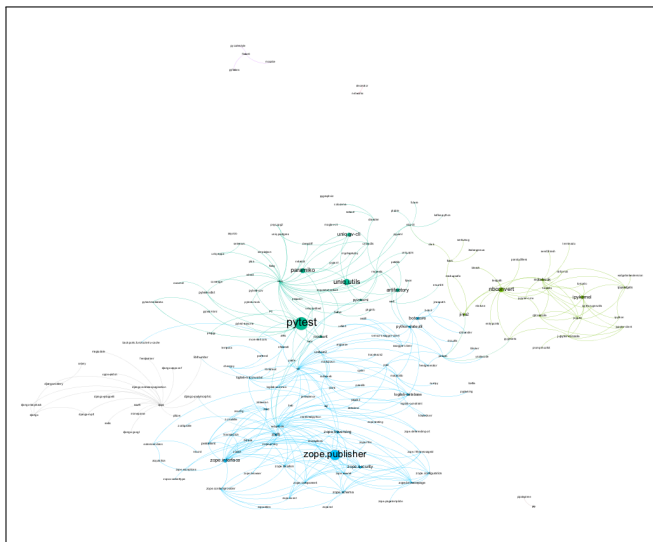


Figure 8: Communities

```
In [115]: package_under_change = "uniq"
impacted_packages = nx.neighbors(GNew, package_under_change)
impacted_packages

Out[115]: ['requests',
'lxm1',
'pyyaml',
'setuptools',
'six',
'more-itertools',
'pytest',
'attrs',
'pluggy',
'py',
'urllib3',
'python-dateutil',
'artifactory',
'pathlib',
'deepdiff',
'et-xml']
```

Figure 9: Impact of Change