

report

June 3, 2019

1 Setup

Collecting coclust

Requirement already satisfied: scipy in /opt/conda/lib/python3.6/site-packages (from coclust)

Requirement already satisfied: scikit-learn in /opt/conda/lib/python3.6/site-packages (from coclust)

Requirement already satisfied: numpy in /opt/conda/lib/python3.6/site-packages (from coclust)

Installing collected packages: coclust

Successfully installed coclust-0.2.1

2 Loading data

A fonction to get the name of a library according to its definition and the depth.

‘java.util.ArrayList’, depth=2 -> ‘java.util’

Define the mapping of imports to an ids.

Some basic vector operators

Count for every commits the number of time a library was imported in all modified files.

Out[12]: {}

Total number of commit : 1993

Made by 8 unique contributors

Represented by 1097 imports

mapped to 174 imports

Sum the commits by users.

Then drop the libraries imported once or less.

Finally sort the data in row and column by the number of modifications.

Out[15]:

		index			
		mean	min	max	count
author					
Anders Nawroth	1403.568627	598	1819	51	
Emil Eifrem	273.608108	0	948	74	

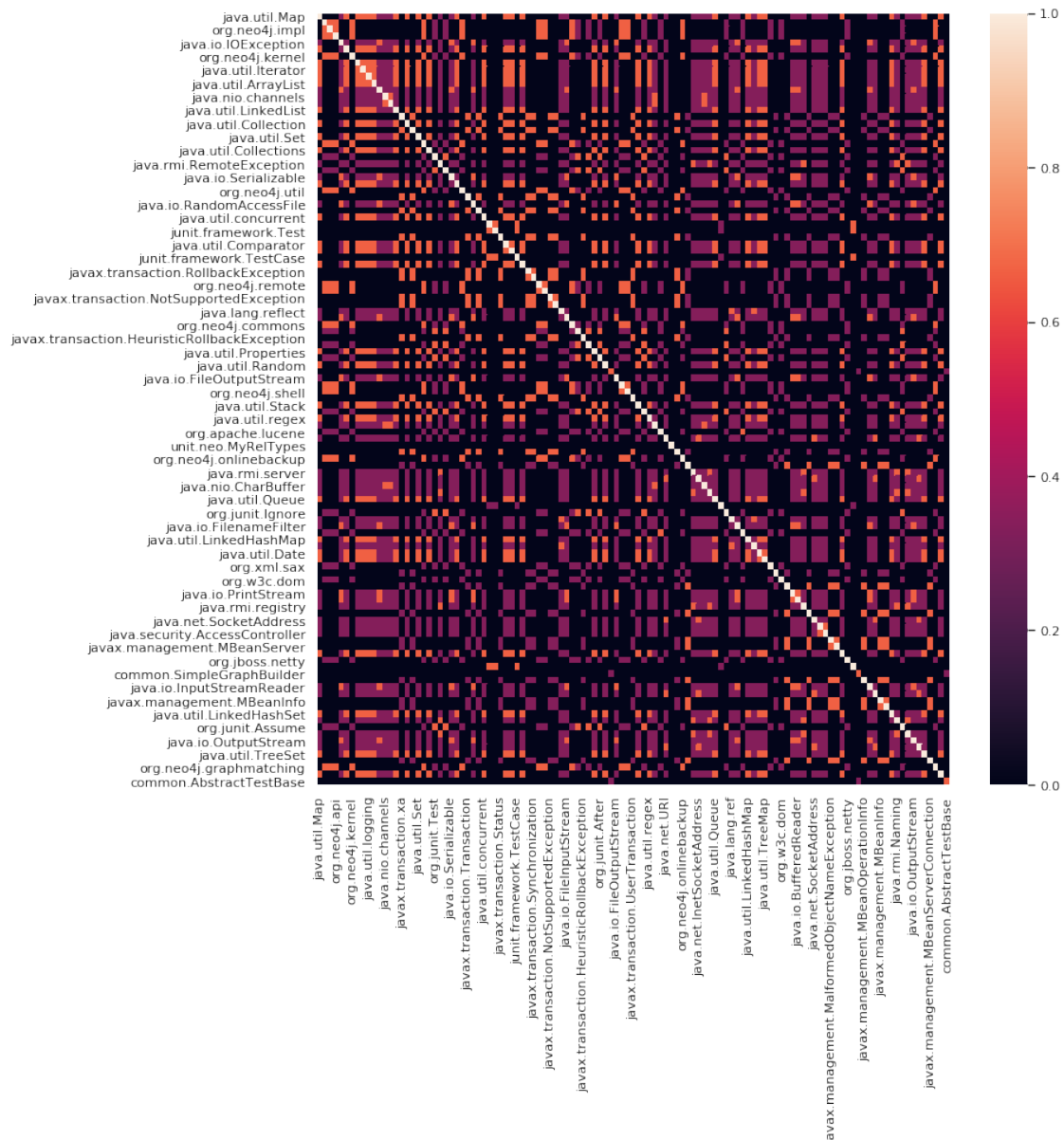
Henrik Larsson	398.000000	398	398	1
Johan Svensson	746.275904	40	1989	830
Mattias Persson	1235.443843	6	1992	739

Out[18]:

	first	last	count	coef
javax.transaction.InvalidTransactionException	1	1411	117	0.416418
java.lang.annotation	612	1823	18	0.774223
javax.transaction.Status	1	1910	666	0.882552
java.io.FilenameFilter	1	1859	66	0.817335
javax.transaction.Transaction	1	1984	858	0.986544

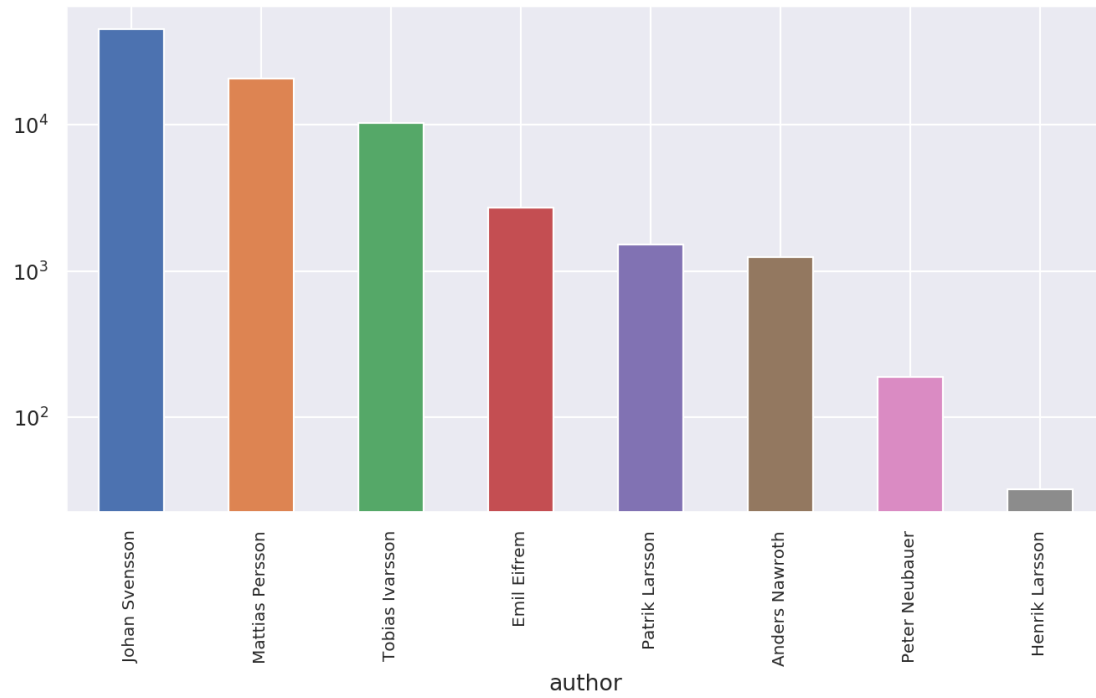
2.1 Distance between libraries

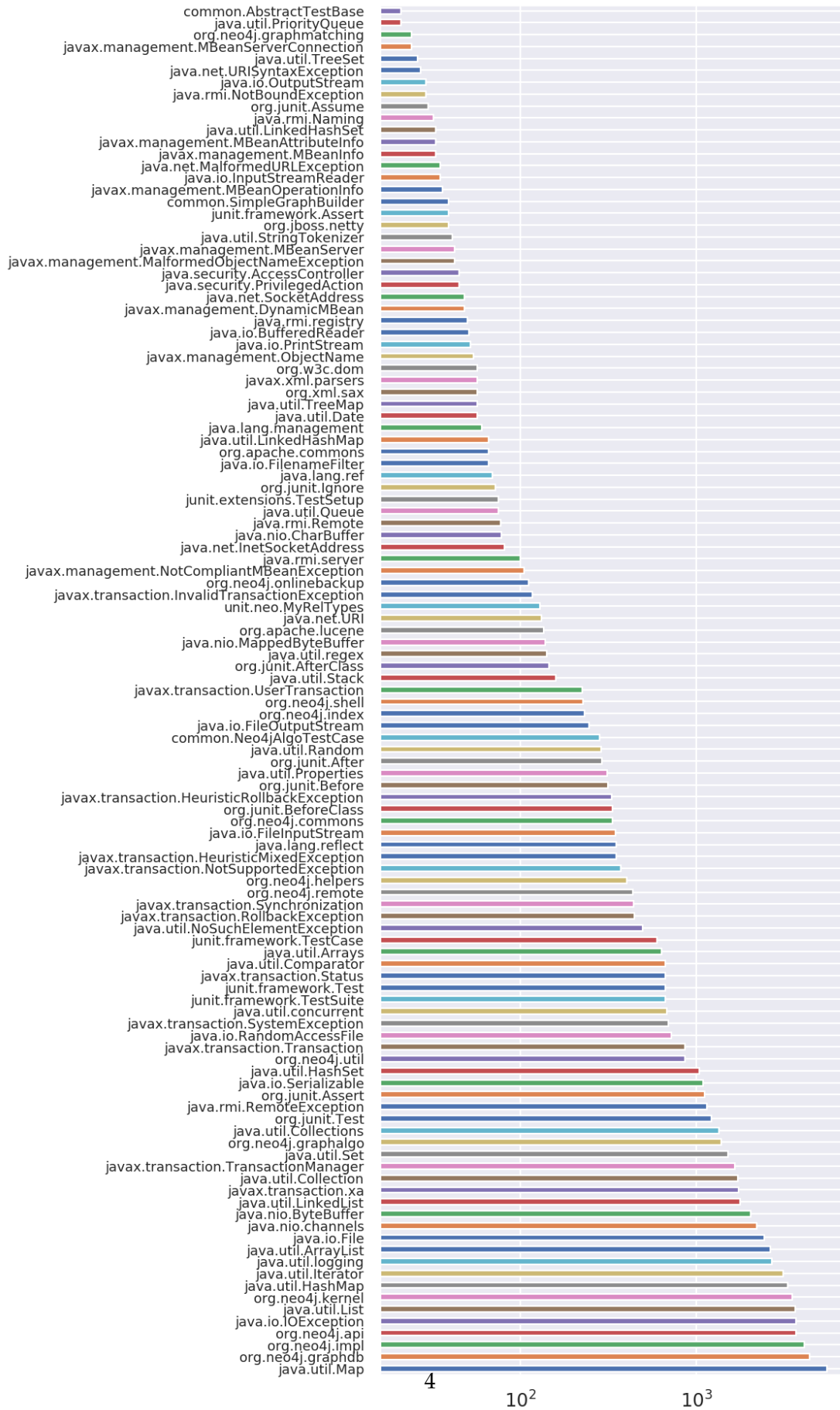
Out[19]: []

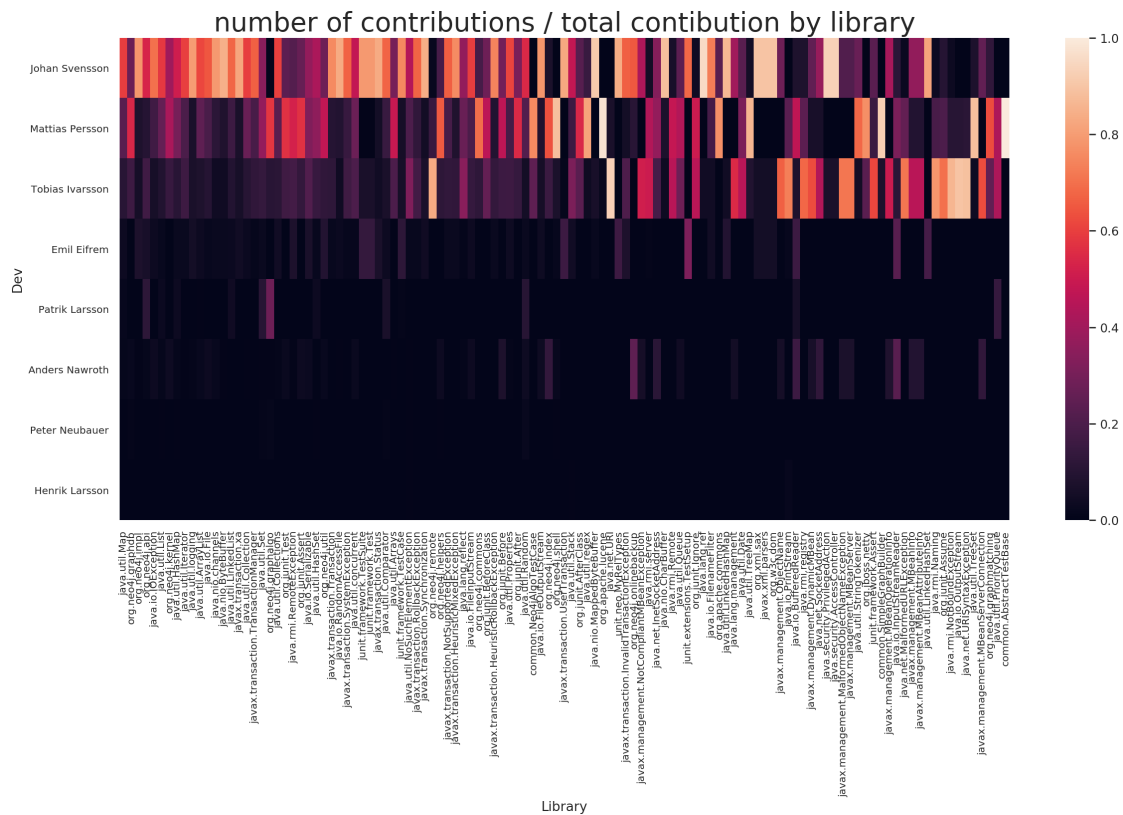


3 Analysis

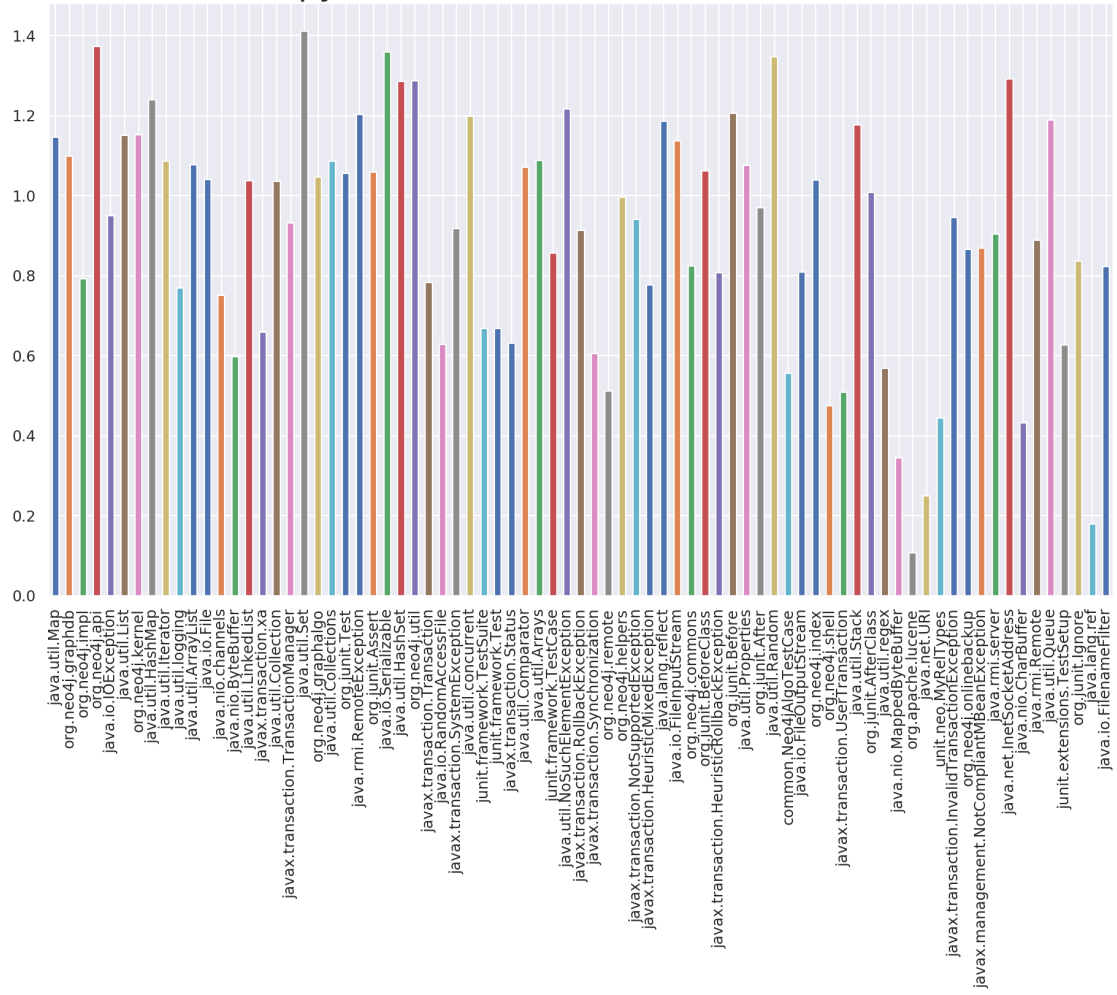
3.1 Summary

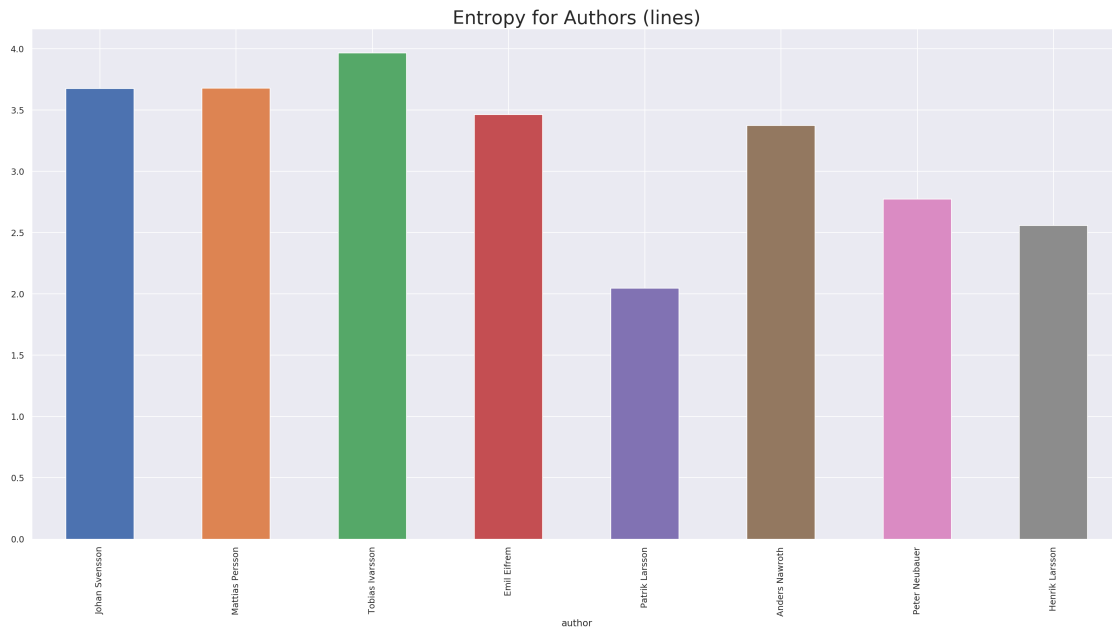






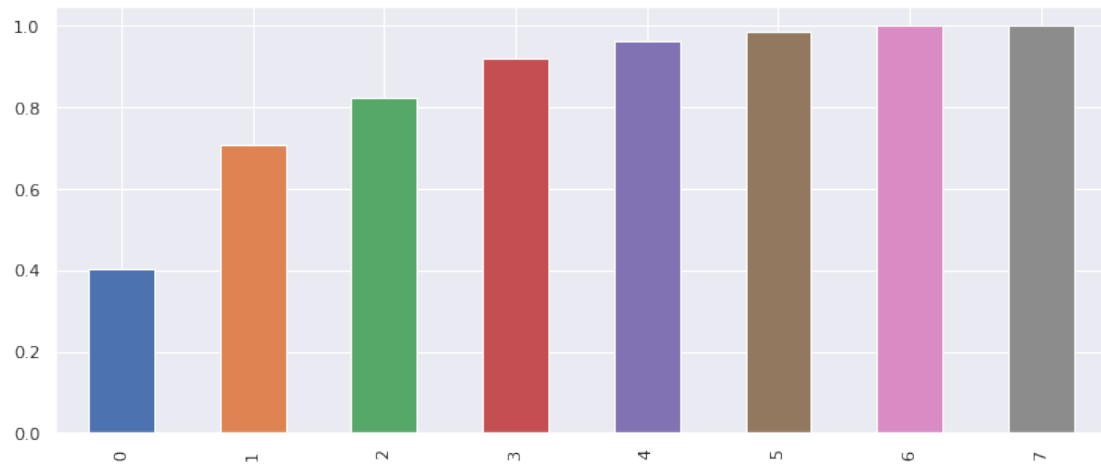
Entropy for libraries (column) (75 most used)





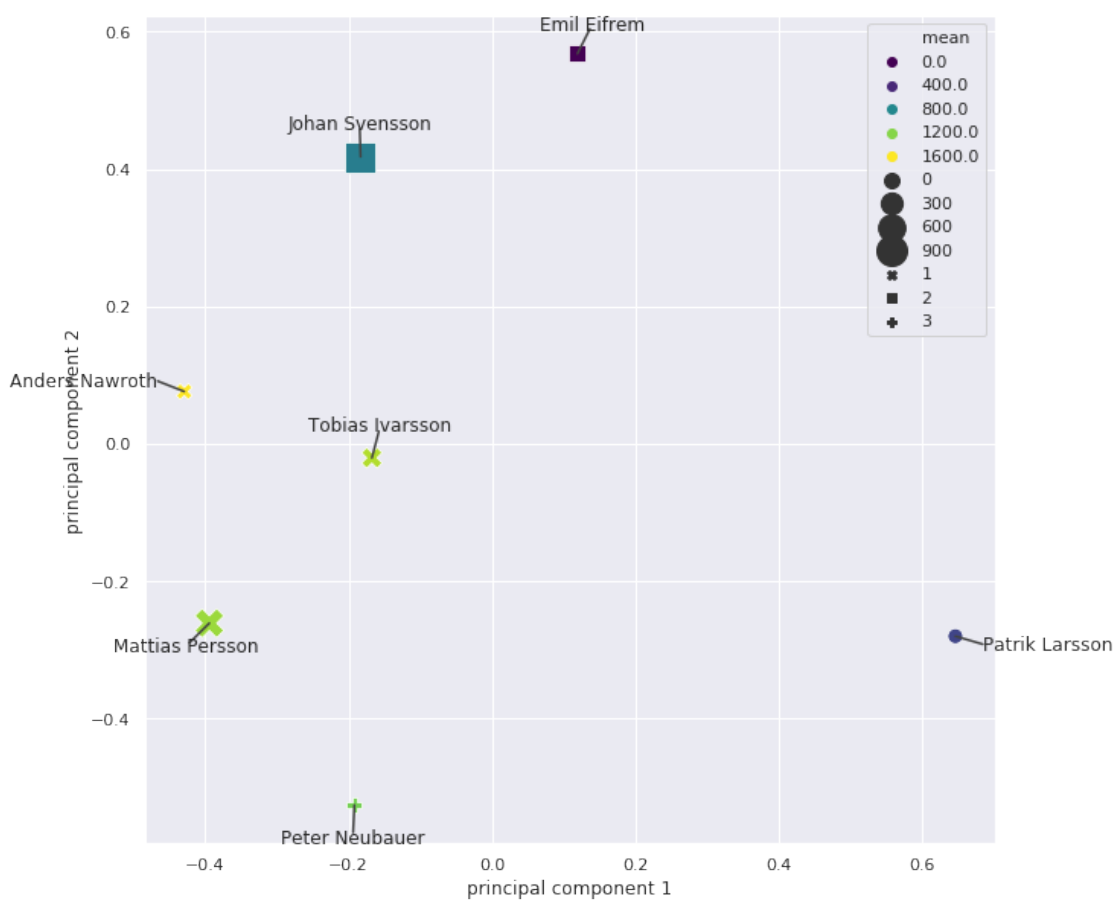
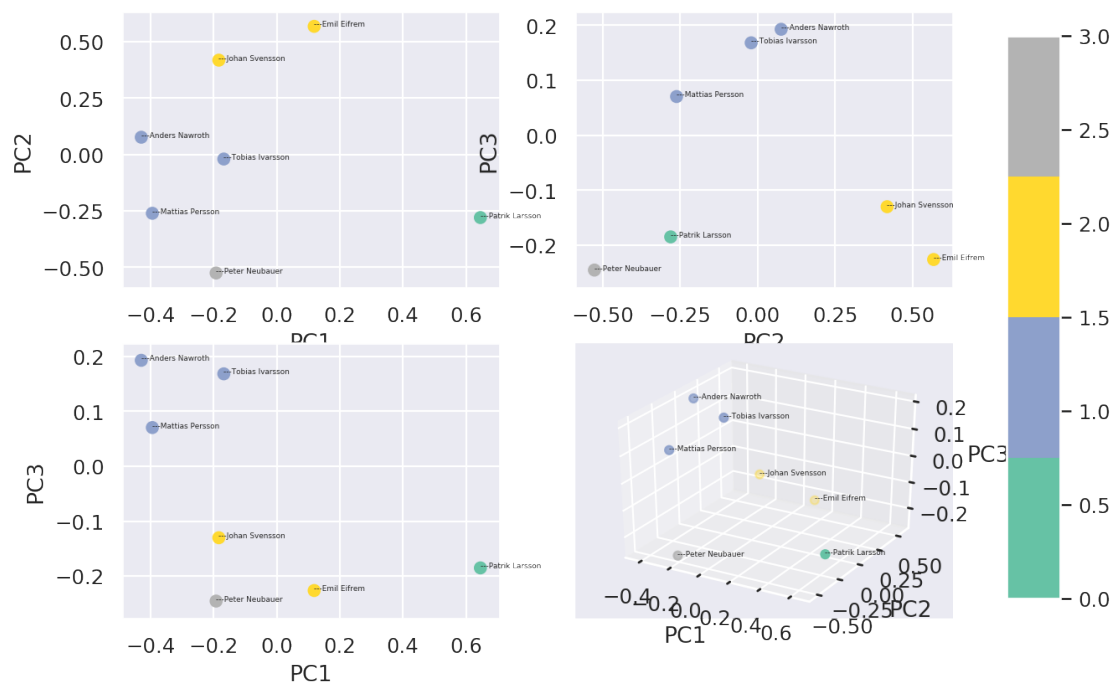
3.2.1 PCA

Cumulative sum of the explained variance

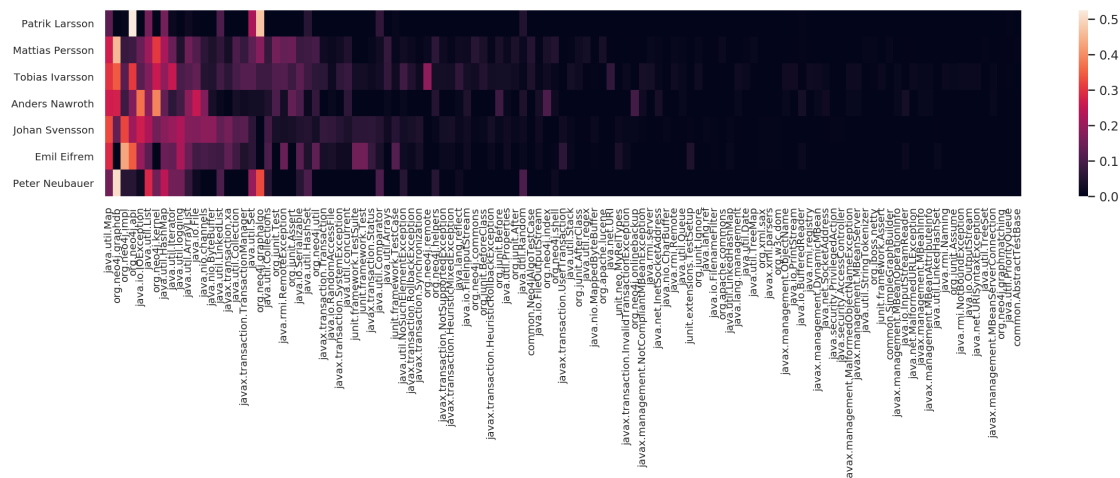


0.707367954956

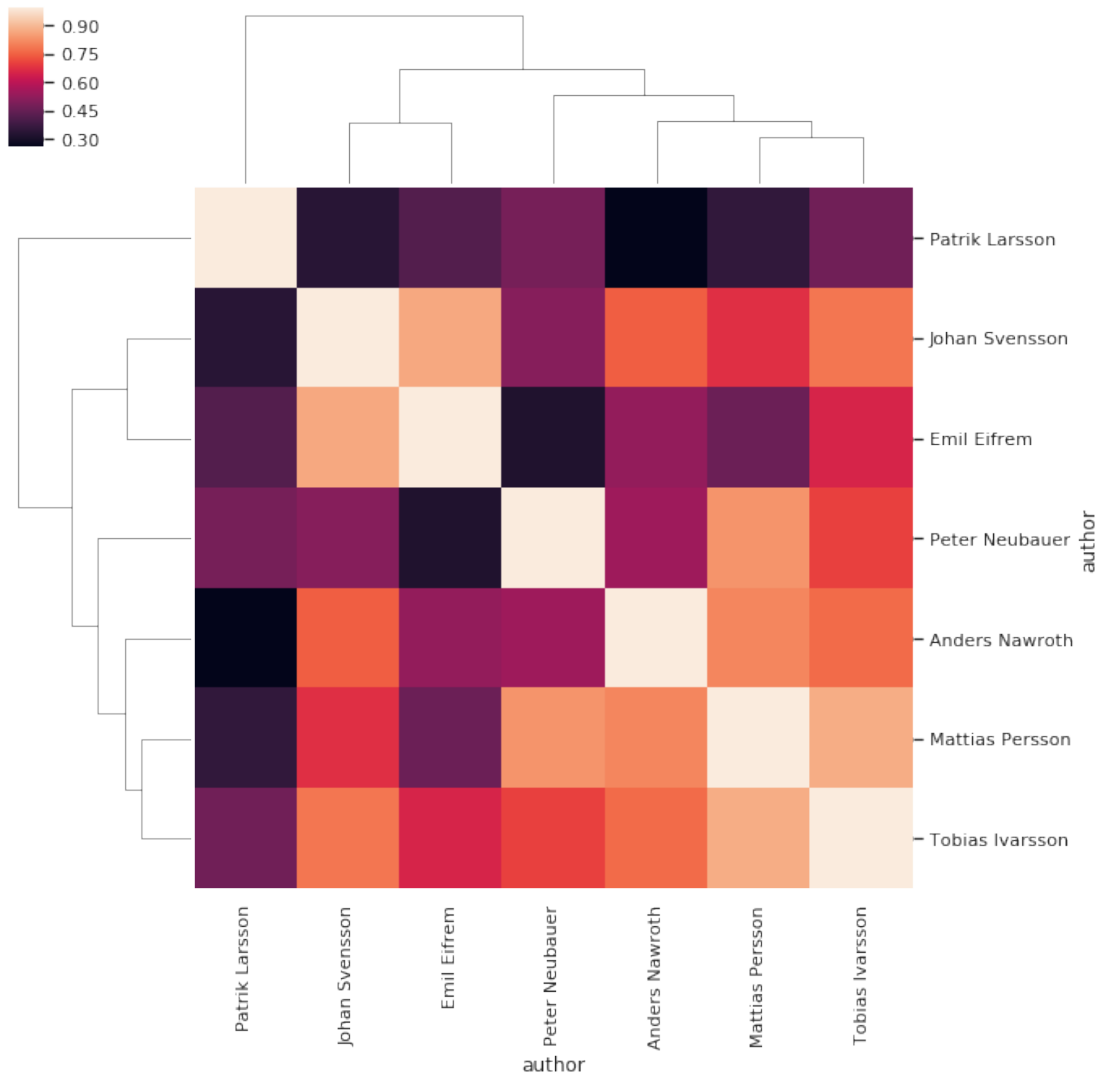
/opt/conda/lib/python3.6/site-packages/matplotlib/figure.py:2117: UserWarning: This figure was using constrained_layout==True, "



3.2.2 Ordered by clusters



3.2.3 Scalar product (similarity of work)

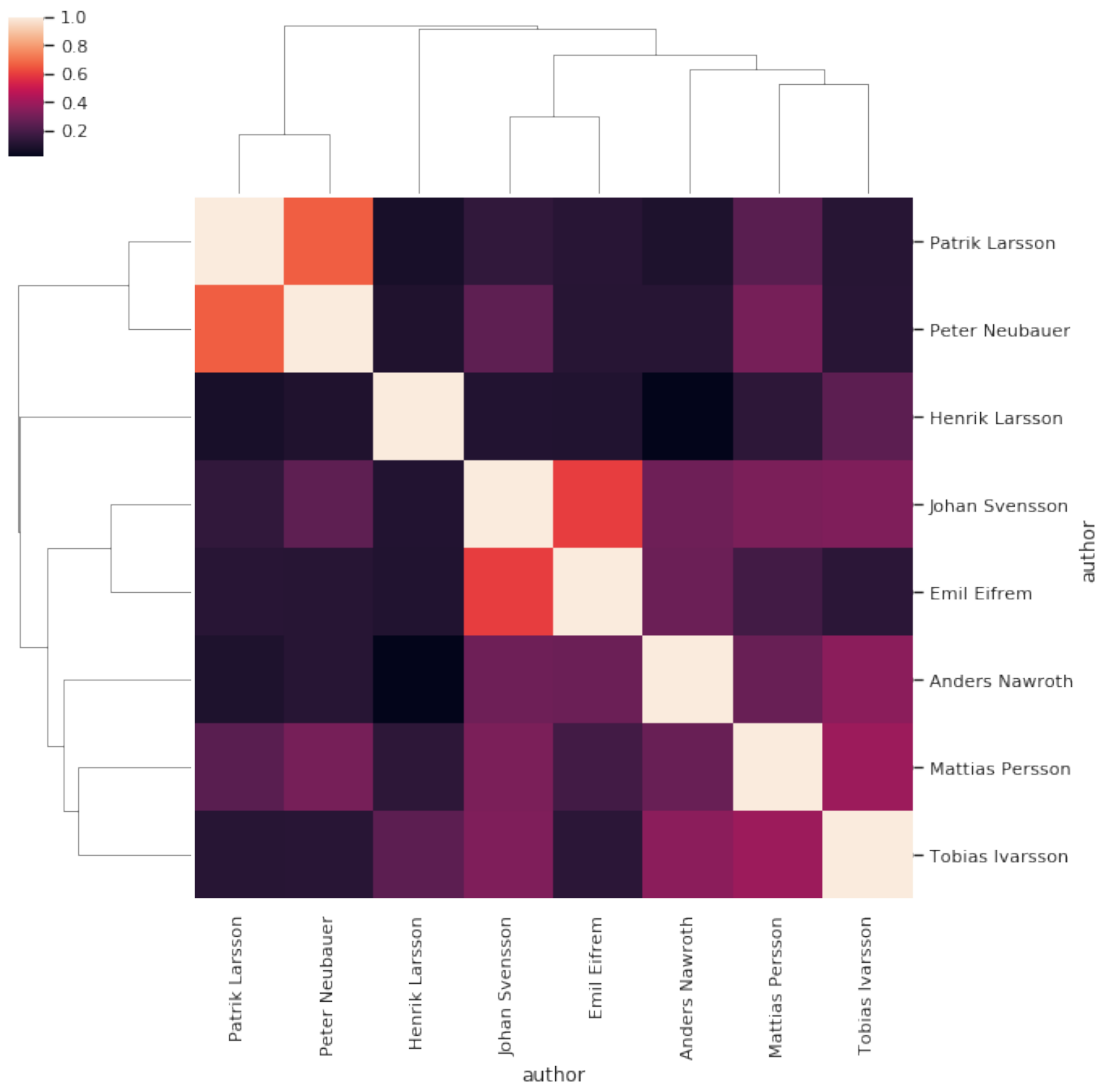


3.3 Normalisation expertise

First we divide by the total number of imports per library and then we normalize (L2) on the devs.
The goal is to have a vector that would represent it's expertise.



3.3.1 Scalar product (similarity of expertise)



3.3.2 Co-Clust

