

Some data from the database:

nObserved: 5001.

nSkipped: 3.

nAborted: 5.

nScanned: 9644.

nProjectsWithRegex: 4873.

nPythonFiles: 550532.

nFilesWithRegex: 50369.

## REPRESENTING A CORPUS OF REGEXES

1. How can we choose a small list of actual regexes to best represent all regexes within a corpus by usage?

- What regexes are most frequently used?
- What features are most frequently used?
- What regexes are most cloned?
- What behavioral clustering can be observed?
- What syntactic clustering can be observed?
- Various synthesis of a-e.

## CONTEXT AND CORPUS

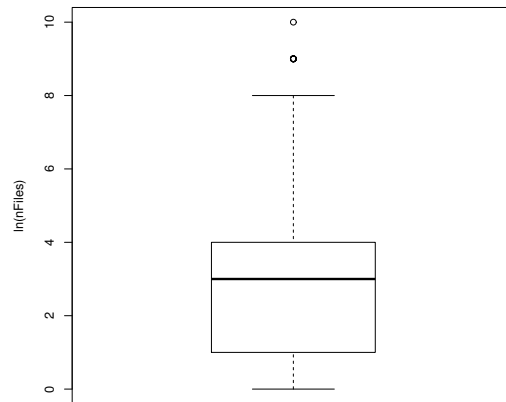


Figure 1: Example diagram: filesPerProject natural log