



Mahak Samim

A Corpus of Persian Academic Texts for Evaluating
Plagiarism Detection Systems

Document Collection



جستجو ...

ورود به سامانه

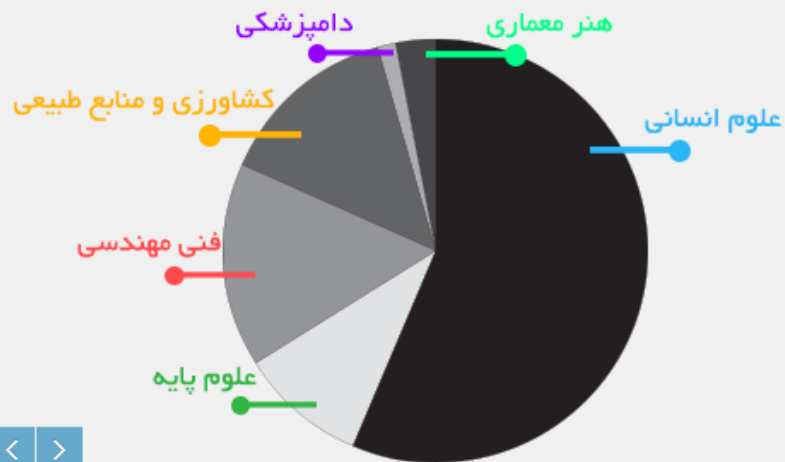
راهنمای داوران نشریات

تماس با ما

راهنمای سامانه

رتبه‌بندی نشریات

صفحه اصلی



دامپزشکی: ۱۳
علوم انسانی: ۵۹۳
علوم پایه: ۱۰۰
فنی مهندسی: ۱۶۵
کشاورزی و منابع طبیعی: ۱۴۶
هنر معماری: ۳۱

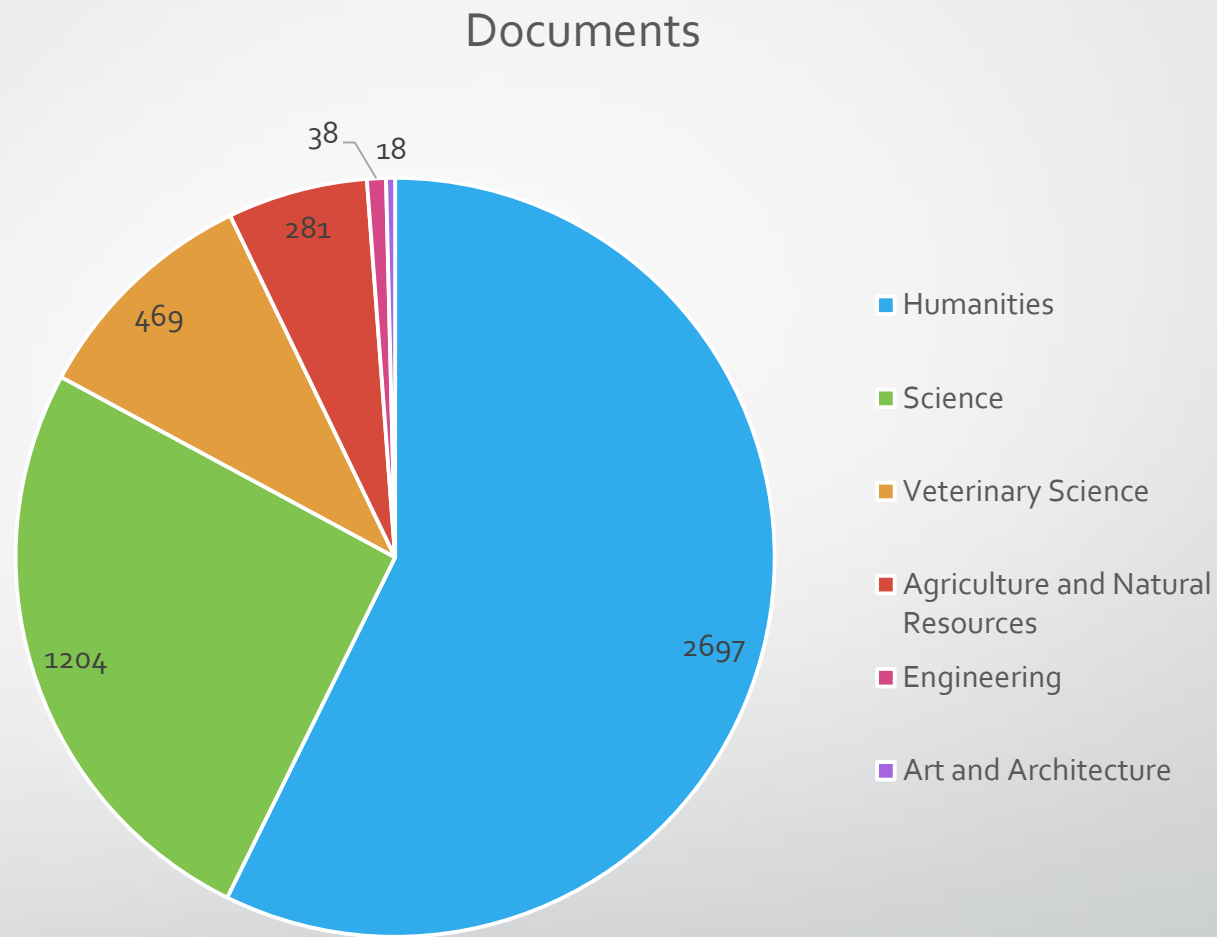
تعداد نشریات علمی
به تفکیک گروه

Documents Source

We crawled the websites of journals

Subject distribution

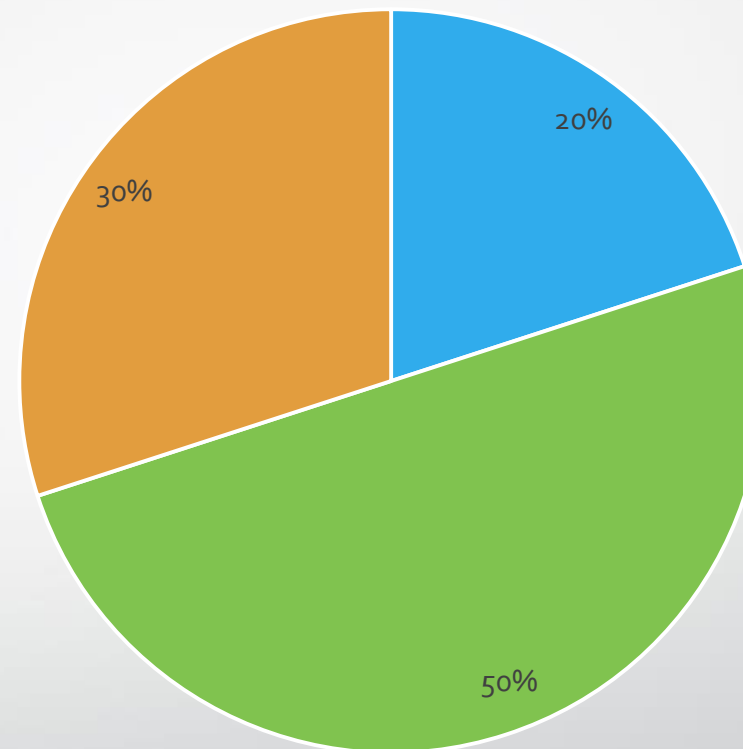
Statistics of number of documents per subject



Length distribution

Statistics of document lengths

Documents



■ short (1-3000 words) ■ medium (3000-6000 words) ■ long (6000-30000 words)

Source / suspicious documents

Source / suspicious distribution

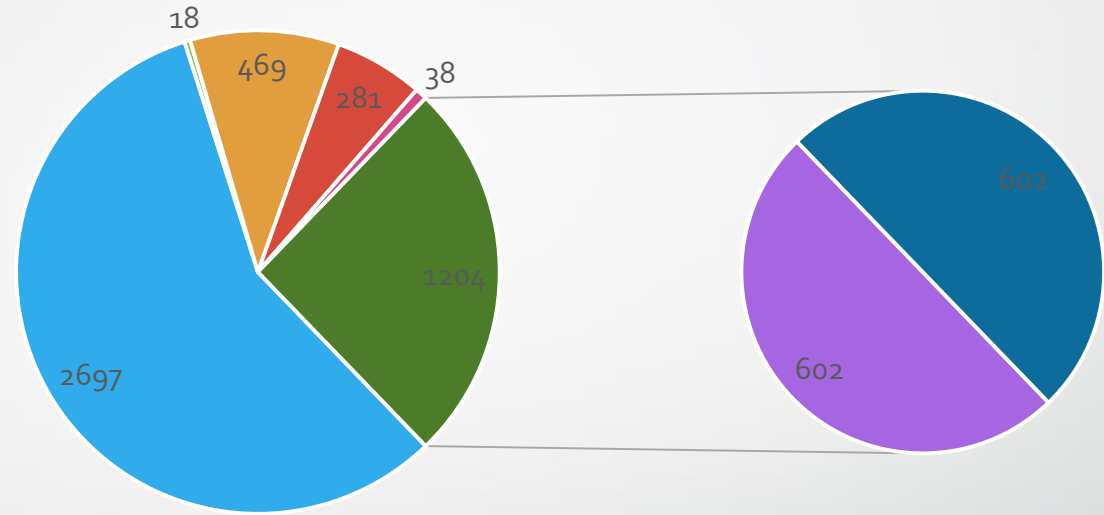
We follow PANs tradition of using half of the documents as source documents and half as suspicious documents.



Subject distribution in source / suspicious

Subjects of the papers were taken into consideration while dividing the collection into halves.

Documents



- Humanities
- Veterinary Science
- Engineering
- Science - Suspicious Documents
- Art and Architecture
- Agriculture and Natural Resources
- Science - Source Documents

Plagiarism per document

Suspicious documents without plagiarism

The documents without plagiarism allow to determine whether or not a detector can distinguish plagiarism cases from overlaps that occur naturally between random documents.

Documents



■ Source documents

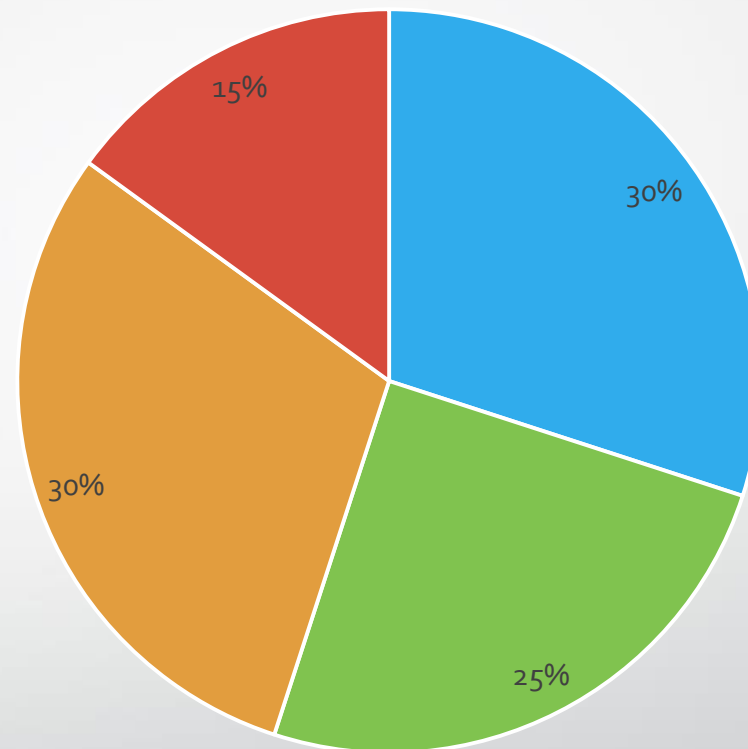
■ Suspicious documents with plagiarism

■ Suspicious documents without plagiarism

Suspicious documents with plagiarism

Statistics of plagiarism per document in the suspicious documents with plagiarism, i.e. 25 percent of the whole corpus.

Suspicious documents with plagiarism



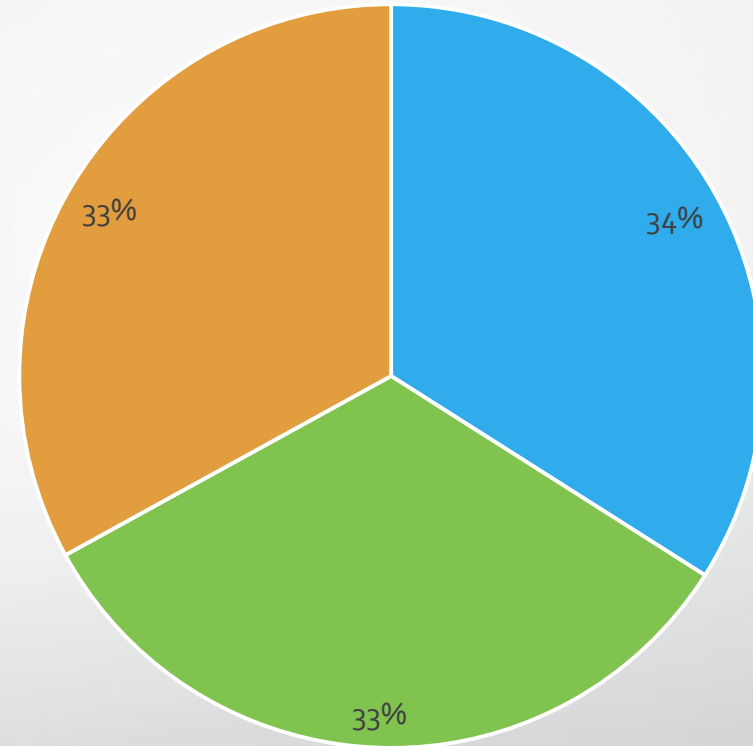
■ hardly (5%-20%) ■ medium (20%-50%) ■ much (50%-80%) ■ entirely (>80%)

Plagiarism case length

Statistics of lengths of plagiarism cases

Our corpus consists of a total of 5862 plagiarism cases with lengths between 50 and 5000 words. Long plagiarism cases may include more than one sentence.

Plagiarism cases



■ Short (50-150 words) ■ Medium (300-500 words) ■ Long (3000-5000 words)

Topic match

Intra-topic & inter-topic cases

Fifty percent of the plagiarism cases were made between papers with same topics (intra-topic cases) and fifty percent between papers with different topics (inter-topic cases).

Plagiarism cases



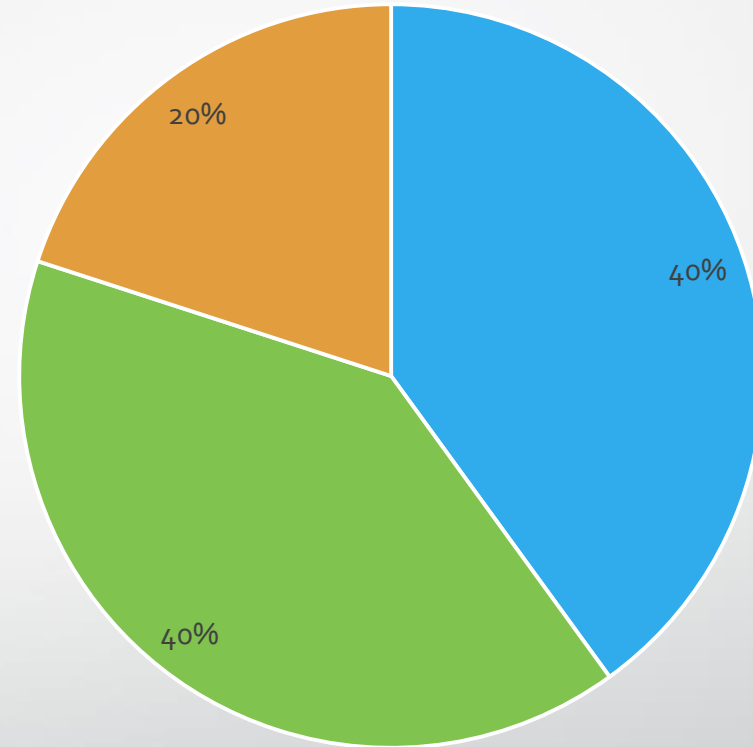
■ intra-topic cases ■ inter-topic cases

Obfuscation types

Types of obfuscation in plagiarism cases

40 percent of the plagiarism cases have no obfuscation. These cases are especially appropriate for evaluating intrinsic plagiarism detection.

Plagiarism cases

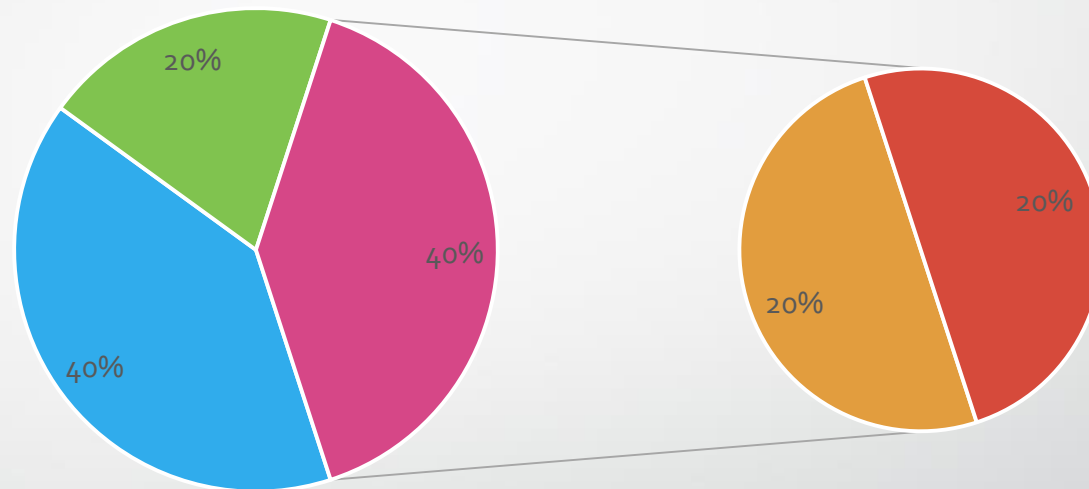


■ None ■ Random Text Operations ■ Semantic Word Variation

Random Text Operations

Random text operations are operations such as adding, deleting and substituting words, which are all done randomly.

Plagiarism cases

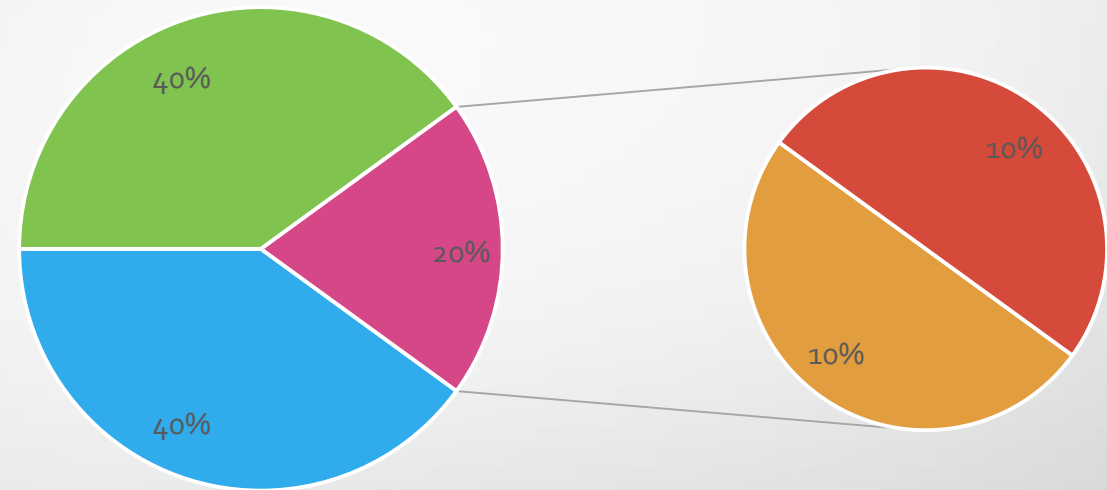


■ None ■ Semantic Word Variation ■ low obfuscation ■ high obfuscation

Semantic Word Variation

Semantic word variation, is the random substitution of words with their synonyms.

Plagiarism cases



■ None ■ Random Text Operations ■ Low obfuscation ■ High obfuscation

SUMMARY

- Mahak Samim is a plagiarism corpus which can be used for evaluating both intrinsic and external plagiarism detection systems.
- In order to preserve overall balance, many factors – plagiarism per document, plagiarism case length, topic match, obfuscation type, and obfuscation degree – were taken into consideration while preparing each plagiarism case.