

# Семинарское задание №3

## Page Rank

Мосиенко Константин Викторович

2019

В данном задании вам необходимо будет вычислить *ранг господина Пейджа* для страниц симпл википедии. Тут вам пригодятся результаты выполнения первого задания - ссылочный граф можно получить из скачанных документов. Вот только необходимо обратить внимание на качество полученной коллекции:

- Удалите ссылки вида:
  - /wiki/Help:, /wiki/Help\_talk:
  - /wiki/File:, /wiki/Media:, /wiki/MediaWiki:, /wiki/MediaWiki\_talk:
  - /wiki/Module:, /wiki/Talk:
  - /wiki/Category:, /wiki/Category\_talk:
  - /wiki/User:, /wiki/User\_talk:, /wiki/Special:
  - /wiki/Template:, /wiki/Template\_talk:
  - /wiki/Wikipedia:, /wiki/Wikipedia\_talk:
- Убедитесь, что коллекция содержит достаточное количество статей. Меньше 100000 точно не стоит брать.

Если не получается привести коллекцию к требуемому виду - лучше возьмите её у товарища. Вычислите pagerank (с нормализацией на максимальное значение) и отсортируйте все документы по его убыванию. Мне присылайте код и список документов с числами.