

Predicting the translations of red links

Olga Chernytska , Maksym Gontar , Kateryna Liubonko , and
Oleksandr Zaytsev

Ukrainian Catholic University
Faculty of Applied Sciences
Lviv, Ukraine
{chernytska,hontar,aloshkina,oleks}@ucu.edu.ua

Abstract

The goal of this project is to build the graph of all Wikipedia pages and learn from it to find pages that correspond to the non-existing pages referenced by red links in other languages. This would allow to create the so-called "red pages" by translating the existing ones.

1 Introduction

Wikipedia allows its users to create links to the pages that were not yet created. Such links become red and if you click on them, you get redirected to an empty page and asked to write the missing article. It is much easier for contributors to translate articles from other language than to write a new article from scratch.

We propose a solution of finding the potential translation of the missing page in other language. This would allow Wikipedia to automatically link red links to their translations in other languages and make the contribution process easier.

2 Data collection and preparation

We have downloaded the full dumps of English¹ and Ukrainian Wikipedia² articles from 20/06/2018.

Then we have parsed those articles with regular expressions to get outgoing red and blue links: the link article, text and position in the current article text. There was red links in English Wikipedia and red links Ukrainian Wikipedia. There was blue links in English Wikipedia and blue

¹<https://dumps.wikimedia.org/enwiki/20180620/>

²<https://dumps.wikimedia.org/ukwiki/20180620/>

links Ukrainian Wikipedia. This data is stored in `enwiki-20180620-pages-links.csv` and `ukwiki-20180620-pages-links.csv` files. Format of those files is following:

`id` - id of a page

`link_id` - id of a linked page

`link_pos` - position of link in a page markup text

`link_pos_perc` - relative position of link in a page markup text, range from 0 (at the beginning of page text) to 1 (at the end of page text)

`link_val` - title of a linked page

`link_text` - link text, if available

`is_red_link` - boolean, whether a link is red or not

Then we have downloaded Wiki interlanguage link records ³ and parsed out all interlingual links between En and Uk Wiki articles. There was 441928 pairs of Uk-En Wiki articles. This data is stored in `20180620-langlinks_uk_en.csv` file. Format of this file is following:

`id_uk` - id of a page in Uk Wiki

`id_en` - id of a linked page in En Wiki

From dumps we collected data about pages aliases (redirects) in En and Uk Wiki. The alias page is the page user can come upon by searching the article not by it's original name, but by it's alias name, then user is redirected to the original page. Alias data is important for our task, since links may lead not to the original article, but to it's redirect page. This data is stored in set of files:

`enwiki-20180620-id_alias_title_alias.csv`, `ukwiki-20180620-id_alias_title_alias.csv` with format:

`id_alias` - id of alias page

`title_alias` - title of alias page

`enwiki-20180620-id_alias_id_orig.csv`, `ukwiki-20180620-id_alias_id_orig.csv` with format:

`id_alias` - id of alias article

`id_orig` - id of original article

`enwiki-20180620-id_orig_title_alias.csv`, `ukwiki-20180620-id_orig_title_alias.csv` with format:

`id_orig` - id of original article

`title_alias` - title of alias pages

³<https://dumps.wikimedia.org/ukwiki/20180620/ukwiki-20180620-langlinks.sql.gz>

Also we composed a list of all pages in En and Uk Wiki. There were 5669865 in En and 581098 articles Uk Wiki. This data is stored in enwiki-20180620-id_name.csv, ukwiki-20180620-id_name.csv files. Format of those files is following:

id - id of a page

title - title of a page

length - length of a page markup text

Besides that, we did a statistical analysis for red links, we calculated how many times each red link was used and saved results in the enwiki-20180620-red_name_count.csv, ukwiki-20180620-red_name_count.csv files with format:

link_title - title of a red link

in_count - number of times it was used

Also we calculated how many red links there are with a certain number of use, and stored results into the enwiki-20180620-red_count_by_count.csv, ukwiki-20180620-red_count_by_count.csv files with format:

count - number of red link usage

in_count - number of this count case

For example, in En Wiki red link was used only once for 4354094, and twice for 811612 times.

For every eng red link in the matrix we calculate similarity (using similarity metrics selected on the previous step) to those ukr articles that do not have eng version. The most similar ukr article is the one that corresponds to this red link. Possibly, we will find several ukr articles with the same similarities, so add them as candidates for further preprocessing.

3 Baseline assumption

$$p(a)p(b)p(c) < p(a)p(b)$$