

# TeDDi: Crawling issues

Olga Pelloni (Sozinova)

October 1, 2022

In this document, I describe the known issues in the crawled TeDDi files. In addition to the issues, I suggest potential ways of fixing them and report on the bugs I already fixed.

## 1 Mixed languages

It is not always purely one language per text file. In case of OpenSubtitles, we notice sometimes double subtitles in English in parallel to the main language (often in case of Mandarin).

```
33 你就安定下来了吗 Settled down, have you?  
34 没错 Yeah.  
35 你们做过吗 Have you done it yet?  
36 这是我同事亨利 This is Henry from work.  
37 你总是这样 You always did that.  
38 从不出轨 只玩暧昧 Never had affairs, just obsessions.  
39 我是兰斯·沙利文 合作教育科科长 Lance Sullivan, Head of Corporate and Education.  
40 -你好 小心点  
41  
42 -你想干什么  
43  
44 他偷了我的东西 He just stole from me.  
45 我没想偷窃 I wasn't going to steal it.  
46 你能跟我们一起回去 只打炮不麻烦 We could go back to our house.  
47  
48 跟我们俩做 你就能住一晚上 With the both of us.  
49  
50 兰斯 我们该停手 Lance, we should stop!  
51 有个男人在我家里 There is a...  
52  
53 没必要伤害他 There's no need to hurt him!  
54 我能住这吗 can I stay?  
55 每对恩爱的情侣都有秘密 Every happy couple has a secret.  
56 他们努力对此视而不见 They try to ignore it, they turn a blind eye,  
57 但其实他们的感情都岌岌可危 but every happy couple is in danger,  
58 因为每对恩爱的情侣都有手机 because every happy couple has a mobile phone.  
59 我发明了一个测试 I've invented a test.  
60 起名为“亨利测试” I call it the Henry test.  
61 测试流程是这样的 It works like this.  
62 如果你觉得你正处于恋爱的甜蜜期 If you think you're happy in your relationship,
```

Figure 1: Several languages in one text file.

In case of Gutenberg texts, we notice some texts to be dictionaries, thus, containing text data in two languages as well. We try to avoid them while crawling.

**Potential solution:**

- **Potential language identification tool** <https://polyglot.readthedocs.io/en/latest/>

### 1.1 Solving the issue

I've tried using the Polyglot identification tool. It returns the relative number of up to 3 languages identified in a text. I set the following threshold: language 1 less than 99%, language 2 more than 5%. The results are written down in a text file.

Most of the texts seem to mix together with English, but not always. There are less obvious cases, such as Chinese - Spanish, Finnish - Swedish, Chinese simplified - Chinese traditional. Next step is to look closer at this results file and to decide what to do with the text files: keep the most present language inside a file or remove the file. Note that mixed languages can appear also in a single line, which means that if we want to clean the file, we have to identify a language token by token, which might not work that well.

Some files in the results belong to low-resource languages, so they probably have languages identified by mistake. It is better to filter out my text file only to the crawled languages.

## 2 Non-linguistic noise

Any non-linguistic items, such as tags (HTML, XML), some Latex-like commands, etc.

[illegible]

Figure 2: Non-linguistic noise example.

**Potential solution:**

- Identify the least frequent symbols or N-grams, look at the context of those symbols

## 2.1 Solving the issue

I've implemented an algorithm, which creates a frequency dictionary of the characters in a text file, filters out all the characters which have a frequency lower than 20.

Then, we search for sequences in a text of length 5, where every character has a frequency lower than 20. These sequences are potential non-linguistic noise strings.

This is a results file, where I write down a filename and an identified sequence. Many of them point out to mixed languages (mostly parts of English words), rather than to noise. A potential way to filter this file, is to look at the sequences containing non-letter symbols (punctuation marks, slashes, etc.). However, some sequences containing letters, and yet belong to noise, e.g.: `heb_nfi_1064.txt, 080ff`, which is a color code, included in an HTML tag. Quotation marks might help with such cases.

Once we find a way to filter out the actual noise in the results file, we can dig into that file and look at the context, in order to clean the noise. This, most probably, will require manual cleaning (the context can be the whole line, part of the line, whole tag or part of a tag, etc), or, maybe, the most frequent noise items can be cleaned out automatically (e.g. `color = "0FFFFFF"` and alike).

## 3 Malformed sentences

Strange appearance of the letters, Unicode problems (symbols are not combined together as they should). Example of a Thai text:



Figure 3: Malformed sentences in Thai.

### Potential solution:

- Check the Unicode sequence somehow (?), Unicode decomposition

## 3.1 Solving the issue

After looking closer at the suspicious Thai files, it turns out that Unicode is fine there, however, the text itself contains nonsense words and many control characters (non-printing characters), which is probably a sign of a strange encoding or machine, where the text was created. We found such files by filtering the segmentations obtained with BPE-200, length of words > 30 (with very regular segments' lengths).

The list of corrupted files: `tha_nfi_67.txt, tha_nfi_688.txt, tha_nfi_915.txt`. These files are to be deleted.

The file `tha_nfi_46.txt` has to be filtered from noise (tag `"/c.bg_transparent"`), otherwise, the text is fine.

**Status: DONE (pull request accepted)**

## 4 Meta in the text body

In some text files from the Gutenberg Project, there is meta information inside the text body. It often includes either the start or the ending of the Gutenberg copyright information, e.g. \*\*\*START OF THE PROJECT GUTENBERG EBOOK IT WAS MARLOWE\*\*\*.

```
65 [21] Unfortunately for mankind, these were the last words pronounced
66 by this great Painter from the Academical chair. He died about
67 fourteen months after this Discourse was delivered.
68
69 ***END OF THE PROJECT GUTENBERG EBOOK SIR JOSHUA REYNOLDS' DISCOURSES***
70 ***** This file should be named 47610-h.htm or 47610-h.zip *****
71 This and all associated files of various formats will be found in:
72 http://www.gutenberg.org/4/7/6/1/47610
73
74 Updated editions will replace the previous one--the old editions will
75 be renamed.
76 Creating the works from print editions not protected by U.S. copyright
77 law means that no one owns a United States copyright in these works,
78 so the Foundation (and you!) can copy and distribute it in the United
79 States without permission and without paying copyright
80 royalties. Special rules, set forth in the General Terms of Use part
81 of this license, apply to copying and distributing Project
82 Gutenberg-tm electronic works to protect the PROJECT GUTENBERG-tm
83 concept and trademark. Project Gutenberg is a registered trademark,
84 and may not be used if you charge for the eBooks, unless you receive
85 specific permission. If you do not charge anything for copies of this
86 eBook, complying with the rules is very easy. You may use this eBook
```

Figure 4: Meta related information in the text body.

### Potential solution:

- Find the key phrases (Start/End of the Project Gutenberg...), remove the text before START and after END.

**Files:** eng\_fic\_97.txt, eng\_fic\_54.txt, eng\_fic\_68.txt, eng\_fic\_57.txt, eng\_fic\_92.txt, eng\_fic\_51.txt, eng\_fic\_50.txt, eng\_fic\_23.txt, eng\_fic\_20.txt, eng\_fic\_32.txt, eng\_fic\_17.txt, eng\_fic\_15.txt, eng\_fic\_8.txt, eng\_fic\_12.txt, eng\_fic\_48.txt, eng\_fic\_89.txt, eng\_fic\_67.txt, spa\_fic\_63.txt, spa\_fic\_66.txt, spa\_fic\_67.txt, spa\_fic\_108.txt, spa\_fic\_104.txt, spa\_fic\_23.txt, spa\_fic\_111.txt, spa\_fic\_35.txt, spa\_fic\_31.txt, spa\_fic\_96.txt, fra\_fic\_84.txt, fra\_fic\_79.txt, fra\_fic\_78.txt, fra\_fic\_83.txt, fin\_fic\_9.txt, fin\_fic\_120.txt, fin\_fic\_136.txt, fin\_fic\_10.txt, fin\_fic\_38.txt, fin\_fic\_123.txt, fin\_fic\_132.txt, fin\_fic\_118.txt, fin\_fic\_142.txt, fin\_fic\_71.txt, fin\_fic\_73.txt, fin\_fic\_140.txt, fin\_fic\_62.txt, fin\_fic\_151.txt, fin\_fic\_147.txt, fin\_fic\_153.txt, fin\_fic\_61.txt, fin\_fic\_79.txt, fin\_fic\_160.txt, fin\_fic\_47.txt, fin\_fic\_91.txt, fin\_fic\_159.txt, fin\_fic\_80.txt, fin\_fic\_82.txt, fin\_fic\_83.txt, fin\_fic\_27.txt, fin\_fic\_115.txt, fin\_fic\_101.txt, fin\_fic\_129.txt, fin\_fic\_19.txt, fin\_fic\_21.txt, fin\_fic\_7.txt, fin\_fic\_6.txt, fin\_fic\_107.txt, fin\_fic\_4.txt, fin\_fic\_23.txt, fin\_fic\_138.txt, fin\_fic\_110.txt, deu\_fic\_134.txt, deu\_fic\_36.txt, deu\_fic\_18.txt, deu\_fic\_118.txt, deu\_fic\_54.txt, deu\_fic\_96.txt, deu\_fic\_43.txt, deu\_fic\_151.txt, deu\_fic\_52.txt, deu\_fic\_44.txt, deu\_fic\_74.txt, deu\_fic\_17.txt, deu\_fic\_111.txt

## 4.1 Solving the issue

Removing meta info file by file. There is sometimes a long text of a full Gutenberg license. It might be a good idea to rerun the reports after such cleanings (line count, tokens).

In some files I modify the meta field `short_description`, if the author and title were missing; the same applies to the field `year_published`.

General observation: the issue often appeared in the samples, where we sample text in the end of a novel, and then continue from the beginning. The crawler could be improved by taking into account this scenario, and cleaning the license in the end and the meta in the beginning. However, this was also the case in the "whole" type of samples.

Note: some files have a tab in the `short_description` field. Maybe it needs to be checked.

Files done:

`eng_fic_97.txt`, `eng_fic_54.txt`, `eng_fic_68.txt`, `eng_fic_57.txt`, `eng_fic_92.txt`, `eng_fic_51.txt`,  
`eng_fic_50.txt`, `eng_fic_23.txt`, `eng_fic_20.txt`, `eng_fic_32.txt`, `eng_fic_17.txt`, `eng_fic_15.txt`,  
`eng_fic_8.txt`, `eng_fic_12.txt`, `eng_fic_48.txt`, `eng_fic_89.txt`, `eng_fic_67.txt`

Pull request done, will be continued with the other files.

### UPD 30.09.22:

Finished all the files (pull request finished too):

Spanish:

`spa_fic_63.txt`, `spa_fic_66.txt`, `spa_fic_67.txt`, `spa_fic_108.txt`, `spa_fic_104.txt`, `spa_fic_23.txt`,  
`spa_fic_111.txt`, `spa_fic_35.txt`, `spa_fic_31.txt`, `spa_fic_96.txt`

Finnish:

`fin_fic_9.txt`, `fin_fic_120.txt`, `fin_fic_136.txt`, `fin_fic_10.txt`, `fin_fic_38.txt`, `fin_fic_123.txt`,  
`fin_fic_132.txt`, `fin_fic_118.txt`, `fin_fic_142.txt`, `fin_fic_71.txt`, `fin_fic_73.txt`, `fin_fic_140.txt`,  
`fin_fic_62.txt`, `fin_fic_151.txt`, `fin_fic_147.txt`, `fin_fic_153.txt`, `fin_fic_61.txt`, `fin_fic_79.txt`,  
`fin_fic_160.txt`, `fin_fic_47.txt`, `fin_fic_91.txt`, `fin_fic_159.txt`, `fin_fic_80.txt`, `fin_fic_82.txt`,  
`fin_fic_83.txt`, `fin_fic_27.txt`, `fin_fic_115.txt`, `fin_fic_101.txt`, `fin_fic_129.txt`, `fin_fic_19.txt`,  
`fin_fic_21.txt`, `fin_fic_7.txt`, `fin_fic_6.txt`, `fin_fic_107.txt`, `fin_fic_4.txt`, `fin_fic_23.txt`,  
`fin_fic_138.txt`, `fin_fic_110.txt`

German:

`deu_fic_134.txt`, `deu_fic_36.txt`, `deu_fic_18.txt`, `deu_fic_118.txt`, `deu_fic_54.txt`, `deu_fic_96.txt`,  
`deu_fic_43.txt`, `deu_fic_151.txt`, `deu_fic_52.txt`, `deu_fic_44.txt`, `deu_fic_74.txt`, `deu_fic_17.txt`,  
`deu_fic_111.txt`

Comments to the files:

- `fra_fic_84.txt`, `fra_fic_79.txt`, `fra_fic_78.txt`, `fra_fic_83.txt`: look fine, no changes done; `gutenberg` is used as a word;
- `spa_fic_108.txt`, `spa_fic_31.txt`: a lot of English comments and there is a Spanish-English dictionary at the end of texts;
- `deu_fic_44.txt`: looks like Old German text or as dialectal German variety;

- Maybe it's better to explore NBSP symbols (like whitespace). I removed them in Spanish files (substituted with whitespace), in the other language files there are many NBSP symbols, I didn't change them. I think they need a check, how they would influence preprocessing in Python.

**Status: DONE (pull request waiting for approval)**

## 5 Missing meta information

In some text files from the Gutenberg Project, there is no value in the `short_description` meta field, which usually includes the author name and the title.

```

1 # language_name_wals: English
2 # language_name_glottol: English
3 # iso639_3: eng
4 # year_composed: NA
5 # year_published: 2006
6 # mode: written
7 # genre_broad: fiction
8 # genre_narrow: general_fiction
9 # writing_system: Latn
10 # special_characters: NA
11 # short_description: NA
12 # source: http://www.gutenberg.org/files/8600/8600-h/8600-h.htm
13 # copyright_short: http://www.gutenberg.org
14 # copyright_long: http://www.gutenberg.org
15 # sample_type: part
16 # comments: NA
17
18 Gervaise flushed. She thought she would have felt less shame if he had
19 taken her in his arms and kissed her. Goujet was an odd fellow, proposing
20 to elope, just the way it happens in novels. Well, she had seen plenty of
21 workingmen making up to married women, but they never took them even as
22 far as Saint-Denis.

```

Figure 5: Missing meta information (author name and title).

### Potential solution:

- Search for the author name and the title at the HTML page

### Files:

- eng\_fic\_108.txt
- eng\_fic\_106.txt
- eng\_fic\_107.txt
- spa\_fic\_17.txt
- spa\_fic\_22.txt
- spa\_fic\_21.txt

- spa\_fic\_20.txt
- spa\_fic\_18.txt
- spa\_fic\_19.txt

## 5.1 Solving the issue

Found the name of author and the title on the web pages. These samples were taken from two novels: Émile Zola, L'Assommoir and Leopoldo Alas, La Regenta. Copy-pasted this info manually.

**Status: DONE (pull request accepted)**

## 6 Manually added data

Consistency of tags (what goes to which layer, glosses or not in the segmentation layer), usage of symbols (asterix as a sound in Piraha), other consistency issues.

## 7 Status on 30.09.2022

Here is an update about the TeDDi cleaning status as of 30.09.2022.

- **Mixed languages:** I did not change anything in the TeDDi github repo. What should be done here next: The file should be filtered; 1) leave there only the crawled languages (remove low-resouce languages that we collected manually); 2) find empirically a threshold for confidence and filter according to that. For example, the file kor\_nfi\_271.txt has Korean with confidence 52 and English with confidence 47. Most likely, this is a good candidate to remove from the corpus. On the other hand, kor\_nfi\_275.txt has Korean with confidence 94 and English with confidence 5; probably, such file can be left unchanged.
- **Non-linguistic noise:** I did not change anything in the TeDDi github repo. What should be done here next: The results file has to be manually processed and filtered as well. This task will take more time than the mixed languages problem. First, one needs to leave only crawled languages as in the mixed languages problem. Then, one needs to go through all examples and identify which ones belong to noise and which do not. Some examples of how this could be automatized are given above in the respective section.
- **Malformed sentences:** Done in Thai, pull request accepted.
- **Meta in the text body:** Done in 4 languages, pull request waiting for approval. See comments left above in the respective section for potential things to check later.
- **Missing meta information:** Done in English and Spanish, pull request accepted.

- **Manually added data:** To be done later, probably after consulting Chris and Steve.