

Manual for the Tagging Optimiser

Version 1.0

Author:

Martin Weisser

September 2018

Contents

1	Introduction	1
2	Overview of the Tool	1
3	Selecting Input Files.....	2
4	Optimising Tagged Files	2
5	Viewing/Editing Optimised or Input Files.....	3
6	The Optimised Tagset (Ver. 1.0)	3
7	References	9

1 Introduction

The Tagging Optimiser is a tool designed to allow corpus users to enhance and improve the output of freeware taggers. It takes one or more PoS-tagged files as input, tries to correct some of the errors produced by the probabilistic tagging systems that tend to form the basis for these taggers, and diversifies the tag set used to allow for more fine-grained grammatical analyses. At the same time, it makes the tags in the tagset more readable by expanding the traditional minimalistic tags, thereby making them more mnemonic and easier to understand.

It has been developed for and tested with the output from 3 different taggers, the TreeTagger (Schmid 1994), The Stanford POS Tagger (Toutanova et al. 2003), and my own Simple PoS Tagger (Weisser 2014), including the version integrated into my Simple Corpus Tool (Weisser 2018a).

2 Overview of the Tool

The layout of the tool is quite simple. It consists of a workspace for handling input and output files, as well as an editor pane for viewing/editing in- or output files. Figure 1 shows an image of the Tagging Optimiser startup screen.

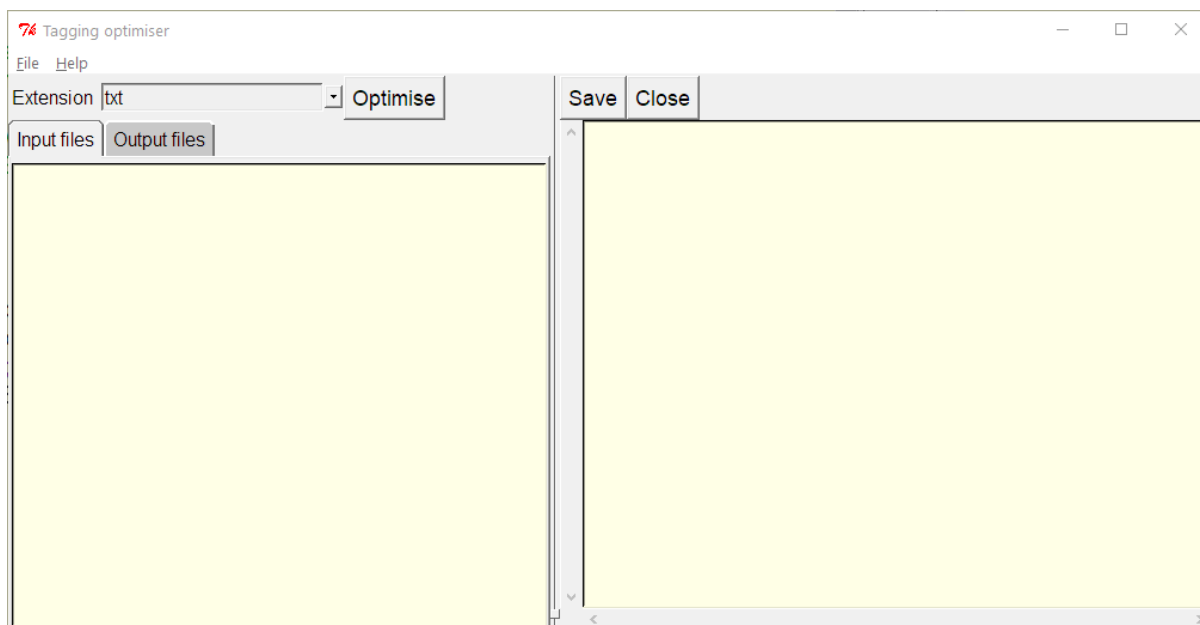


Figure 1 – The Tagging Optimiser Startup screen

The input/output workspace on the left contains 2 tabs, one where input files will be listed, and the other where any output files will be shown after optimisation. The editor pane on the right

consists of the editor window itself and 2 buttons, one for saving a file opened in the editor and the other for closing it.

3 Selecting Input Files

To select input files for optimising or editing, you first need to set an extension for the input files via the dropdown list above the input/output workspace. This option defaults to 'txt', but 'xml' can also be chosen, or user-defined extensions typed into the selection box.

There are 2 options for selecting files, either by choosing a whole directory/folder, or individual single or multiple files. To achieve the former, either use 'File → Select source → Directory' or press 'Ctrl + Alt + d' or 'F2' on the keyboard, then navigate to the relevant folder and press enter. All files that fit the extension defined earlier will then be listed under the 'Input files' tab in the input/output workspace.

For individual files, select 'File → Select source → File(s)' or press 'Ctrl + Alt + f'. After navigating to the appropriate folder, you can then either select a single file or multiple ones, using the appropriate multi-selection mechanisms of your operating system. As before, all selected files will be listed in the input workspace.

Files listed in either workspace can also be selected and removed by pressing 'Del' in order to exclude them from processing.

The default file location for data files in the Tagging Optimiser is the 'data' sub-folder of the program folder, but of course you can also select any folder on your system that contains your data to be optimised.

4 Optimising Tagged Files

Once one or more files have been loaded into the input workspace, you can optimise them by clicking the 'Optimise' button. Provided you have write-access to the folder your data originates from, a new sub-folder called 'optimised' will then automatically be created and the optimised files output there, as well as listed under the 'Output files' tab, which will also automatically be activated. Should the latter not occur, it's most likely that you don't have write-access and you need to move your files to a different folder that you can write to.

5 Viewing/Editing Optimised or Input Files

Files loaded into the input or output workspace can be viewed or edited by double-clicking the filename in the respective workspace tab. They will then open in the editor window on the right. Any changes made to the file can be saved either via the ‘Save’ button or by pressing ‘Ctrl + s’. You can close a file by clicking on the ‘Close’ button or pressing ‘Ctrl + w’.

If you want to edit a number of files in succession, there’s no need to close any open file first, though, as double-clicking a new file will automatically close the file before opening the new one, also prompting you to save or discard any unsaved changes. The same check will also be performed if you close the program window.

To facilitate the editing process, the editor window also has a number of keyboard shortcut bindings. The standard copy (Ctrl + c), cut (Ctrl + x), and paste operations (Ctrl + v) are automatically implemented via the editor widget, but, in addition, you can also use ‘Ctrl + f’ to find expressions and ‘F3’ to find the next occurrence of the same expression (provided it is still highlighted).

6 The Optimised Tagset (Ver. 1.0)

The table below list the entries in the optimised tagset. The right-hand column (labelled Example/Word(s)) either contains examples drawn from sample files of the FLOB or FROWN corpora I’ve used for testing, or, if the particular tag has a closed set of representatives, the complete list.

Please note that the final tagset for version 1.0 is slightly different from the one presented in Weisser (2018b), especially with regard to the tagging of participles, which are now labelled with the ‘suffixes’ *presPart* instead of *ing* and *pastPart* instead of *ppart* to make them easier to identify and search for in concordances.

Label	Explanation	Example/Word(s)
Det~def	determiner, definite (underspecified for number)	<i>the</i>
Det~indef	determiner, indefinite	<i>a, an</i>
Det~demSing	determiner, demonstrative singular	<i>this, that</i>

Label	Explanation	Example/Word(s)
Det~demPl	determiner, demonstrative plural	<i>these, those</i>
Noun~Sing	noun, general singular	<i>ability, creature, meeting</i>
Noun~Pl	noun, general plural	<i>consequences, people, workshops</i>
Noun~propSing	noun, proper singular	<i>English, Wales</i>
Noun~propPl	noun, proper plural	<i>Proms, Republicans</i>
Noun~propDaySing	noun, proper name of weekday singular	<i>Friday</i>
Noun~propDayPl	noun, proper name of weekday plural	<i>Fridays</i>
Verb~BEbase	verb, BE base form	<i>be</i>
Verb~BEpresPart	verb, BE present participle	<i>being</i>
Verb~BE1Sing	verb, BE 1 st person singular	<i>am</i>
Verb~BE3Sing	verb, BE 3 rd person singular	<i>is</i>
Verb~BENot3Sing	verb, BE not 3 rd person singular (underspecified for number & person)	<i>are</i>
Verb~BESingpast	verb, BE singular past tense	<i>was</i>
Verb~BEPlpast	verb, BE plural past tense	<i>were</i>
Verb~BEen	verb, BE past participle	<i>been</i>
Verb~HAVEbase	verb, HAVE base form	<i>have</i>
Verb~HAVEpresPart	verb, HAVE present participle	<i>having</i>
Verb~HAVE3Sing	verb, HAVE 3 rd person singular	<i>has</i>
Verb~HAVENot3Sing	verb, HAVE not 3 rd person singular	<i>have</i>
Verb~HAVEpast	verb, HAVE past tense	<i>had</i>

Label	Explanation	Example/Word(s)
Verb~HAVEPerf	verb, HAVE perfect tense	<i>have</i> + Verb~pastPart
Verb~HAVEPastPerf	verb, HAVE past perfect	<i>had</i> + Verb~pastPart
Verb~DObase	verb, DO base form	<i>do</i>
Verb~DOpresPart	verb, DO present participle	<i>doing</i>
Verb~DO3Sing	verb, DO 3 rd person singular	<i>does</i>
Verb~DOnot3Sing	verb, DO (underspecified for number & person)	<i>do</i>
Verb~DOpast	verb, DO past tense	<i>did</i>
Verb~DOen	verb, DO past participle	<i>done</i>
Verb~MOD	verb, modal	<i>may, can</i>
Verb~base	verb, general base form	<i>bring, go, tell</i>
Verb~presPart	verb, general present participle	<i>improving, making</i>
Verb~3Sing	verb, general 3 rd person singular	<i>asks, gives</i>
Verb~not3Sing	verb, general not 3 rd person singular (underspecified for person & number)	<i>call, say</i>
Verb~past	verb, general past tense	<i>asked, gave</i>
Verb~pastPart	verb, general past participle	<i>(had) looked, (had been) placed</i>
Adj	adjective, base form	<i>bumpy, glamorous</i>
Adj~comp	adjective, comparative form	<i>greater, less</i>
Adj~sup	adjective, superlative form	<i>best, clearest</i>
Num~card	number, cardinal	<i>eight, 1986</i>
Num~ord	number, ordinal	<i>second, 15th</i>
Quant	quantifier	<i>all, few, many</i>

Label	Explanation	Example/Word(s)
Adv	adverb, base form	<i>already, slightly</i>
Adv~comp	adverb, comparative form	<i>earlier, more</i>
Adv~sup	adverb, superlative form	<i>earliest, most</i>
Neg	negation operator	<i>not, n't</i>
Pron~1SingSub	pronoun, 1 st person singular subject	<i>I</i>
Pron~1SingObj	pronoun, 1 st person singular object	<i>me</i>
Pron~1PlSub	pronoun, 1 st person plural subject	<i>we</i>
Pron~1PlObj	pronoun, 1 st person plural object	<i>us</i>
Pron~2	pronoun, 2 nd person (underspecified for number)	<i>you</i>
Pron~3SingMSubj	pronoun, 3 rd person singular male subject	<i>he</i>
Pron~3SingMObj	pronoun, 3 rd person singular male object	<i>him</i>
Pron~3SingFSubj	pronoun, 3 rd person singular female subject	<i>she</i>
Pron~3SingFObj	pronoun, 3 rd person singular female object	<i>her</i>
Pron~3SingN	pronoun, 3 rd person singular neuter	<i>it</i>
Pron~3PlSub	pronoun, 3 rd person plural subject underspecified for gender	<i>they</i>
Pron~3PlObj	pronoun, 3 rd person plural object	<i>them</i>
Pron~1SingRefl	pronoun, 1 st person singular reflexive	<i>myself</i>

Label	Explanation	Example/Word(s)
Pron~1PIRefl	pronoun, 1 st person plural reflexive	<i>ourselves</i>
Pron~2Refl	pronoun, 2 nd person reflexive underspecified for number	<i>yourself</i>
Pron~2PIRefl	pronoun, 2 nd person plural reflexive	<i>yourselves</i>
Pron~3SingMRefl	pronoun, 3 rd person singular reflexive	<i>himself</i>
Pron~3SingFRefl	pronoun, 3 rd person singular reflexive	<i>herself</i>
Pron~3SingNRefl	pronoun, 3 rd person singular reflexive	<i>itself</i>
Pron~3PIRefl	pronoun, 3 rd person plural reflexive	<i>themselves</i>
Pron~rel	pronoun, relative	<i>that, what, which, who, whom</i>
Pron~relPoss	pronoun, relative reflexive	<i>whose</i>
Poss	possessive marker	<i>' , 's</i>
Prep	preposition, general	<i>about, by, of</i>
Prep~comp	preposition, comparative	<i>than</i>
Part	particle	<i>down, out, up</i>
Inf	infinitive marker	<i>to</i>
Exist	existential	<i>there</i>
Con~co	conjunction, coordinating	<i>and, &, but, or</i>
Con~sub	conjunction, subordinating	<i>also, although, because, since</i>
Qword	question word	<i>what, whatever, who</i>
Resp~yes	response, yes	<i>yes, yep, yeah, aye</i>
Resp~no	response, no	<i>no, nope</i>
DM	discourse marker	<i>ok, okay, well</i>

Label	Explanation	Example/Word(s)
Interj	interjection	<i>ah, aw, oh</i>
Punc~stop	punctuation mark, full stop	.
Punc~exclam	punctuation mark, exclamation mark	!
Punc~query	punctuation mark, question mark	?
Punc~comma	punctuation mark, comma	,
Punc~semi	punctuation mark, semi-colon	;
Punc~colon	punctuation mark, colon	:
Punc~lquot	punctuation mark, left quotation mark	‘, ’, “, ”
Punc~rquot	punctuation mark, right quotation mark	’, ’, ’’, ’’
Punc~lbr	punctuation mark, opening bracket	(
Punc~rbr	punctuation mark, closing bracket)
Sym	symbol	\$, £, €

7 References

- Schmid, H. (1994). Probabilistic Part-of-Speech Tagging Using Decision Trees. *Proceedings of International Conference on New Methods in Language Processing*, Manchester, UK.
- Toutanova, K., Klein, D., Manning, C., and Singer, Y. (2003). Feature-Rich Part-of-Speech Tagging with a Cyclic Dependency Network. In *Proceedings of HLT-NAACL 2003*. pp. 252–259.
- Weisser, M. (2014). The Simple PoS Tagger (Version 1.0) [Computer Software]. Downloadable from http://martinweisser.org/ling_soft.html#tagger.
- Weisser, M. (2018a). The Simple Corpus Tool (Version 2.0) [Computer Software]. Downloadable from http://martinweisser.org/ling_soft.html#viewer.
- Weisser, M. (2018b). Automatically Enhancing Tagging Accuracy and Readability for Common Freeware Taggers. In Y. Tono & H. Isahara (Eds.). *Proceedings of APCLC 2018*, Takamatsu, Japan. 502–505.