

Redpen を使った文書自動添削ツール

PM コース 矢吹研究室 1442031 氏名 小山 隆太郎

1. 研究の背景

ソフトウェアエンジニアはプログラムを組むだけではなく、たくさんの技術文書を書く。そして専門的なチュートリアル、マニュアル等を読み手に理解してもらうことが大切である。ソフトウェアを開発するには多くの実装テストをする必要がある。そのため CheckStyle や lint 等の静的解析ツールを導入することで、フォーマットのエラーを自動で検知することが出来る。しかし静的解析ツールは文書のチェックを行えるものがなく、文中のミスを修正することに時間を割く等で作業の本質に取り組む時間以外に無駄が生じてしまうことがある。このことはエンジニアにとって悩みの種である。このような状況に対して Redpen と呼ばれる文書静的解析ツールが開発され、日々エンジニアが改良を加えている [1]。Redpen は一部の静的解析ツール (CheckStyle, lint 等) に相当する機能を文書に与えるものであり、文書作成でも最低限の検査を自動で行いたいという動機の下、改良が進んでいる。大学の授業等で文書を提出することがあるため、Redpen の開発に着目することでこれからの文書作成の質が向上するのではないかと考えた。

2. 研究の目的

Redpen を利用し、文書作成に割く時間を短縮できるようにすることで、もっと大きな粒度の問題 (実作業、着目すべき文や質等) に集中できるようにすることが目的である。また、文中のミスを少なくした状態で文書を提出できるようにする [2]。

3. 研究の手法

- 1, Redpen が動作する環境を構築する [3]。
- 2, 動作することを確認したら、日本語の文書” ja” が添削できるように設定ファイルを書き換える。
- 3, 参考文献のサンプル文 [4] で添削を行い、添削結果が本の改善文に準拠するよう設定ファイルを書き換える。

4. 研究の結果

Redpen の設定ファイルを書き換え、以下の表現を添削する機能にした。

- SuccessiveWord 同一の単語が連続して使用されていないかを検査する。
- KatakanaEndHyphen カタカナ単語末尾の長音検査する。
- KatakanaSpellCheck カタカナ単語のスペルチェックをする。
- InvalidExpression 顔文字や感情的な表現など技術文書で利用すべきでない句の使用を検査する。
- JapaneseStyle ですます調、である調が混じっていないかを検査する。
- DoubleNegative 二重否定を検知する。
- FrequentSentenceStart 同じ文頭表現が過度に利用されていないか検査する。例えばどの文頭も”私は～”からはじまる等。
- JapaneseAmbiguousNounConjunction 格助詞の「の」+ 名詞連続 + 格助詞の「の」というパターンを発見するとエラーを出力する。
- JapaneseNumberExpression 日本語の数値表現を検査する機能。「ひとつ、ふたつ」、「1 つ、2 つ」などのゆらぎを検知する。
- LongKanjiChain 長過ぎる漢字の連続をみつけるとエラーを出力する。
- SuccessiveSentence 同一内容の文が連続して出現するとエラーを出力する。

- DoubledConjunctiveParticleGa 一文に二回以上, 接続助詞または二重否定の「が」が出現するとエラーを出力する.

Redpen がコマンド上で動作し, 文中ミスを正しく抽出できたので以下に実行結果を載せる.

```

2016-12-08 17:53:00.691 [INFO] cc.redpen.config.ConfigurationLoader - Succeeded to load configuration file
2016-12-08 17:53:00.691 [INFO] cc.redpen.config.ConfigurationLoader - Language is set to "ja"
2016-12-08 17:53:00.692 [WARN] cc.redpen.config.ConfigurationLoader - No variant configuration...
2016-12-08 17:53:00.694 [INFO] cc.redpen.config.ConfigurationLoader - No "symbols" block found in the configuration
2016-12-08 17:53:00.708 [INFO] cc.redpen.config.SymbolTable - "ja" is specified.
2016-12-08 17:53:00.709 [INFO] cc.redpen.config.SymbolTable - "zenkaku" variant is specified
2016-12-08 17:53:01.477 [INFO] cc.redpen.parser.SentenceExtractor - "[ , ? , !]" are added as a end of sentence characters
2016-12-08 17:53:01.478 [INFO] cc.redpen.parser.SentenceExtractor - "[ , , ? , !]" are added as a right quotation characters
2016-12-08 17:53:01.681 [INFO] org.reflections.Reflections - Reflections took 87 ms to scan 1 urls, producing 4 keys and 50 values
2016-12-08 17:53:01.798 [INFO] org.reflections.Reflections - Reflections took 10 ms to scan 1 urls, producing 159 keys and 162 values
2016-12-08 17:53:01.834 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load katakana word dictionary.
2016-12-08 17:53:01.836 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load InvalidExpressionValidator default dictionary.
2016-12-08 17:53:01.845 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load double negative expression rules.
2016-12-08 17:53:01.855 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load double negative words.
draft.txt:1: Validation Error[CommaNumber], The number of commas (5) exceeds the maximum of 3. at line: 野菜ジュースは濃度が上がるほど喉目繊維が豊富になりますが、消化機能低下した人が誤嚥をする危険が増えますが、そこで濃度の異なる野菜ジュースを作り、食物繊維の摂取量、安全、嗜好を考慮するとの濃度が適切なものを研究している。
draft.txt:1: Validation Error[DoubledConjunctiveParticleGa], Found multiple conjunctive particle: "が" at line: 野菜ジュースは濃度が上がるほど喉目繊維が豊富になりますが、消化機能が低下した人が誤嚥をする危険が増えますが、そこで濃度の異なる野菜ジュースを作り、食物繊維の摂取量、安全、嗜好を考慮するとの濃度が適切なものを研究している。
[2016-12-08 17:53:01.922][ERROR] cc.redpen.Main - The number of errors "2" is larger than specified (limit is "1").

```

図 1

図 1 は 2 ケ所の文中ミスを指摘している. Validation Error[CommaNumber], The number of commas (5) exceeds the maximum of 3. at line は文中の句点 (コマ数) が多いことを指摘している. この場合” 摂取量、安全、嗜好” の句点を”・” にすることで見やすい表記になる. Validation Error[DoubledConjunctiveParticleGa], Found multiple conjunctive particle: "が" at line は文中に二重否定があることを指摘している. この場合” 野菜ジュースは濃度が上がるほど喉目繊維が豊富になりますが、消化機能が低下した人が誤嚥をする危険が増えますが、” と二重否定が使われていることがわかる.

```

2016-12-08 18:22:36.127 [INFO] cc.redpen.config.ConfigurationLoader - Succeeded to load configuration file
2016-12-08 18:22:36.128 [INFO] cc.redpen.config.ConfigurationLoader - Language is set to "ja"
2016-12-08 18:22:36.128 [WARN] cc.redpen.config.ConfigurationLoader - No variant configuration...
2016-12-08 18:22:36.130 [INFO] cc.redpen.config.ConfigurationLoader - No "symbols" block found in the configuration
2016-12-08 18:22:36.145 [INFO] cc.redpen.config.SymbolTable - "ja" is specified.
2016-12-08 18:22:36.146 [INFO] cc.redpen.config.SymbolTable - "zenkaku" variant is specified
2016-12-08 18:22:36.399 [INFO] cc.redpen.parser.SentenceExtractor - "[ , ? , !]" are added as a end of sentence characters
2016-12-08 18:22:36.399 [INFO] cc.redpen.parser.SentenceExtractor - "[ , , ? , !]" are added as a right quotation characters
2016-12-08 18:22:37.099 [INFO] org.reflections.Reflections - Reflections took 88 ms to scan 1 urls, producing 4 keys and 50 values
2016-12-08 18:22:37.213 [INFO] org.reflections.Reflections - Reflections took 8 ms to scan 1 urls, producing 159 keys and 162 values
2016-12-08 18:22:37.261 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load katakana word dictionary.
2016-12-08 18:22:37.269 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load InvalidExpressionValidator default dictionary.
2016-12-08 18:22:37.280 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load double negative expression rules.
2016-12-08 18:22:37.289 [INFO] cc.redpen.util.DictionaryLoader - Succeeded to load double negative words.
draft.txt:1: Validation Error[SuccessiveWord], Found word "た" repeated twice in succession. at line: 気が付くとレベルが20になっていた。
draft.txt:1: Validation Error[SuccessiveSentence], Found similar two sentences in succession: "ポケモンGoはおもしろい。" and "ポケモンGoはおもしろい。" at line: ポケモンGoはおもしろい。
[2016-12-08 18:22:37.353][ERROR] cc.redpen.Main - The number of errors "2" is larger than specified (limit is "1").

```

図 2

図 2 は 2 ケ所の文中ミスを指摘している. Validation Error[SuccessiveWord], Found word "た" repeated twice in succession. at line: は文末の” た” が 2 つあることを指摘している. Validation Error[SuccessiveSentence], Found similar two sentences in succession は文中に同じパラグラフが 2 つあることを指摘している. この場合” ポケモン Go はおもしろい。” が 2 つあることを指摘している.

5. 考察と今後の計画

Redpen には設定ファイルが設けられており, 設定を書き換えるごとに出力結果も変わる. その為きちんとした設定が確立されておらず, 現在も研究が続いている. Redpen のファイルの設定を確立し, 文中のミスを抽出することは出来たが, 中には不要である間違いを指摘されることがあった. これを取り除く設定を確立できれば, より精度の高い文書添削ツールになり得るだろう. 今後は皆が共通して使えるように添削マシンを拡張する予定である.

参考文献

- [1] Redpen を使って技術文書を気軽に校正しよう. <http://gihyo.jp/lifestyle/serial/01/redpen>.
- [2] Redpen でお手軽文書校正. <http://tech.respect-pal.jp/try-redpen/>.
- [3] 阿部紘久. シンプルに書く 伝わる文章術. 飛鳥新社, 2012.
- [4] Redpen1.7 ドキュメント. http://redpen.cc/docs/latest/index_ja.html.