

PubAnnotation - テキスト 注釈リポジトリ

PubAnnotation - Text Annotation Repository

金 進東 *1

Jin-Dong Kim

*1情報・システム研究機構/ライフサイエンス統合データベースセンター

Rsearch Organization of Information and Systems / Database Center for Life Science

Text annotation is expensive but indispensable for development of text mining technology. While there are a few projects developing text annotation data, their reusability is becoming more and more a serious issue. PubAnnotation is developed to address the reusability issue. Developed primarily as an annotation storage system that can easily scale to tera-bytes, PubAnnotation features, among others, text alignment, which enables comparative and integrative use of text annotation. Together with its satellite tools, including PubDictionaries and TextAE, it also serves as a platform of crowd-sourcing annotation.

1. はじめに

テキストマイニング (text mining) に関する研究開発においては、人がテキストを読むさいに取得すべき情報を、機械が抽出しやすいように構造化された形でテキストに追加する作業が行われており、テキスト注釈 (text annotation) と呼ばれる。テキスト注釈データ (注釈付きコーパス) は手動あるいは自動で作られるが、自動でテキスト注釈 (以下、自動注釈) を高精度で行うことができれば、テキストマイニングの性能向上が期待される。手動で作られた注釈データは一般的に品質が高く、ベンチマーク (benchmark) 用のデータとして自動注釈の性能を測るため使われたり、自動注釈のモデルとして活用されたりと重要な資源である。しかし、高品質のテキスト注釈データは生産コストが高いことから、既存データを再利用することが期待される。

世界的にテキスト注釈データを構築するプロジェクトは多数あり、成果は貴重な資源として広く活用されている。しかし、注釈データは各プロジェクト独自のやり方で作られ、形式 (format) やテキスト前処理 (pre-processing) 方法等がばらばらになってしまい互換性が一般的に低く、再利用を困難にする大きな原因になっている。これらのテキストデータを統一し、標準的な手順でアクセスできるようにすることでデータの再利用性が高まると期待される。

このような状況を踏まえ、テキスト注釈の統合的な管理のためのストーリージシステム (storage system) として PubAnnotation を開発した。本論文では注釈レポジトリ (repository) とクラウドソーシングプラットフォーム (crowd-sourcing platform) としての PubAnnotation の機能に関して紹介する。

2. PubAnnotation

2.1 注釈レポジトリ

PubAnnotation は標準化されたテキスト注釈データのレポジトリとして機能できるよう開発している。そのため複数の研究グループもしくは個人によって作られた様々な注釈データに対応する必要がある。そのため、以下の機能を実装した。

- **テキストアライメント**: 異なるグループやユーザが作成した注釈データを集めると、その中には同一のテキストに付けられた注釈もある。しかし、多くの注釈プロジェクトは注釈作業を容易にする目的でテキストを前処理するためにテキストが変わる場合が多い。例えば、ギリシャ文字の展開 (「 α 」を「alpha」に変換するなど) やトークン化が行われる。その結果、同一テキストに対して作られた注釈データであっても、異なるプロジェクトで作られたもの同士では互換性のない場合が多い。この問題を解決するため、PubAnnotation は Generalized LCS algorithm [Kim 2013] を用いたテキストアライメント (text alignment) 機能を実装している。

- **プロジェクト管理**: 複数のグループ、ユーザが注釈データを登録する場を提供するためにプロジェクト管理機能は必須である。PubAnnotation では誰もがアカウントを作り、注釈プロジェクトを始めることができる。自身のプロジェクトに自身の注釈データを格納することができるだけでなく、副管理者を指定して共同作業することもできる。

- **関係データベースの利用**: PubAnnotation は安定性とスケーラビリティの面で技術的に熟成されている関係データベースを基盤として用いている。そのため大量の注釈データを効率的に管理することが可能である。

- **検索機能**: 大量の注釈データを格納するため、必要な時に必要な部分にアクセスできる仕組みが必要になる。PubAnnotation は関係データベースに格納しているため、SQL を使った検索が行える。

- **REST API**: 様々な条件でデータにアクセス可能な REST API を提供している。

- **形式変換**: PubAnnotation は注釈データを格納するために独自のテーブル形式を用いているが、様々な要求に対応するため形式変換機能を持つ。現時点での出力可能な形式は、プログラミングのための JSON、交換形式として広く使われる XML、セマンティックウェブのための RDF である。形式変換はより多くの要求に対応するため Plug-in システムとして実装されている。

連絡先: 金 進東, ライフサイエンス統合データベースセンター, 千葉県柏市若柴 178-4-4, Tel:04-7135-5508, Fax:04-7135-5534, e-mail: jdkim@dbcls.rois.ac.jp

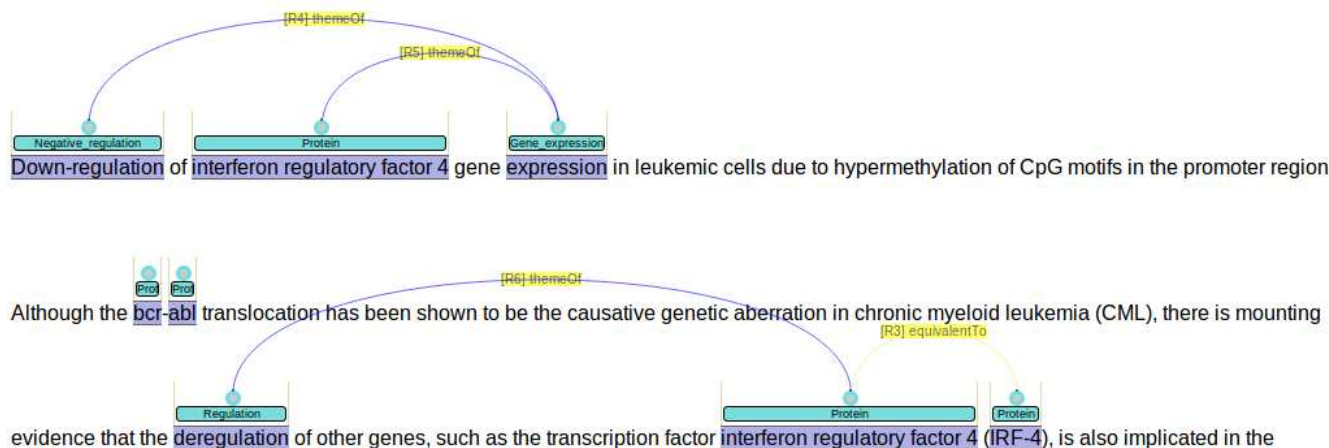


図 1: PubAnnotation が格納している注釈の例。TextAE の編集イメージ。

2.2 クラウドソーシング

PubAnnotation はクラウドソーシングによるテキスト注釈にも対応するため以下の機能を実装している。これらの機能を使って作られた注釈データはすぐに誰とでも共有できる。

○ 手動注釈作業のためのエディター: 人が論文を読みながら簡単に注釈できるよう PubAnnotation の連携ツールとして注釈エディターである TextAE を提供している。TextAE は Javascript で実装されたウェブベースのツールで、グラフィックユーザインタフェース (GUI) を使って注釈の編集ができる。

○ 辞書による自動注釈: 人手で注釈を付ける作業であっても、最初に自動的に付け、その結果の誤りを人手により修正することが最近の一般的な手順になっている。特に注釈のための用語辞書を用意し、それに基づいて自動的に注釈を付けたいという要求は多い。このため、PubAnnotation は連携システムとして辞書の管理と辞書による自動注釈ができる PubDictionaries というシステムを提供している。

○ オープンアーキテクチャ: 自動注釈と注釈エディターは人手による注釈には必須なツールであり、PubAnnotation は連携ツールとして PubDictionaries と TextAE を提供しているが、同様のツールを開発するグループは他にもある。このため、外部のツールも必要に応じて PubAnnotation と一緒に使えるよう REST API を公開している。PubAnnotation と TextAE も当該 API を通じて PubAnnotation と繋がる仕組みになっている。

3. 終わりに

テキスト注釈データは構築にコストが掛かるものであるが、信頼できるテキストマイニングシステムの開発のためには必須の資源である。従って、構築されたデータを共有することでテキストマイニング研究開発コミュニティにおける開発コストが軽減できると思われる。このため、多様な注釈を統合的に管理でき、かつ、スケーラブルなレポジトリとして PubAnnotation を開発した。連携ツールとして辞書ベースの自動注釈システムである PubDictionaries やウェブベースの注釈エディターである TextAE も開発し、誰もが気軽に注釈データを作り、共有できるクラウド注釈にも対応している。図 1 は PubAnnotation

に格納されている注釈の例である。PubAnnotation とその連携システムは以下の URL からアクセス可能である。

- PubAnnotation: <http://www.pubannotation.org>
- PubDictionaries: <http://www.pubdictionaries.org>
- TextAE: <http://textae.dbcls.jp>

参考文献

- [Kim 2013] Kim, Jin-Dong: A Generalized LCS Algorithm and Its Application to Corpus Alignment (2013), Proceedings of the Sixth International Joint Conference on Natural Language Processing, 1112–1116.