

Parsing Tweets into Universal Dependencies

First Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Second Author

Affiliation / Address line 1
Affiliation / Address line 2
Affiliation / Address line 3
email@domain

Abstract

1 Introduction

Analyzing the syntax of tweets is challenging for traditional NLP tools because most of the tweet texts are informal and noisy.

In this paper, we propose to parse the tweets in the convention of universal dependencies and built up the whole pipeline to parse the tweets from the raw text form.

Contribution of this paper includes:

- We create a new version of tweet Treebank Tweepbank 2.0
- We propose a neural network method to parse tweets into universal dependencies
- We study the adaptation of universal dependencies for analyzing tweets

2 Related Work

Eisenstein (2013) reviewed NLP approaches for analyzing text on social media, especially for tweets and showed that there are two major directions for NLP community to handle the tweets, including normalization and domain adaptation. He also pointed out that normalization can be problematic because precisely defining the normalization task is difficult.

Kong et al. (2014) argues that the Penn Treebank approach to annotation is poorly suited to more informal genres of text, as some of the annotation challenges for tweets, including token selection,

multiword expressions, multiple roots, and structure within noun phrases diverge significantly from conventional approaches. They believe that rapid, small scale annotation efforts performed by imperfectly-trained annotators should provide enough evidence to train an effective parser, given the rapidly changing nature of tweets (Eisenstein, 2013), the attested difficulties of domain adaptation for parsing (Dredze et al., 2007), and the expense of creating Penn Treebank-style annotations (Marcus et al., 1993). Therefore, they build a new corpus of tweets (Tweepbank), with conventions informed by the domain, using new syntactic annotations that can tackle all the forementioned problems annotated in a day by two dozen annotators, most of whom had only cursory training in the annotation scheme. Then, they modify the decoder of the TurboParser, a graph-based dependency parser, which is open-source and has been found to perform well on a range of parsing problems in different languages (Martins et al., 2013) to adapt to the Tweepbank dataset, and incorporate new features such as Brown Clusters and Penn Treebank features and changes to specification in the output space into TurboParser.

3 Data

3.1 Linguistic Phenomena of Twitter

As the representative of the web language, different from the standard languages such as newswire, twitter has its own linguistic phenomena, that could be categorized into word level and structure level phenomena, and we will discuss them separately.

3.1.1 Token Level

There are special tokens that will appear only or much more frequently in tweets, including

- Retweet discourse marker: RT
- URL: <http://bit.ly/xyz>
- Hashtag: #ACL
- At-mentioned username: @user
- Emoticon, emoji and symbol: :), - - - <<<, :-),
- Acronym: wtf (what the fuck), smh (shake my hand), mfw (my face when), ima (i am going to), rn (right now), af (as fuck)
- Contraction: he's, buy'em, gonna, trma, gimme, im
- Truncated token: because each tweet can have only 140 characters, excessive characters will be cut off, and usually there could be token in the last of the tweet that is partially elided

Consider the following example

RT @Yijia : #ACL2017 im heading to Canada for ACL can u gimme some money :) <http://url1> I would love to do it with th

From the example, we argue that most of the the retweet discourse markers, URLs, hashtags, usernames and the emoticons/emojis/symbols usually do not have clear syntactic functions in the tweet, and therefore we should not include them in the analysis, and treat them all equally as non-syntactic tokens. Especially, for truncated tokens, although it might be possible to recover the original tokens from the context in some cases, and we can infer that the truncated token “th” could be probably the token “the” it is extremely hard to further predict the rest of the elided sentence. For the simplicity of the annotation, we will consider all the truncated tokens as non-syntactic tokens.

For acronyms and contractions, it is obviously better if we could recover their original forms before analyzing the whole tweet, and we believe that all of these tokens will always be syntactically part of the tweet.

However, consider another example

@Yi u gonna ♡ it . #NAACL will be #Awesome .
Check my new paper with @nlpnoah on <http://url2>
RT it !!

Similar in the analysis of Kong et al. (2014), in some cases, retweet marker (RT), URLs (<http://url2>), hashtags (#NAACL, #Awesome), usernames (@Yi, @nlpnoah) and emoticons (♡) can also have syntactic functions, and we need to take them into the account in the syntactic analysis.

We define acronyms and contractions as full syntactic tokens, truncated tokens as zero syntactic tokens, and all of the other tokens as partial syntactic tokens.

3.1.2 Structure Level

Besides the token level phenomena of twitter, there are also structural patterns appearing in the twitter.

- Retweet structure: RT @user : < tweet content > is a typical structural pattern when a user is retweeting from other users
- Parataxis: Very often, one tweet is comprised of many sentences or phrases without any delimiting punctuations, such as < sentence1 > < phrase1 > < sentence2 > ...

We treat “RT @user” as the structural patterns in tweet and keep their annotations consistent across tweets. For parataxis case

3.2 POS Tagging

Gimpel et al. (2011) proposed a set of part-of-speech tags that handle most of token level phenomena of tweet. We are inspired by their work but argue that, first different abbreviation can have different syntactic functions. Like mfw is usually followed by an adverbial clause and ima is usually followed by a clausal complement. It is not reasonable to treat them in the same part-of-speech. In this paper, when annotating the POS tagging for abbreviations, we first try to recover their original forms, then use the POS of the core-word as the POS for the abbreviation. Second, four special POS tags (S, L, M, Y) were designed to handle contraction words in Gimpel et al. (2011). Major concern of designing such tags is to minimize the effort of tokenization. However, contractions of common nouns and

pronouns are casted into the same category which increase the difficulty of distinguishing their syntactic function (say, there's and book'll are treated with the same syntactic function). What's more, only a small proportion of words can be categorized into these tags (2.7 % in total), which cast a doubt of the usefulness of these certain tags. In this paper, we believe such contraction can be properly handled by tokenization module, so we suggest to tokenize the contraction word and annotate POS tag accordingly. Besides the contraction that be conventionally tokenized, tweets also witness a set of unconventional contraction like iv (I've), whatis (what is). In this paper, we follow the same idea of annotation abbreviation to handle the unconventional contractions and use the POS of core word of the original form as their POS. Third, special POS was designed to handle emoticon in Gimpel et al. (2011). However, in most cases, emoticon plays the same role as most of the symbolic tokens. In this paper, we follow the UD guideline to annotate the emoticon as symbol (SYM). At last, its arguable that some of the hashtags, URLs can work as a nominal in tweets. Whether treating them as the same part-of-speech or different ones according to their context is an open question. A preliminary survey on the standard UD English data shows that URL, email address are all tagged as the foreign language (X), so we also tag them as X and leave the disambiguation of their syntactic function to the annotation of parse tree.

We use the Universal POS tags (Petrov et al., 2012) to tag the tweet tokens.

3.3 Dependency Annotation

4 Pipeline

4.1 Tokenization

We use the UDPipe¹ to tokenize the tweets and then detokenize the wrongly tokenized usernames and hashtags.

4.2 POS Tagging

4.2.1 Tweets-level Special Construction

There are several tweet-level constructions which are unconventional to standard text, including:

- Retweet: RT @user : \langle sentence \rangle

¹<https://github.com/ufal/udpipe>

- Leading or ending topic marked as hashtag: #topic #topic #topic \langle sentence \rangle #topic #topic #topic
- Leading or ending complementary URL: \langle complementary URL \rangle \langle sentence \rangle \langle complementary URL \rangle

4.3 Sentence Segmentation

5 Model

6 Experiments

7 Conclusion

References

- Mark Dredze, John Blitzer, Partha Pratim Talukdar, Kuzman Ganchev, João Graca, and Fernando Pereira. 2007. Frustratingly hard domain adaptation for dependency parsing. In *Proceedings of the CoNLL Shared Task Session of EMNLP-CoNLL 2007*, pages 1051–1055, Prague, Czech Republic, June. Association for Computational Linguistics.
- Jacob Eisenstein. 2013. What to do about bad language on the internet. In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 359–369, Atlanta, Georgia, June. Association for Computational Linguistics.
- Kevin Gimpel, Nathan Schneider, Brendan O'Connor, Dipanjan Das, Daniel Mills, Jacob Eisenstein, Michael Heilman, Dani Yogatama, Jeffrey Flanigan, and Noah A. Smith. 2011. Part-of-speech tagging for twitter: Annotation, features, and experiments. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: Short Papers - Volume 2*, HLT '11, pages 42–47, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Lingpeng Kong, Nathan Schneider, Swabha Swayamdipta, Archana Bhatia, Chris Dyer, and Noah A. Smith. 2014. A dependency parser for tweets. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1001–1012, Doha, Qatar, October. Association for Computational Linguistics.
- Mitchell P. Marcus, Mary Ann Marcinkiewicz, and Beatrice Santorini. 1993. Building a large annotated corpus of english: The penn treebank. *Comput. Linguist.*, 19(2):313–330, June.
- Andre Martins, Miguel Almeida, and Noah A. Smith. 2013. Turning on the turbo: Fast third-order non-projective turbo parsers. In *Proceedings of the 51st*

Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers), pages 617–622, Sofia, Bulgaria, August. Association for Computational Linguistics.

Slav Petrov, Dipanjan Das, and Ryan McDonald. 2012. A universal part-of-speech tagset. In Nicoletta Calzolari (Conference Chair), Khalid Choukri, Thierry Declerck, Mehmet Uur Doan, Bente Maegaard, Joseph Mariani, Asuncion Moreno, Jan Odijk, and Stelios Piperidis, editors, *Proceedings of the Eight International Conference on Language Resources and Evaluation (LREC'12)*, Istanbul, Turkey, may. European Language Resources Association (ELRA).