

# Biomarker-NLP: A Python Package for Mining Biomarker Information for FDA-Approved Targeted Cancer Therapies

Junxia Lin<sup>1</sup>, Yuezheng He<sup>1</sup>, Subha Madhavan<sup>1, 2</sup>, Chul Kim<sup>3</sup>, and Simina M. Boca<sup>1, 4</sup>

<sup>1</sup> Innovation Center for Biomedical Informatics, Georgetown University Medical Center, Washington, DC, USA <sup>2</sup> Competitive Intelligence, Oncology R&D, AstraZeneca, Gaithersburg, MD, USA <sup>3</sup> Division of Hematology and Oncology, Georgetown Lombardi Comprehensive Cancer Center, Georgetown University, Washington, District of Columbia, USA <sup>4</sup> Early Biometrics & Statistical Innovation, Data Science & Artificial Intelligence, R&D, AstraZeneca, Gaithersburg, MD, USA

DOI: [10.21105/joss.0XXXX](https://doi.org/10.21105/joss.0XXXX)

## Software

- [Review](#)
- [Repository](#)
- [Archive](#)

Editor: [Editor Name](#)

Submitted: 01 January XXXX

Published: 01 January XXXX

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](https://creativecommons.org/licenses/by/4.0/)).

## Summary

This Biomarker-NLP (`biomarker_nlp`) package aims to provide natural language processing (NLP) functionalities for mining and processing biomarker information for targeted cancer therapies approved by the US Food and Drug Administration (FDA). Treatment biomarkers are specific molecular changes, including mutations and gene or protein expression measurements, that are used to decide whether a certain therapy should be prescribed to an individual. Thus, they are often included in the FDA-approved therapy labels, especially for targeted cancer therapies. Our tool pulls information from two webpages in HTML format: 1) The National Cancer Institute (NCI)'s list of FDA-approved targeted cancer therapies [`@Therapy`] and 2) The National Library of Medicine (NLM)'s DailyMed database of drug labels [`@DailyMed`]. Biomarker-NLP parses the NCI and DailyMed HTML pages using tools in the `lxml` library developed by `@Behnel:2005`. It allows users to quickly and easily scrape certain pieces of information from NCI and DailyMed without requiring them to consider the HTML tree structure. The free text biomarker information is mined and structured into fixed entities, including therapy name, disease name (cancer type), gene or protein in biomarker, name of therapies prescribed in combination, etc. For recognizing the biomarker entities, such as genes and proteins, we utilize the pre-trained named-entity recognition (NER) models from `ScispaCy` [`@Neumann:2019`]. In addition, as negated biomarkers can be important but challenging to extract, our package provides tools to detect negations in sentences through two pre-trained negation models from `@Khandelwal:2020` NegBERT program, which applies a transfer learning approach. One model is the negation cue detection model that detects the negation cues in a sentence, while the second is the negation scope detection model that recognizes the scope of negation in a sentence. As the NegBERT program does not provide the output models, we performed the training step and published these two models for free use, integrating them into our package so that users can easily mine the negated biomarker information by using the relevant functions. As an example,

```
>>> from biomarker_nlp import negation_cue_scope
>>> from biomarker_nlp.negation_negbert import *
>>> modelCue = torch.load('/path/to/negation/cue/detection/model') # path
to the model file
>>> modelScope = torch.load('/path/to/negation/scope/detection/model') #
path to the model file
>>> txt = "KEYTRUDA is not recommended for treatment of patients with PMB"
```

```
44 CL who require urgent cytoreductive therapy."
45 # detect negation cue
46 >>> negation_cue_scope.negation_detect(text = txt, modelCue = modelCue)
47 True
48 # extract the negation scope
49 >>> negation_cue_scope.negation_scope(text = txt, modelCue = modelCue, mo
50 delScope = modelScope)
51 ['KEYTRUDA is', 'recommended for treatment of patients with PMBCL who']
```

52 As we can see from the output above, the `negation_detect()` function detects if a negation  
53 cue is presented in a sentence. Afterwards, the `negation_scope()` function extracts the  
54 scope of the associated negation cue from the sentence. Then, we can use other functions  
55 from the package to detect the biomarkers presented in the resulting scope phrases to get the  
56 negated biomarkers.

## 57 Statement of need

58 The NCI website represents a convenient starting location for exploring targeted cancer thera-  
59 pies for patients and physicians, as well as bioinformaticians and biomedical researchers. Dai-  
60 lyMed is a database that provides official label information for about 140,000 FDA-approved  
61 and FDA-regulated products submitted to the FDA [DailyMed]. For a drug or a biological  
62 product, its label contains prescribing information in a structured textual format. Each label  
63 includes various sections, such as the indication and usage and dosage and administration  
64 [DailyMed]. Within each section, information is mostly in free-text format. The informa-  
65 tion contains a variety of biomedical data, including biomarker information. This biomarker  
66 information is valuable for and currently being used by a wide range of stakeholders, such as  
67 doctors, bioinformaticians, healthcare providers, and biomedical researchers. However, as this  
68 labeling information is in free-text format and may be updated with new indications, searching  
69 through multiple labels and reading every word in order to perform curation activities is often  
70 time-consuming, low efficiency labor for bioinformaticians. Here, we present Biomarker-NLP,  
71 a Python package to process biomedical text and extract biomarker information from NCI and  
72 DailyMed efficiently. Curators will still be required to check the NLP output, but are expected  
73 to spend substantially less time on these activities. Thus, Biomarker-NLP will be integrated  
74 into an “AI-augmented curation” workflow, building on similar work by @Mahmood:2017, which  
75 developed a system to extract associations between genomic anomalies and drug responses  
76 from the biomedical literature, but focused on cancer therapy labels.

## 77 Acknowledgements

78 This work was completed as part of a project funded by a pilot award (P30CA051008, PI of  
79 pilot: Kim).

## 80 Declarations

81 Subha Madhavan and Simina M. Boca are currently employees and minor shareholders of  
82 AstraZeneca.

## 83 References