

Spá: a web-based viewer for text mining in Evidence Based Medicine

Kuiper, J¹., Marshall, I.J.², Wallace, B.C.³, and Swertz, M.A.¹

¹ University of Groningen P.O. Box 30001, 9700 RB Groningen
{joel.kuiper,m.a.swertz}@rug.nl

² King's College London, London SE1 3QD, UK
iain.marshall@kcl.ac.uk

³ Brown University, Providence, RI 02906, USA
byron_wallace@brown.edu

Abstract. Summarizing the evidence about medical interventions is an immense undertaking. In part this is because unstructured Portable Document Format (PDF) documents remain the main vehicle for disseminating the results of clinical trials. Clinicians and researchers typically manually extract and synthesise information from these PDFs which is then published as a *systematic review*. We introduce Spá⁴⁵, a web-based PDF viewer that allows for annotation and summarisation Machine Learning modules to semi-automate the task of PDF annotation and summarisation. To illustrate, we demonstrate the use of Spá to semi-automate the assessment of bias in clinical trials. Spá can visualize the output from hybrid models that simultaneously classify documents (e.g., identifying trials as low risk of various biases), and also annotates sentences or words that contribute most to the classification. The Spá architecture is modular, and new models may be trivially added, thus the tool made be widely useful in other domains with a PDF-based literature, including law, physics, and biology.

1 Introduction

Imposing structure on full-text documents (e.g., identifying specific sentences of interest) is an important and practical task in natural language processing and machine learning. For example, *systematic reviews* are fundamental tools in Evidence-based Medicine (EBM) [2], and aim to answer clinical questions by exhaustively searching the literature for relevant studies. During systematic reviewing data on trial design and conduct, intervention effectiveness, adverse effects, and potential biases, is extracted and synthesized to produce a textual and statistical summary. As an example, we demonstrate the task of bias assessment

⁴ From the Old Norse word spá or spæ referring to prophesying and which is cognate with the present English word “spy,” continuing Proto-Germanic *spah- and the Proto-Indo-European root *(s)pek̑ (to see, to observe)

⁵ Source code available under GPLv3 at <https://github.com/joelkuiper/spa> [1]; demo available at <http://spa.clinici.co/>

in which a reviewer assess, for example, whether the participants and personnel were properly blinded. Assessing this is a time-consuming task, a single trial can take a domain expert around ten minutes [3]. Making matters worse, a single review may typically include several dozen trials, and due to empirical low rates of reviewer agreement it is regarded as best practice to have each study assessed twice by independent reviewers who later come to a consensus[4].

Machine Learning (ML) methods provide the machinery to automate such extractions; as they can effectively impose the desired structure onto PDFs. But if such technologies are to be practically useful, we need tools that visualize these model predictions and annotations. Here we describe Spá, which aspires to realize this aim.

Spá is an open-source, web-based tool that can incorporate state-of-the-art machine learning predictors to automatically annotate PDF articles. As a practical demonstration of this technology, we have built a machine learning system that automatically annotates PDFs to facilitate EBM. Specifically, this system leverages a multi-task model that simultaneously assesses the risk of biases across several domains. This tool is unique in that it leverages state-of-the-art ML models applied to full-text articles to assist practitioners of EBM.

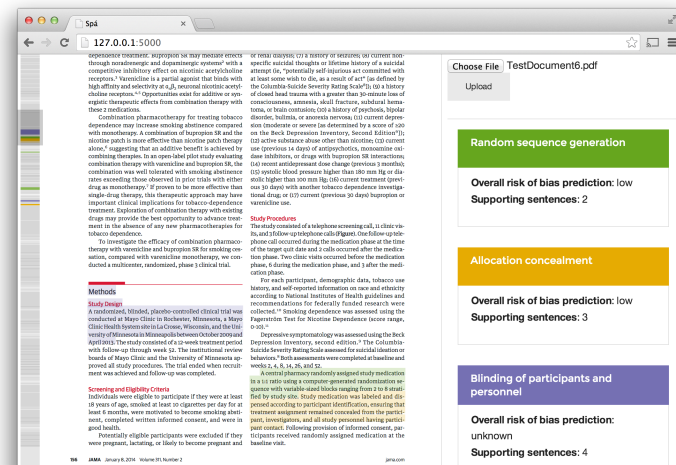


Fig. 1. Screenshot of a PDF with highlighted risk of bias. Here the risk of bias is assessed to be low, for example, and one of the supporting sentences for this assessment describes the randomization procedure (highlighted in green).

While our application of interest is EBM, we emphasize that the visualization tool can be used for any domain in which one wants to annotate PDFs. Thus the contribution of this work is two-fold, as we present: (1) a practical tool that incorporates machine learning to help researchers rapidly assess the risk of biases in published biomedical articles, and, (2) a general open-source web tool

for visualizing the predictions of trained models from full-text articles. These contributions are described further in sections 2 and 3, respectively.

2 Automating Evidence-based Medicine

2.1 Machine Learning Approaches

We briefly describe our model for assessing the study risk of bias (and supporting sentences) across the following domains: random sequence generation, allocation concealment, blinding of participants and personnel, blinding of outcome assessment, incomplete outcome data and selective reporting. To train our model, we have leveraged the Cochrane Database of Systematic Reviews (CDSR) in lieu of manually annotated data, which would be expensive to collect. The CDSR contains descriptions and data about clinical trials reported in existing systematic reviews. Briefly, we match (the full-texts of) studies to entries in the CDSR, which contains risk of bias assessments; this provides document level labels. The CDSR also contains sentences that annotators indicated as supporting their assessments. We match these strings to substrings in the PDFs to provide sentence-level supervision. This can thus be viewed as a *distantly supervised* [5,6] approach.

Concerning the risk of bias predictions, we have two tasks for a given article. First, predict the overall risk of bias for each of the domains of interest, and second extract the sentences that support these assessments. In both cases we leverage standard bag-of-words text encoding and linear-kernel Support Vector Machines (SVMs). Because the risk of bias predictions are correlated (across domains), we take a *multi-task* [7] approach to classification and jointly learn a model for the domains. We accomplish this by way of a feature space construction that includes both shared and domain-specific terms (similar to the domain adaptation approach in [8]).

With respect to the second task of sentence extraction, we *jointly* assess the risk of bias associated with a given article *and* the sentences that support this judgment. Specifically, we first make sentence level predictions (using one set of trained models), and then insert features representing the tokens (words) in the predicted sentences for exploitation by the document level classifier. Due to space constraints here, we omit further specifics of our approach, but we provide further details elsewhere [9]. Figure 1 shows example system output: the overall risk of bias with respect to random sequence generation is judged to be low, and supporting sentence

3 Spá Architecture Overview

Spá relies on Mozilla pdf.js⁶ for visualization of the document and text extraction. The results of the text extraction are processed server-side by a variety of processing topologies, as outlined in figure 2. Results, which could come from

⁶ <http://mozilla.github.io/pdf.js>

complicated machine learning systems, are communicated back to the browser and displayed using React components⁷.

For each of the annotations the relevant nodes in the document are highlighted and a custom scrollbar, inspired by substance.io⁸, that acts as a mini-map is projected to show where it resides within the document. The user can interactively activate and inspect specific results.

4 Future work

We present a web-based tool for interactive visualization of annotations and meta-data on PDF documents. This allows users to see the results from machine learning predictions for a specific document. We have demonstrated the utility of this system within the context of EBM by automatically extracting potential risks of bias (and supporting sentences) via state-of-the-art ML methods.

More generally, we believe the tool to be potentially useful for a much wider range of text mining and machine learning applications. To increase the generality of the tool we are developing a pluggable system for processing topologies, allowing developers to quickly plug in new systems for automated PDF annotation. Furthermore, work is being done on allowing users to save selected annotations, possibly embedded within the document itself, for sharing and off-line use. The ultimate vision is to have an extensible system for semi-automated (machine assisted) screening, data extraction and data summarization for EBM, and to allow rapid development of similar systems for other domains.

Acknowledgments Part of this research was funded by the European Union Seventh Framework Programme (FP7/2007-2013) under grant agreement n° 261433.

References

1. Kuiper, J., Wallace, B.C., Marshall, I.J.: Spa. <http://figshare.com/articles/Spa/997707> (2014)

⁷ <http://facebook.github.io/react>

⁸ <http://substance.io/beta/>

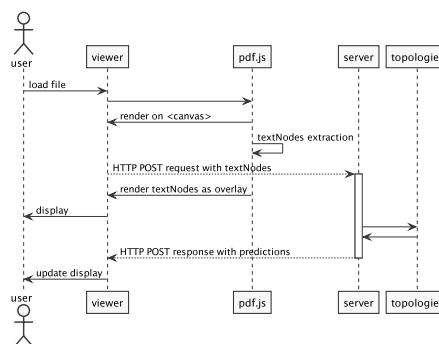


Fig. 2. Sequence diagram of a typical request-response in Spā.

2. Sackett, D.L., Rosenberg, W.M., Gray, J., Haynes, R.B., Richardson, W.S.: Evidence based medicine: what it is and what it isn't. *BMJ: British Medical Journal* **312**(7023) (1996) 71–72
3. Hartling, L., Bond, K., Vandermeer, B., Seida, J., Dryden, D.M., Rowe, B.H.: Applying the risk of bias tool in a systematic review of combination long-acting beta-agonists and inhaled corticosteroids for persistent asthma. *PloS one* **6**(2) (January 2011) e17242
4. Hartling, L., Ospina, M., Liang, Y.: Risk of bias versus quality assessment of randomised controlled trials: cross sectional study. *BMJ* **339** (2009) b4012
5. Mintz, M., Bills, S., Snow, R., Jurafsky, D.: Distant supervision for relation extraction without labeled data. In: *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2*, Association for Computational Linguistics (2009) 1003–1011
6. Nguyen, T., Moschitti, A.: End-to-end relation extraction using distant supervision from external semantic repositories. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies: short papers-Volume 2*, Association for Computational Linguistics (2011) 277–282
7. Evgeniou, T., Pontil, M.: Regularized multi-task learning. In: *Proceedings of the tenth ACM SIGKDD international conference on Knowledge discovery and data mining*, ACM (2004) 109–117
8. Daumé III, H.: Frustratingly easy domain adaptation. In: *Association for Computational Linguistics (ACL). Volume 1785*. (2007)
9. Marshall, I.J., Kuiper, J., Wallace, B.C.: Automatically assessing bias in clinical trials and extracting supporting sentences. In: *under review at Knowledge Discovery in Databases (KDD)*. (2014)