

# Chapter 15

## Machine Learning Applications in Mega-Text Processing

**Marina Sokolova**

*CHEO Research Institute, Canada*

**Stan Szpakowicz**

*University of Ottawa, Canada and Polish Academy of Sciences, Poland*

### **ABSTRACT**

*This chapter presents applications of machine learning techniques to problems in natural language processing that require work with very large amounts of text. Such problems came into focus after the Internet and other computer-based environments acquired the status of the prime medium for text delivery and exchange. In all cases which the authors discuss, an algorithm has ensured a meaningful result, be it the knowledge of consumer opinions, the protection of personal information or the selection of news reports. The chapter covers elements of opinion mining, news monitoring and privacy protection, and, in parallel, discusses text representation, feature selection, and word category and text classification problems. The applications presented here combine scientific interest and significant economic potential.*

### **INTRODUCTION**

The chapter presents applications of Machine Learning (ML) to problems which involve processing of large amounts of texts. Problems best served by ML came into focus after the Internet and other computer-based environments acquired the status of the prime medium for text delivery and exchange. That is when the ability to work extremely large amounts of texts, which ML applications had not previously faced, became a major issue. The resulting set of techniques and practices, which we name *mega-text language processing*, are meant to deal with a mass of informally written, loosely edited text. A case in point is the analysis of opinions expressed in short informal texts written and put on the Web by the general public (Liu 2006). The sheer volume and variety of suddenly available language data has neces-

DOI: 10.4018/978-1-60566-766-9.ch015

sarily invited the use of computing software capable of handling such a mass of data, learning from it and acquiring new information.

Until now, no clearly delineated subfield of Natural Language Processing (NLP) dealt with mega-texts – textual data on the Web, computer-mediated text repositories and in general texts in electronic format. Text Data Mining – a form of Data Mining – concerns itself with deriving new information from texts, but most often restrains from the study of language. Still, many researchers focus on the study of language, for example lexical, grammar and style issues, in such texts (Crystal 2006; Liu 2006). That no overarching NLP discipline has emerged can be explained by the fact that electronic texts and old-fashioned texts in books or newspapers share major characteristics. We discuss these characteristics in the handbook chapter “Machine Learning in Natural Language Processing”.

This chapter will show that ML techniques measure up well to the challenges that mega-texts pose. We focus on applications in aid of the study of language. In all cases which we discuss, an algorithm has ensured a meaningful result, be it the knowledge of consumer opinions, the protection of personal information or the selection of news reports. Although we mostly focus in this chapter on text classification problems, we go beyond document topic classification. English, the most popular language of the Web, is the default language of much of the scientific discourse. We state when problems deal with languages other than English.

In the chapter we cite standard measures used in NLP (Precision, Recall, F-score). Calculated for classifiers produced by an algorithm, they build on the numbers of correctly classified positive examples  $TP$ , incorrectly classified positive examples  $FP$ , and incorrectly classified negative examples  $FN$ .

$$\text{Precision: } P = \frac{TP}{TP + FP} \quad (1)$$

$$\text{Recall: } R = \frac{TP}{TP + FN} \quad (2)$$

F-score is a weighted sum of Precision and Recall:

$$F = \frac{(\beta^2 + 1)P}{(\beta^2 + 1)P + R} \quad (3)$$

In some cases authors use the traditional Accuracy, which we cite:

$$A = \frac{TP + TN}{TP + TN + FP + FN} \quad (4)$$

The remainder of the chapter is organized in four sections and some supplementary material. First, in three consecutive sections, we discuss ML applications in specific NLP mega-text problems. We concentrate on the correspondence between NLP problems and ML algorithms. The NLP problems arise from the need to analyze a high volume of electronic texts in marketing, mass media, and health care. Each problem is well served by different algorithms. Opinion analysis, originally a marketing problem, is mature enough to have become almost a traditional task in text analysis. On the other hand, privacy protection, especially urgent in health care, is the subject of cutting-edge research at the boundary of

21 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:  
[www.igi-global.com/chapter/machine-learning-applications-mega-text/36992](http://www.igi-global.com/chapter/machine-learning-applications-mega-text/36992)

## Related Content

---

### Nature Inspired Methods for Multi-Objective Optimization

Sanjoy Das, Bijaya K. Panigrahi and Shyam S. Pattnaik (2010). *Handbook of Research on Machine Learning Applications and Trends: Algorithms, Methods, and Techniques* (pp. 95-108).

[www.irma-international.org/chapter/nature-inspired-methods-multi-objective/36981](http://www.irma-international.org/chapter/nature-inspired-methods-multi-objective/36981)

### Integrating AI With IoT: Challenges and Solutions

Poornima Pandian, Anu Disney D., Lavanya Devi S. and Gowri Manohari R. (2025). *Merging Artificial Intelligence With the Internet of Things* (pp. 1-32).

[www.irma-international.org/chapter/integrating-ai-with-iot/379403](http://www.irma-international.org/chapter/integrating-ai-with-iot/379403)

### On the Cognitive Complexity of Software and its Quantification and Formal Measurement

Yingxu Wang (2009). *International Journal of Software Science and Computational Intelligence* (pp. 31-53).

[www.irma-international.org/article/cognitive-complexity-software-its-quantification/2792](http://www.irma-international.org/article/cognitive-complexity-software-its-quantification/2792)

### Harnessing Intelligent RIS for Optimized Capacity and Latency in 6G Cooperative NOMA Systems: A Phase Shift Control Approach

Mohamed Hassan, Khalid Hamid, Hashim Elshafie, Elmuntaser Hassan, Rashid A. Saeed, Hesham Alhumyani and Abdullah Alenizi (2024). *International Journal of Software Science and Computational Intelligence* (pp. 1-23).

[www.irma-international.org/article/harnessing-intelligent-ris-for-optimized-capacity-and-latency-in-6g-cooperative-noma-systems/366587](http://www.irma-international.org/article/harnessing-intelligent-ris-for-optimized-capacity-and-latency-in-6g-cooperative-noma-systems/366587)

### Designing a Hybrid Approach for Web Recommendation Using Annotation

Sunny Sharma, Vijay Rana and Vivek Kumar (2022). *Applications of Computational Science in Artificial Intelligence* (pp. 234-247).

[www.irma-international.org/chapter/designing-a-hybrid-approach-for-web-recommendation-using-annotation/302069](http://www.irma-international.org/chapter/designing-a-hybrid-approach-for-web-recommendation-using-annotation/302069)