

From Extractive to Abstractive Summarization: A Journey

Parth Mehta · Prasenjit Majumder

From Extractive to Abstractive Summarization: A Journey

Parth Mehta
Information Retrieval and Language
Processing Lab
Dhirubhai Ambani Institute of Information
and Communication Technology
Gandhinagar, Gujarat, India

Prasenjit Majumder
Information Retrieval and Language
Processing Lab
Dhirubhai Ambani Institute of Information
and Communication Technology
Gandhinagar, Gujarat, India

ISBN 978-981-13-8933-7

ISBN 978-981-13-8934-4 (eBook)

<https://doi.org/10.1007/978-981-13-8934-4>

© Springer Nature Singapore Pte Ltd. 2019

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, expressed or implied, with respect to the material contained herein or for any errors or omissions that may have been made. The publisher remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

This Springer imprint is published by the registered company Springer Nature Singapore Pte Ltd.
The registered company address is: 152 Beach Road, #21-01/04 Gateway East, Singapore 189721,
Singapore

Preface

Research in the field of text summarisation has primarily been dominated by investigations of various sentence extraction techniques with a significant focus towards news articles. In this book, we intend to look beyond generic sentence extraction and instead focus on domain-specific summarisation, methods for creating ensembles of multiple extractive summarisation techniques and using sentence compression as the first step towards abstractive summarisation.

We begin by proposing two new corpora, related to legal and scientific articles, for domain-specific summarisation. The first corpus is a collection of judgements delivered by the Supreme Court of India, with corresponding handwritten summaries written by legal experts. The second dataset is a collection of scientific articles related to the domain of computational linguistics and indexed in the ACL anthology. These two tasks highlight the challenges in domain-specific summarisation. The legal summaries are recall-oriented (trying to retain as much information as possible) and semi-extractive (reuses phrases and sentences from the original document). In comparison, the abstracts of ACL articles are more precision-oriented and abstractive. Both collections have a reasonable number of article-summary pairs, enabling us to use data-driven techniques. We annotate both these corpora through a completely automatic process described without requiring any manual intervention. Excluding newswire corpora, where the summaries are usually article headlines, the proposed collections are among the largest openly available collections for document summarisation. We propose a completely data-driven technique for sentence extraction for legal and scientific articles. In both legal and ACL corpora, the summaries have a predefined format. Hence, it is possible to identify *summary worthy* sentences depending on whether they contain certain key phrases. Our proposed approach uses attention-based neural network to automatically identify these key phrases from pseudo-labelled data, without requiring any annotation or handcrafted rules. The proposed model outperforms the existing baselines and state-of-the-art systems by a large margin.

Next, we explore methods for combining several sentence extraction techniques to build an ensemble. A large number of sentence extraction techniques exist, none of which guarantee better performance than the others. As a part of this book, we

explore whether it is possible to leverage this variance in performance for generating an ensemble of several extractive techniques. In the first model, we study the effect of using multiple sentence similarity scores, ranking algorithms and text representation techniques. We demonstrate that such variations can be used for improving rank aggregation, leading to better sentence rankings. Using several sentence similarity metrics, with any given ranking algorithm, always generates better abstracts. We also propose several content-based aggregation models. Given the variation in performance of extractive techniques across documents, the a priori knowledge about which technique would give the best result for a given document can drastically improve the result. In such case, an oracle ensemble system can be made which chose the best possible summary for a given document. In the proposed content-based aggregation models, we estimate the probability of a summary being good by looking at the amount of content it shares with other candidate summaries. We present a hypothesis that a good summary will necessarily share more information with another good summary, but not with a bad summary. We build upon this argument to construct several content-based aggregation techniques, achieving a substantial improvement in the ROUGE scores.

In the end, we propose another attention-based neural model for sentence compression. We use a novel context encoder, which helps the network to handle rare but informative terms better. We contrast our work to some sentence compression and abstractive techniques that have been proposed in the past few years. We present our arguments for and against these techniques and build a further road map for abstractive summarisation. In the end, we present the results on an end-to-end system which performs sentence extraction using stand-alone summarisation systems as well as their ensembles and then uses the sentence compression technique for generating the final abstractive summary. Overall, this book is aimed at providing alternates to generic extractive summarisation in the form of domain-specific as well as ensemble techniques, gradually leading to sentence compression and abstractive summarisation.

Gandhinagar, India

Parth Mehta
Prasenjit Majumder

Contents

1	Introduction	1
1.1	Book Organisation	1
1.2	Types of Summarisation Techniques	3
1.3	Extractive Summarisation	3
1.4	Information Fusion and Ensemble Techniques	4
1.5	Abstractive Summarisation	6
1.6	Main Contributions	7
	References	8
2	Related Work	11
2.1	Extractive Summarisation	11
2.2	Ensemble Techniques for Extractive Summarisation	13
2.3	Sentence Compression	16
2.4	Domain-Specific Summarisation	18
2.4.1	Legal Document Summarisation	18
2.4.2	Scientific Article Summarisation	19
	References	21
3	Corpora and Evaluation for Text Summarisation	25
3.1	DUC and TAC Datasets	26
3.2	Legal and Scientific Article Dataset	27
3.3	Evaluation	29
3.3.1	Precision and Recall	29
3.3.2	BLEU	30
3.3.3	ROUGE Measure	31
3.3.4	Pyramid Score	31
3.3.5	Human Evaluation	32
	References	33

4 Domain-Specific Summarisation	35
4.1 Legal Document Summarisation	36
4.1.1 Boosting Legal Vocabulary Using a Lexicon	36
4.1.2 Weighted TextRank and LexRank	37
4.1.3 Automatic Keyphrase Identification	38
4.1.4 Attention-Based Sentence Extractor	39
4.2 Scientific Article Summarisation	42
4.3 Experiment Details	43
4.3.1 Results	44
4.4 Conclusion	47
References	47
5 Improving Sentence Extraction Through Rank Aggregation	49
5.1 Introduction	49
5.2 Motivation for Rank Aggregation	51
5.3 Analysis of Existing Extractive Systems	52
5.3.1 Experimental Setup	54
5.4 Ensemble of Extractive Summarisation Systems	57
5.4.1 Effect of Informed Fusion	58
5.5 Discussion	63
5.5.1 Determining the Robustness of Candidate Systems	64
5.5.2 Qualitative Analysis of Summaries	65
References	67
6 Leveraging Content Similarity in Summaries for Generating Better Ensembles	69
6.1 Limitations of Consensus-Based Aggregation	69
6.2 Proposed Approach for Content-Based Aggregation	71
6.3 Document Level Aggregation	72
6.3.1 Experimental Results	73
6.4 Sentence Level Aggregation	74
6.4.1 SentRank	75
6.4.2 GlobalRank	75
6.4.3 LocalRank	76
6.4.4 HybridRank	77
6.4.5 Experimental Results	77
6.5 Conclusion	79
References	80
7 Neural Model for Sentence Compression	83
7.1 Sentence Compression by Deletion	84
7.2 Sentence Compression Using Sequence to Sequence Model	85
7.2.1 Sentence Encoder	86
7.2.2 Context Encoder	86

7.2.3 Decoder	87
7.2.4 Attention Module	87
7.3 Exploiting SMT Techniques for Sentence Compression	87
7.4 Results for Sentence Compression	88
7.5 Limitations of Sentence Compression Techniques	89
7.6 Overall System	92
References	94
8 Conclusion	97
References	98
Appendix A: Sample Document-Summary Pairs from DUC, Legal and ACL Corpus	99
Appendix B: The Dictionary Built Using Legal Boost Method	105
Appendix C: Summaries Generated Using Rank Aggregation	107
Appendix D: Summaries Generated Using Content-Based Aggregation	111
Appendix E: Visualising Compression on Sentences from Legal Documents	115

About the Authors

Dr. Parth Mehta completed his M.Tech. in Machine Intelligence and his Ph.D. in Text Summarization at Dhirubhai Ambani Institute of ICT (DA-IICT), Gandhinagar, India. At the DA-IICT he was part of the Information Retrieval and Natural Language Processing Lab. He was also involved in the national project “Cross Lingual Information Access”, funded by the Govt. of India, which focused on building a cross-lingual search engine for nine Indian languages.

Dr. Mehta has served as reviewer for the journals Information Processing and Management and Forum for Information Retrieval Evaluation. Apart from several journal and conference papers, he has also co-edited a book on text processing published by Springer.

Prof. Prasenjit Majumder is an Associate Professor at Dhirubhai Ambani Institute of ICT (DA-IICT), Gandhinagar and a Visiting Professor at the Indian Institute of Information Technology, Vadodara (IIIT-V). Professor Majumder completed his Ph.D. at Jadavpur University in 2008 and worked as a postdoctoral fellow at the University College Dublin, prior to joining the DA-IICT, where he currently heads the Information Retrieval and Language Processing Lab. His research interests lie at the intersection of Information Retrieval, Cognitive Science and Human Computing Interaction. He has headed several projects sponsored by the Govt. of India.

He is one of the pioneers of the Forum for Information Retrieval Evaluation (FIRE), which assesses research on Information Retrieval and related areas for South Asian languages. Since being founded in 2008, FIRE has grown to become a respected conference, drawing participants from across the globe. Professor Majumder has authored several journal and conference papers, and co-edited two special issues of Transactions in Information Systems (ACM). He has co-edited two books: ‘Multi Lingual Information Access in South Asian Languages’ and ‘Text Processing,’ both published by Springer.