

Metagenomic sample clustering enhanced by SRA metadata and NLP

Rikke Møller Larsen

201607214@post.au.dk

Supervised by Associate Professor Thomas Bataillon
&

Co-supervised by Patrick Ettenhuber and Fernando Tavares at QIAGEN

Bioinformatics Research Centre,
Universitetsbyen 81
8000 Aarhus C

Aarhus University

June 1, 2022



AARHUS UNIVERSITY



Preface

The work within this Master's thesis was carried out under the supervision of Associate Professor Thomas Bataillon at the Bioinformatics Research Centre at Aarhus University. The project was in collaboration with QIAGEN Aarhus, with Patrick Ettenhuber and Fernando Tavares, (Senior Product Owner and Project Manager/Senior Software Developer, respectively) from QIAGEN as co-supervisors.

Thank you to my supervisor Thomas Bataillon for trusting me to work independently, while also offering helpful counsel and direction, whenever I needed it.

A heartfelt thank you to QIAGEN and its employees for their hospitality through these four months, in particular to Patrick Ettenhuber and Fernando Tavares, and the rest of the Microbial Team. Thank you for your guidance and advice, and for believing in me enough to invite me to stay for the coming year.

I also thank the professors, staff and my fellow Master's students at BiRC for making it two wonderful years despite an ongoing pandemic.

I would like to thank Rikke Myrhøj Jensen, Stine Katrine Rye Østergaard, Janne Auning Engestoft and Benjamin Nichum for their continued support, both moral and academic, and most important of all, their friendship.

Abstract

The data available from metagenomic sequencing experiments are exponentially increasing, in accordance with an increased focus on microbiomes and their effect on health and disease. Databases, such as NCBI's Sequence Read Archive (SRA) seek to make this data publicly and easily accessible with the use of standardized metadata. This metadata can contain important information for the description and distinction of experiments, but given its text-based form, it is not directly available for computer-based methods.

Utilizing biomedical natural language processing (NLP) to harness this data, has the potential to improve the ability of automatically classifying new experiments relative to using only analysis data, which ultimately can help facilitate novel research and discoveries.

While the metadata may contain valuable information, it is also prone to be faulty or incomplete, so there is a need for the development of unsupervised corrective methods, in order to obtain training data of high enough quality.

Abbreviations

| | |
|---|--|
| CBOW - Continuous Bag-of-words | OOV - Out-of-vocabulary |
| CRF - conditional random field | OTU - Operational taxonomic unit |
| LCA - Lowest common ancestor | QMI-PTDB - QIAGEN Microbial Insights - Prokaryotic Taxonomy Database |
| LSTM - Long short term memory | SRA - Sequence read archive |
| MeSH - Medical subject headings | S2V - sent2vec |
| MIMIC-III - Medical Information Mart for Intensive Care | UMAP - Uniform Manifold Approximation and Projection |
| NCBI - National Center for Biotechnology Information | WGS - Whole genome sequencing |
| NGS - Next-generation sequencing | W2V - word2vec |
| NLP - Natural language processing | ΔWB - Delta within-between |

Contents

| | |
|---|-----------|
| Abstract | ii |
| Abbreviations | ii |
| 1 Introduction | 1 |
| 1.1 Metagenomic analysis | 1 |
| 1.1.1 Next-generation sequencing | 2 |
| 1.1.2 Taxonomic profiling | 3 |
| 1.1.3 Sequence Read Archive | 5 |
| 1.2 NLP | 6 |
| 1.2.1 Word embedding | 6 |
| 1.2.2 Sentence embedding | 8 |
| 1.3 Cluster evaluation | 8 |
| 1.3.1 Visualization | 9 |
| 1.3.2 Distance metric | 9 |
| 1.3.3 Silhouette score | 11 |
| 2 Materials and methods | 12 |
| 2.1 Data acquisition | 12 |
| 2.2 Data preprocessing | 12 |
| 2.2.1 Sparse vectors | 13 |
| 2.2.2 Label assignment | 14 |
| 2.3 Text embedding | 15 |
| 2.4 Clustering | 16 |
| 2.4.1 The alternative method | 16 |
| 3 Results | 18 |
| 3.1 Abundance data | 18 |
| 3.2 Text embeddings | 21 |
| 3.3 Combining abundance and word vectors | 23 |
| 4 Discussion | 26 |
| 4.1 Choice of distance metric and data clean-up | 26 |
| 4.2 Improving metadata coverage | 28 |
| 4.3 Future NLP work | 28 |
| 5 Conclusion | 30 |
| References | 31 |
| A Hierarchical labels | 34 |
| B UMAP plots | 36 |