

Domain Understanding

Predictive Analysis of Antibiotic Resistance in Bacterial Strains via DNA Sequence Motifs

Author: Michał Raczkowski
Student number: 4465024
Fontys ICT - AI-core-AI4-RB02

September 22, 2023

Contents

1	Introduction	2
2	Background	2
3	Domain Definition	2
4	Problem Statement	2
5	User & Stakeholder Analysis	3
6	Domain Dynamics	3
7	Data Understanding	3
8	Existing Solutions and Tools	3
9	Domain Challenges and Pain Points	3
10	Opportunities for Improvement	4
11	Recommendations and Next Steps	4
12	Appendices	4
13	Analytic Approach	4
14	Potential Data Sources and File Formats	5
15	File Format Descriptions	5
16	Expected Results and Deliverables	5

1 Introduction

The evolution of antibiotic resistance in bacteria poses a significant challenge to modern medicine. This research aims to ascertain whether specific DNA sequence motifs can be used as markers to predict antibiotic resistance in bacterial strains.

2 Background

Antibiotics have historically been instrumental in combating bacterial infections. However, with the emergence of antibiotic-resistant strains, there is an urgent need to understand the genetic markers responsible for this resistance. This work seeks to investigate the genetic blueprints of these strains for potential predictive patterns.

3 Domain Definition

Below are key terminologies pertinent to this study:

DNA (Deoxyribonucleic Acid): The molecular carrier of genetic instructions in all known living organisms and many viruses.

Genes: DNA segments responsible for encoding functional products, typically proteins.

DNA Motifs: Recurrent DNA patterns with potential biological significance, often recognized and bound by specific proteins.

Antibiotic Resistance: A phenomenon where bacteria evolve mechanisms to mitigate the effects of antibiotics.

Chromosome: DNA molecules containing an organism's genetic material.

Genome: The entirety of an organism's genetic content.

4 Problem Statement

The rise of antibiotic-resistant bacterial strains necessitates innovative methods for timely detection and effective treatment. By examining DNA sequence motifs, this study endeavors to predict the likelihood of antibiotic resistance in various bacterial strains.

5 User & Stakeholder Analysis

Users: Molecular biologists, microbiologists, and medical practitioners.

Needs: A robust method for predicting antibiotic resistance based on genetic markers.

Expectations: An efficient tool that provides accurate predictions with minimal errors.

6 Domain Dynamics

The landscape of genomics and molecular biology is rapidly evolving. Advancements in sequencing technologies and bioinformatics tools continually provide novel insights into bacterial genomics and the dynamics of antibiotic resistance.

7 Data Understanding

The success of this research hinges on obtaining comprehensive DNA sequences from diverse bacterial strains. Fortunately, there are vast repositories where such sequences are collated and made available for research.

8 Existing Solutions and Tools

Currently, bacterial antibiotic resistance is often determined via antibiotic susceptibility tests. Some bioinformatics tools can analyze specific DNA sequences, but a comprehensive tool for predicting antibiotic resistance remains elusive.

9 Domain Challenges and Pain Points

The genetic mechanisms underpinning antibiotic resistance are complex and multifaceted. Moreover, environmental factors can further complicate the predictive analysis. Multi-drug resistant strains present additional challenges due to their resistance to multiple antibiotics.

10 Opportunities for Improvement

By leveraging modern computational methods and large genomic datasets, it's possible to develop a more comprehensive tool that can predict antibiotic resistance based on DNA motifs, which would revolutionize bacterial diagnostics and treatment strategies.

11 Recommendations and Next Steps

- Collaborate with genomic repositories to access a wide array of bacterial DNA sequences.
- Catalog DNA motifs associated with known antibiotic resistance genes.
- Evaluate computational models' efficacy in predicting antibiotic resistance using collected DNA sequences.

12 Appendices

Supplementary information, including detailed descriptions of DNA motifs and other pertinent findings, will be included in this section.

13 Analytic Approach

1. **Data Collection:** Accumulate relevant data, emphasizing sequences associated with antibiotic resistance genes.
2. **Data Cleaning:** Ensure dataset consistency, noise reduction, and handle missing values.
3. **Feature Extraction:** Pinpoint DNA motifs consistently linked with specific antibiotic resistance.
4. **Model Building:** Utilize machine learning algorithms for predictive model training. Algorithms such as Random Forests, SVMs, or deep learning models like CNNs could be apt given the data's sequential nature.
5. **Validation and Testing:** Employ a rigorous validation process using separate training and test datasets to ensure model accuracy.
6. **Deployment:** Integrate a validated model into diagnostic platforms for real-time predictions.

14 Potential Data Sources and File Formats

NCBI's Antibiotic Resistance Genes Database (ARDB): Comprehensive database of antibiotic resistance-related sequences.

- **File Extensions:** '.fasta', '.gb' (GenBank format)

ResFinder: Database for identifying antibiotic resistance genes in bacterial DNA.

- **File Extensions:** '.fasta'

PATRIC: Extensive bacterial genome database.

- **File Extensions:** '.fasta', '.gff' (General Feature Format)

Local Health Departments: Region-specific bacterial DNA sequences.

- **File Extensions:** Can vary, but common formats include '.fasta', '.fastq', and '.vcf'.

15 File Format Descriptions

.fasta: A common format for DNA, RNA, or protein sequences. Each entry has a description line followed by sequences.

.gb or .genbank: GenBank format. Used for coding sequences, RNA sequences, and provides annotations related to the sequences.

.gff: General Feature Format. Used to describe genes and other features of DNA, RNA, and protein sequences.

.fastq: Provides both the sequence of the read fragments and their corresponding quality scores. Commonly used in next-generation sequencing.

.vcf: Variant Call Format. Used in bioinformatics for storing gene sequence variations.

16 Expected Results and Deliverables

At the culmination of this research, the anticipated outcomes are:

- A comprehensive catalog of DNA motifs associated with antibiotic resistance.

- A validated predictive model capable of ascertaining antibiotic resistance likelihood based on analyzed sequences.
- A user-friendly platform enabling researchers and medical practitioners to utilize the predictive tool in real-time.