





ARSENAL:

Antimicrobial resistance prediction by machine learning approach

Ulysse Guyet

Macha Nikolski & Alexis Groppi

Computational Biology and Bioinformatics team, IBGC, UMR 5095

Machine Learning for Microbial Genomics Workshop – 23/09/2022

TABLE OF CONTENTS

- 1. Introduction
- 2. Dataset description
- 3. Genome assembly & annotation
- 4. ARSENAL : machine learning pipeline for MIC prediction and linking phenotype and genotype
- 5. ARSENAL : Preliminary results for *Streptococcus pneumoniae*
- 6. Future work and questions

TABLE OF CONTENTS

1. Introduction

- 2. Dataset description
- 3. Genome assembly & annotation

4. ARSENAL : machine learning pipeline for MIC prediction and linking phenotype and genotype

5. ARSENAL : Preliminary results for *Streptococcus pneumoniae*

6. Future work and questions

Antibiotic resistance: a major health problem



Antibiotic resistance mechanisms



- Resistance mechanisms are mainly acquired by genome modifications: e.g. mutations, horizontal gene transfers
- Misuse or overuse of antibiotics increases the drug resistance by selection pressure
- Resistance is measured by MIC (Minimum Inhibitory Concentration) which is the minimum concentration of an antibiotic needed to inhibit bacterial growth

Antibiotic resistance mechanisms



Predicting antibiotic resistance can allows:

- Fast screening of antibiotics to determine the most effective antibiotic and the right dose against a specific bacterial infection
- Highlight new antibiotic resistance genes and biomarkers

vlacrolides	Macrolides
Aminoglycosides	Rifamycins

- Resistance mechanisms are mainly acquired by genome modifications: e.g. mutations, horizontal gene transfers
- Misuse or overuse of antibiotics increases the drug resistance by selection pressure
- Resistance is measured by MIC (Minimum Inhibitory Concentration) which is the minimum concentration of an antibiotic needed to inhibit bacterial growth

TABLE OF CONTENTS

1. Introduction

- 2. Dataset description
- 3. Genome assembly & annotation

4. ARSENAL : machine learning pipeline for MIC prediction and linking phenotype and genotype

5. ARSENAL : Preliminary results for *Streptococcus pneumoniae*

6. Future work and questions

Antimicrobial resistant bacteria models

Streptococcus pneumoniae

- Genome size: ~2 Mbp
- Gram-positive
- WHO antimicrobial (AMR) priority pathogen
- Resistance acquired by either **mutations** or by DNA **transformation** leading to mosaic genes

Pseudomonas aeruginosa

- Genome size: 5.5-7 Mbp
- Gram-negative
- WHO AMR priority pathogen
- Resistance mainly mediated by **mutations** or horizontal gene transfer (mostly plasmid **conjugation**)

Antimicrobial resistant bacteria models

Streptococcus pneumoniae

- Genome size: ~2 Mbp
- Gram-positive
- WHO antimicrobial (AMR) priority pathogen
- Resistance acquired by either **mutations** or by DNA **transformation** leading to mosaic genes

Pseudomonas aeruginosa

- Genome size: 5.5-7 Mbp
- Gram-negative
- WHO AMR priority pathogen
- Resistance mainly mediated by **mutations** or horizontal gene transfer (mostly plasmid **conjugation**)

Datasets

1312 *Streptococcus pneumoniae* newly sequenced genomes of strains isolated from patients

- Different sequence types (ST) and 64 distinct serotypes
- many strains are multidrug resistant (MDR)
- Geographical origin (isolated from 28 hospitals in 18 provinces of China between 2007 and 2020)

□ MIC values (Minimum Inhibitory Concentration) of each strain for 11 various antibiotic classes :

- β-lactams
- macrolides
- Lincosamide
- Fluoroquinolones
- Oxazolidinone
- glycopeptides

- Rifampin
- Sulfonamides
- Chloramphenicol
- Phosphonic
- aminosides

TABLE OF CONTENTS

- 1. Introduction
- 2. Dataset description
- 3. Genome assembly & annotation

4. ARSENAL : machine learning pipeline for MIC prediction and linking phenotype and genotype

5. ARSENAL : Preliminary results for *Streptococcus pneumoniae*

6. Future work and questions

Genome assembly & annotation pipeline



PATRIC (Pathosystems Resource Integration Center) Database :

- Prediction of CDS and functional annotation
- Assignment of each gene to a PLFam: group of genus-specific genes, wich share the same function and high sequence homology

5770 PLFams by considering all genomes of *Streptococcus pneumoniae*

TABLE OF CONTENTS

- 1. Introduction
- 2. Dataset description
- 3. Genome assembly & annotation

4. ARSENAL : machine learning pipeline for MIC prediction and linking phenotype and genotype

5. ARSENAL : Preliminary results for *Streptococcus pneumoniae*

6. Future work and questions



Machine learning pipeline for MIC prediction and linking phenotype and genotype





One k-mer occurrences table for each of the PLFam

2

Genome n





1

8

2

0.5

... 2

Genome 1

Genome 2

Genome 3

Genome 4

Genome n

ARSENAL pipeline



Machine learning pipeline for MIC prediction and linking phenotype and genotype





Genome-distance-based cross-validation

- Maximize genome distance between the test sets of cross-validation folds
- Improve independence of test sets by segregating samples based on a known dependence structure in the data (here genomic distance)

rirrors the application of the trained model towards independently sampled datasets

TABLE OF CONTENTS

- 1. Introduction
- 2. Dataset description
- 3. Genome assembly & annotation

4. ARSENAL : machine learning pipeline for MIC prediction and linking phenotype and genotype

5. ARSENAL : Preliminary results for *Streptococcus pneumoniae*

6. Future work and questions

Streptococcus pneumoniae

β-lactams:

Antibiotic Class	Antibiotic	Accuracy (±1 cat.)
Penicillins	Penicillin	0.832
	Amoxicillin	0.897
Cephalosporins	Cefuroxime	0.829
	Cefepime	0.848
	Ceftriaxone	0.859
Carbapenem	Imipinem	0.935
	Meropenem	0.898
	Ertapenem	0.939

	Antibiotic Class	Antibiotic	Accuracy (±1 cat.)
	Fluorquinolones	Ciprofloxacin	0.848
		Levofloxacin	0.949
		Moxifloxacin	0.909
	Rifampin	Rifampicin	0.838
	phosphonic	phosphonomycin	0.838

SHAPley value (SHapley Additive exPlanations):

Determine the top genes that have a positive effect on the model construction

Among genes of interest :

- The 3 Penicillin-binding proteins of *Streptococcus pneumoniae* are found in 10th, 14th and 16th position
- Significant number of mobile elements (1st, 2nd; 12 genes in the first 100 most important genes)
- choline binding proteins- associated with drug and cell death (4th, 13th, 31th position)
- Cell wall surface anchor family protein (15th)



Identification of these resistance-related genes provides confidence in model predictions

TABLE OF CONTENTS

- 1. Introduction
- 2. Dataset description
- 3. Genome assembly & annotation
- 4. ARSENAL : machine learning pipeline for MIC prediction and linking phenotype and genotype
- 5. ARSENAL : Preliminary results for *Streptococcus pneumoniae*
- 6. Future work and questions

Conclusions

- The model is accurate for β-lactams antibiotics in the case of *Streptococcus pneumoniae*
- Validation of the method by the identification of genes known to be involved in resistance mechanisms
- Highlight some genes without functional annotation potentially linked with resistance

Functional characterization via CRISPR/Cas9

Future computational work

- Test the method for other antibiotics classes and for *Pseudomonas aeruginosa* (data from V. Dubois lab, Bordeaux & J. Feng lab, Tianjin - China)
- Add tRNAs and promoters to the model since they may influence some resistance mechanisms

Long-term perspectives

• Rapid identification of resistance markers in patients based on genomic & metabolic markers



Integration of metabolic data (in process of generation)

Acknowledgements

MAGITICS project collaborators

Jacques Corbeil (Université Laval, Québec, Canada) Jie Feng (Institute of microbiology, Tianjin, China) Véronique Dubois (CNRS, Bordeaux) Léa Bientz (CNRS, Bordeaux)

CBiB/IBGC

Macha Nikolski

Alexis Groppi

Aurélien Barré

Emmanuel Bouilhol

Edgar Lefevre

Grigorii Sukhorukov

Thanks for your attention !





