

BenchmarkDR

A modular and expandable benchmarking pipeline for machine learning based antimicrobial resistance prediction

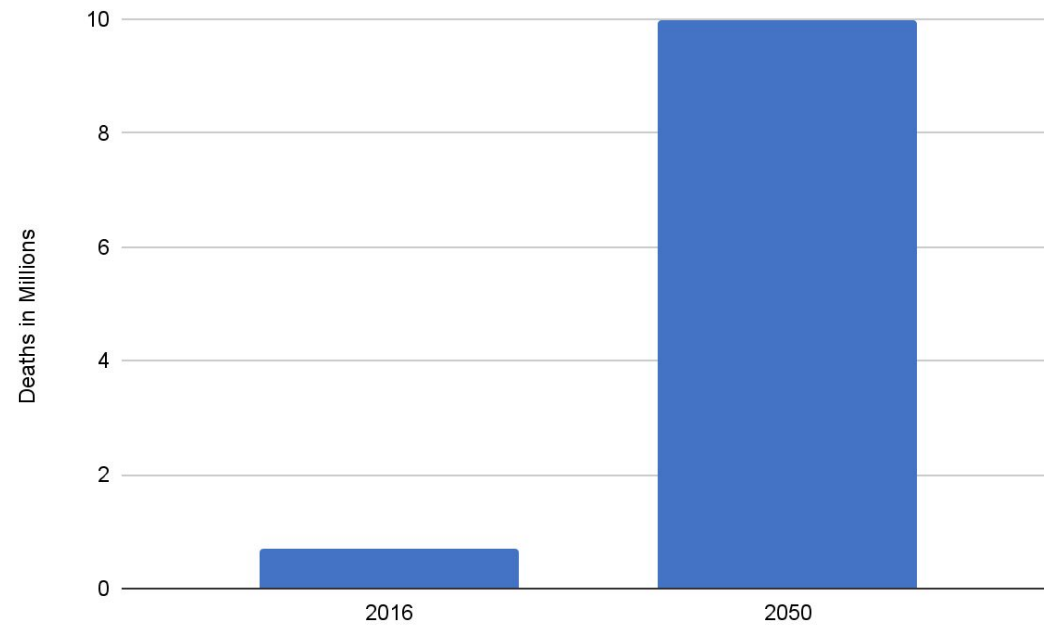
Niklas Stotzem*, Fernando Guntoro*, Leonid Chindelevitch*

*Imperial College London

Multi-drug-resistant pathogenic bacteria are an increasing global threat

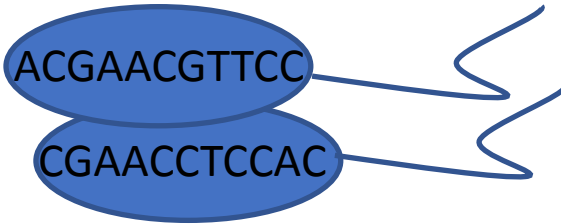


Deaths linked to multi-drug-resistant (MDR) microbes (J. O'Neill, 2016)



Machine Learning (ML) is a useful tool to fight pathogenic drug resistant bacteria

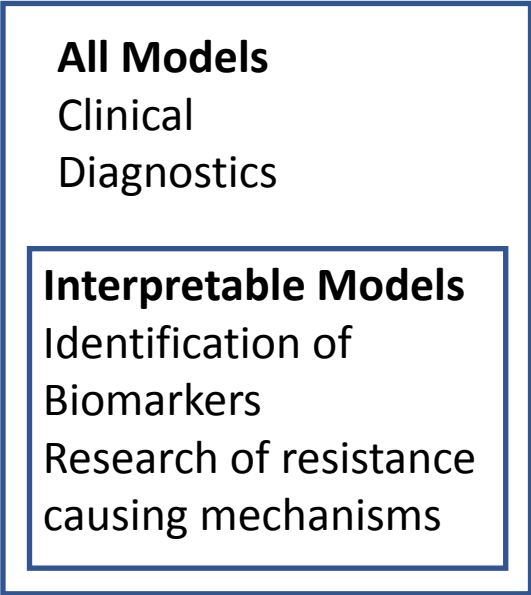
Genomic NGS Data



Drug Susceptibility Test Data

Bacterial Isolate	Antibiotic 1	...	
1	Resistant	...	
2	Bacterial Isolate	Antibiotic 1	...
...			
1	MIC X	...	
2	MIC Y	...	
...	

ML Methods for
Classification and
Regression



Different features of the genome as well as different ML methods have been used by researchers

Representations

Genes

...ATCAAATCCGTTTCAAGGTCCCTTGCCAACCGGTTGGAACG

Single Nucleotide Polymorphisms (SNPs)

...ATCAAATCGTTTCAAGGTCCCTTGCCAACCGGTTGGAACG

A G

K-MERs

...ATCAAATCCGTTTCAAGGTCCCTTGCCAACCGGTTGGAACG

ATCA

TCAA

.....

ML Methods

Logistic Regression

Linear Regression

Random Forests

Support Vector Machines

...

New Ones

How to compare methods performance on different representations of different data?

We built an end-to-end pipeline to allow users to compare representations and methods



Genomic Raw Data

ACGAACGTTCC
CGAACCTCCAC

User Configuration

Mode: Classification/ Regression
Methods: Machine Learning Method 1 (ML1), ...
Method Parameters: ML1: Param 1, Param 2, ...
Hyperparameter-Tuning: Grid Search, ...
Representation: Gene Presence/ Absence, ...

Drug Susceptibility Test Data

Bacterial Isolate	Antibiotic 1	...	
1	Resistant	...	
2	Bacterial Isolate	Antibiotic 1	...
.. 1	MIC X	...	
2	MIC Y	...	
...	

Bacterial Isolate	K-MER 1	...
1
2	SNP 1	...
1
Bacterial Isolate	Gene 1	...
1	Present	...
2	Absent	...
...

Method	Evaluation Parameter 1	...
Method 1	Score 1.1	...
Method 2	Score 1.2	...
...

Input

Representation

Evaluation

What does the pipeline have under the hood?

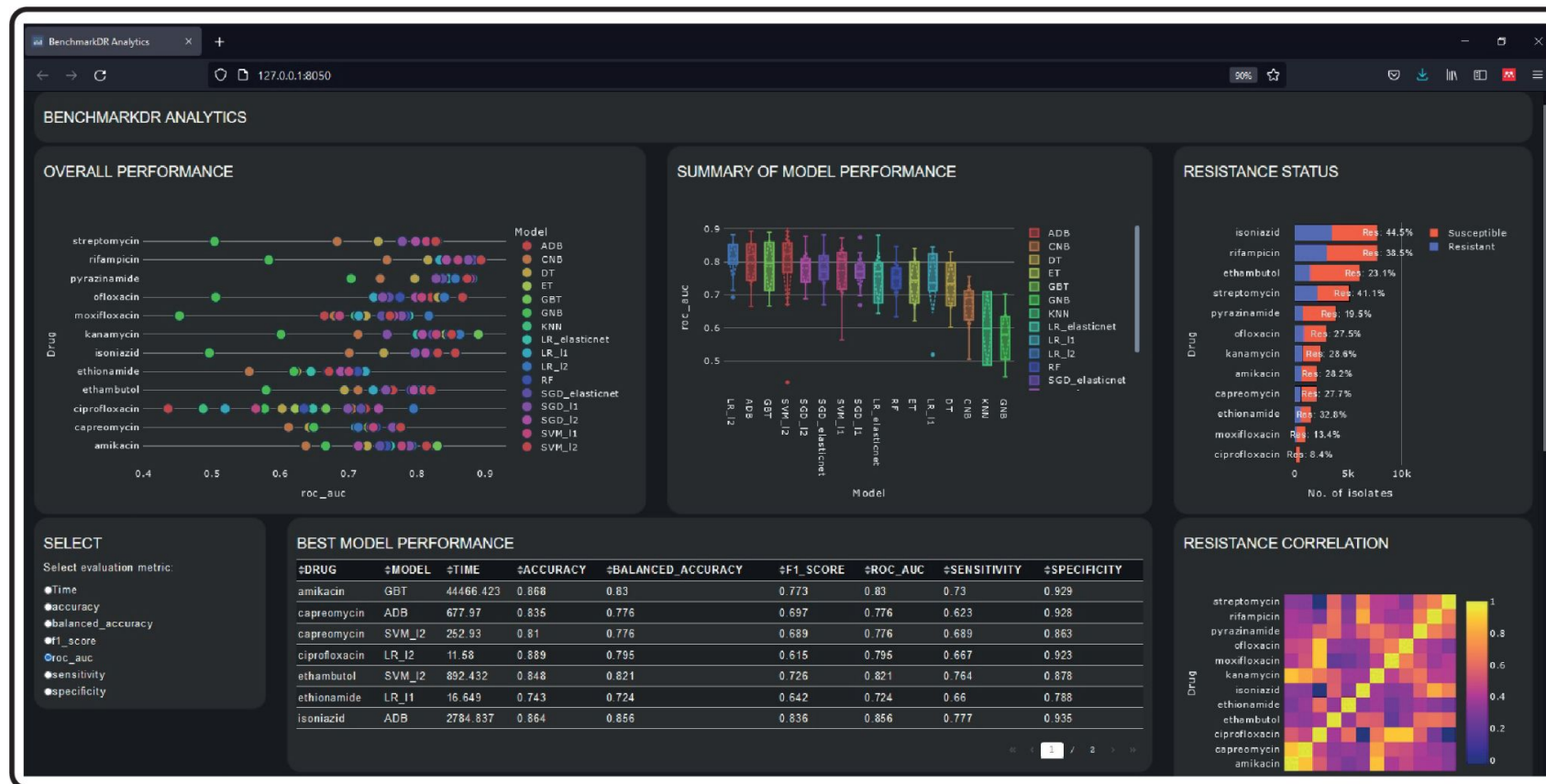
- **Representations**
 - Tools based on benchmark papers, popularity and ease of use
- **ML Methods**
 - 17 for Classification (Scikit-Learn (Pedregosa, F. et al., 2011) & INGOT-DR (Zabeti, H. et al., 2021))
 - 9 for Regression
- **Evaluation**
 - Time and Prediction Performance

Configuration can be adapted at different levels of detail

```
1 PATH_DATA: # Path to folder in which pneumonia_mic folder is located as string
2
3 BACTERIA: ["pneumonia_mic"]
4
5 OUTPUT_DIR: # Path to folder in which all further downstream files will be saved as string
6
7 REPRESENTATION: ["gene_presence_absence", "snps", "kmer"]
8
9 MODE: "MIC" # "Classification" or "MIC"
10
11 DRUGS: ['Tobramycin',
12        'Cefuroxime sodium',
13        'Ciprofloxacin',
14        'Gentamicin',
15        'Ampicillin',
16        ]
17
18 METHODS: ["sklearn_LinR", "sklearn_LinR_l1", "sklearn_LinR_l2", "sklearn_LinR_elasticnet"]
19
20 OPTIMIZATION: "GridSearchCV" # "GridSearchCV", "RandomizedSearchCV" and "None"
```

```
127 sklearn_ADBC:
128     module: sklearn.ensemble
129     model: AdaBoostClassifier
130     params:
131         n_estimators: 50
132         learning_rate: 0.0001
133     cv:
134         n_estimators: [50, 100, 500, 1000]
135         learning_rate: [0.0001, 0.001, 0.01, 0.1, 1.0]
136
137 sklearn_GBTC:
138     module: sklearn.ensemble
139     model: GradientBoostingClassifier
140     params:
141         max_features: 'auto'
142         n_iter_no_change: 20
143     cv:
144         learning_rate: [0.001, 0.01, 0.1, 1]
145         n_estimators: [100, 300, 500, 1000]
146         min_samples_split: [2, 5, 10, 15]
147         max_depth: [5, 10, 15, 30]
```


A dashboard prototype provides an convenient overview of results and further details



Conclusions & Future Outlook

- ➔ Easily useable and extensible end-to-end pipeline to benchmark 26 different ML methods on genomic data represented in 3 ways
- ➔ Applicable for other microbial phenotype predictions

Future Work

- Further Extensions (Methods, Representations)
- Provision of comprehensive gold standard datasets
- Adding explainability approaches, e.g. SHAP
- Include long-read sequencing data
- Grow the tool with the community

Acknowledgements

Repository: <https://github.com/WGS-TB/BenchmarkDR>

- Fernando Guntoro
- Dr. Leonid Chindelevitch (Supervision)
- Dr. John Lees (Supervision)
- Dr. Hooman Zabeti (Supervision)



Key References

Guntoro, F. (2021), Benchmarking antimicrobial resistance prediction with an automated pipeline for machine learning models using BenchmarkDR [Master's Thesis, Imperial College London].

J. O'Neill (2016). Review on Antimicrobial Resistance. *Tackling Drugresistant Infections Globally: Final Report and Recommendations*. Review on Antimicrobial Resistance, 2016.

Pedregosa, F. et al. (2011), *Scikit-learn: Machine Learning in Python*, Journal of Machine Learning Research, 12.

Stotzem, N. (2021), An End-to-end Machine Learning Pipeline for Drug Resistance Prediction in Bacteria [Master's Thesis, Imperial College London].

Zabeti, H., Dexter, N., Safari, A. H., Sedaghat, N., Libbrecht, M., & Chindelevitch, L. (2021). *INGOT-DR: an interpretable classifier for predicting drug resistance in M. tuberculosis*. Algorithms for Molecular Biology, 16(1), 17.

Images

<https://theconversation.com/antibiotic-resistance-new-discovery-could-change-the-future-of-treatment-131262>

<https://github.com/snakemake>

Backup

We selected the tools used to create the representations based on published benchmarks, popularity and ease of use

Representation	Tools
Gene Presence/ Absence	SPAdes v3.15.2 (Assembly), Prokka v1.13.4 (Annotation)
Single Nucleotide Polymorphisms (Requires Reference Genome)	BWA v0.1.17 (Alignment), Samtools v1.12 (Sorting), Picard v2.25.6 (Duplicate Removal), VarScan v2.4.4 (SNP calling)
<i>K</i> -mers	KMC v3.1.2.rc1

The pipeline includes - so far - a variety of standard ML methods and INGOT-DR

Binary Classification (Resistant/ Susceptible)	Regression (MIC)
Logistic Regression* Support Vector Machine Classification* Decision Trees Random Forests Extremely Randomized Trees AdaBoost Decision Tree Classifier Gradient Boosted Decision Trees Stochastic Gradient Descent Classifier* K-Nearest Neighbours Gaussian/ Complement Naive Bayes INterpretable GrOup Testing for Drug Re- sistance (INGOT-DR) (Zabeti, H. et al., 2021)	Linear Regression* Support Vector Machine Regression Decision Tree Regressor Random Forest Regressor Gradient Boosted Trees Regressor AdaBoost Decision Tree Regressor

* incl. variations with l1 and l2 regularization

Several evaluation metrics are provided

Classification	Regression
Accuracy, Balanced Accuracy, F1-Score, AUC, Sensitivity, Specificity	Mean Squared Error, Mean Squared Log Error, Coefficient of Determination