AMR and machine-learning

Prediction of AMR from metagenomes among other things

Finlay Maguire finlaymaguire@gmail.com December 3, 2019

Faculty of Computer Science, Dalhousie University

- 1. Genomic Phenotype Prediction
- 2. Non-Bioinformatics Interlude
- 3. AMRtime

Genomic Phenotype Prediction

Antibiotic Susceptibility Testing



Bradley et al. (2015)

AAFC Salmonella Data-set



Genomic RGI Predictions



Linking AMR determinants to Phenotype



Logistic Regression



6

Set-Covering Machines



AST Prediction Performance



A: RGI, B: RGI-efflux, C: Logistic Regression, D: Set Covering Machines. Major Disagreement is overprediction of resistance, Very Major Disagreement is underprediction

Learnt features/weights



Extending beyond Salmonella



ARO Predictions (Kara Tsang)

Extending beyond Salmonella



Logistic Regression

• Using direct annotations works very poorly across different organisms and resistance mechanisms.

- Using direct annotations works very poorly across different organisms and resistance mechanisms.
- Even very simple logistic regression models greatly improve predictions.

- Using direct annotations works very poorly across different organisms and resistance mechanisms.
- Even very simple logistic regression models greatly improve predictions.
- Investigation of learnt weights and features can be very scientifically informative.

Non-Bioinformatics Interlude

• Non-profits have data and lots of contextualising knowledge.

- Non-profits have data and lots of contextualising knowledge.
- No time or resources to analyse or use it

- Non-profits have data and lots of contextualising knowledge.
- No time or resources to analyse or use it
- Informaticians have the skills and resources but no specific understanding of the context.

- Non-profits have data and lots of contextualising knowledge.
- No time or resources to analyse or use it
- Informaticians have the skills and resources but no specific understanding of the context.
- Many low-hanging fruit that can make big differences.

Refugee Women's Health Clinic





Language Development in Autism



Qualitative Social Media Analysis (Tamara Sorenson-Duncan)



Unique Subreddit Posting Activity Diversity

other_submissions/jan_25_052518.json other_submissions/jan_27_013436.json other_submissions/jan_27_164903.json other_submissions/jan_27_031552.json autism.json

Submission

other_submissions/jan_27_130001.json other_submissions/jan_25_014256.json other_submissions/jan_27_063048.json other_submissions/jan_27_000142.json other_submissions/jan_27_152049.json

Beta Diversity of Posting Activity



18

- Halifax Community Learning Network
- Shelter Nova Scotia
- 211 Nova Scotia

AMRtime

AMR-metagenomics



Why is this difficult?



AMR Reads in Metagenome (0.643%)

2184 CARD-Prevalence Genomes at 1-10X abundance

AMR genes have wildly different abundances



1236 AMR PATRIC genomes

AMR genes have highly variable diversity



AMR sequence space overlaps



MDS of CARD Proteins BLASTP-%ID

Insufficient Signal in 250bp Fragments



NDM Multiple Sequence Alignment

Insufficient Signal in 250bp Fragments



NDM Multiple Sequence Alignment
- No point doing what we do if people can't use it.
- Limited hardware requirements (a standard workstation or instance < 8 12Gb, 1 8 cores).
- Fast enough (< 12 hours).
- Easy to install/configure.
- Easy to use.
- Easy to update.

AMRtime

AMRtime structure



Read filtering

Homology Filter Approaches



Relative Computational Demands

Precision-Recall of Homology Search



Optimising for recall



Sensitive Homology Classification

Dealing with imbalanced training data



- Different gene lengths within families (coverage vs read number)?
- Different family sizes?
- Different family diversity?
- Using a generator to improve on SMOTE.





NB 7-mer Average Precision: 0.63



NB 7-mer Average Precision: 0.63 %

Revised classifier structure: exploiting the ARO





Advantages: read length invariant, low dimensionality, uses filtering data computation

- Encodings:
 - Raw sequence
 - Filtering homology search family similarity/dissimilarity
 - Manual feature extraction (GC/TNF/compositional)
 - One-hot K-mer representation
 - K-mer embeddings (DNA2vec/BioVec)
- Classifiers:
 - Random Forests
 - Naive Bayes
 - Logistic Regression
 - Neural Networks of varying architecture (Torch)







Family diversity as explanation?



Within family label imbalance



- Multiset prediction when insufficient signal.
- Systematic benchmarking.
- Full end-to-end comparisons with other approaches (soliciting ideas!)
- rRNA and variant models (not discussed here).
- Integration into CARD platform and IRIDA.

Acknowledgements



- McMaster University: Kara Tsang, Brian Alcock, and Andrew McArthur
- Simon Fraser University: Fiona Brinkman
- Dalhousie University: Robert Beiko
- Funding: Genome Canada, NERC Undergraduate Student Research Award, Donald Hill Family Fellowship

References

- Bradley, P., Gordon, N. C., Walker, T. M., Dunn, L., Heys, S., Huang,
 B., Earle, S., Pankhurst, L. J., Anson, L., De Cesare, M., et al. (2015).
 Rapid antibiotic-resistance predictions from genome sequence data for staphylococcus aureus and mycobacterium tuberculosis. *Nature communications*, 6:10063.
- Huang, Y., Gilna, P., and Li, W. (2009). Identification of ribosomal rna genes in metagenomic fragments. *Bioinformatics*, 25(10):1338–1340.
- Kopylova, E., Noé, L., and Touzet, H. (2012). Sortmerna: fast and accurate filtering of ribosomal rnas in metatranscriptomic data. *Bioinformatics*, 28(24):3211–3217.

- McArthur, A. G., Waglechner, N., Nizam, F., Yan, A., Azad, M. A., Baylay, A. J., Bhullar, K., Canova, M. J., De Pascale, G., Ejim, L., et al. (2013). The comprehensive antibiotic resistance database. *Antimicrobial agents and chemotherapy*, 57(7):3348–3357.
- Schmieder, R., Lim, Y. W., and Edwards, R. (2011). Identification and removal of ribosomal rna sequences from metatranscriptomes. *Bioinformatics*, 28(3):433–435.

Backup

Variant Models

Ribosomal Variant Models



- MetaRNA (Huang et al., 2009)
- Ribopicker (Schmieder et al., 2011)
- SortmeRNA (Kopylova et al., 2012)
- 77 models
- Reads simulated from the underlying 30 species reference genomes

Identifying Ribosomal Reads



Identifying Ribosomal Reads



Identifying Ribosomal Reads



Species Specific Recall

Proportion of Reads

Identifying Taxonomy

True label

Borrelia burgdorferi	26	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0		- 150
Brachyspira hyodysenteriae	1																												0		
Chlamydia psittaci	0																												0		
Chlamydia trachomatis	0																												0		
Chlamydomonas reinhardtii	1			0																									0		
Escherichia coli	0					129	0																						0		120
Halobacterium sp.	0																												0		- 120
Helicobacter pylori	1																												0		
Moraxella catarrhalis	0								68																				0		
Mycobacterium abscessus	0																														
Mycobacterium avium	0																												0		
Mycobacterium chelonae	0																												0		- 90
Mycobacterium intracellulare	0																												0		
Mycobacterium kansasii	0																												1		
Mycobacterium smegmatis	0														46														0		
Mycobacterium tuberculosis	0														0														0		
Mycoplasma fermentans	0															0													0		
Mycoplasma gallisepticum	0																0	41											0		- 60
Mycoplasma hominis	0																	0											0		
Mycoplasma pneumoniae	0																			22	0								0		
Neisseria gonorrhoeae	0																			0		8							0		
Neisseria meningitidis	0																			0	23	76	1						0		
Pasteurella multocida	0																					1	151	0					0		
Propionibacterium acnes	0																						0	62	0				0		- 30
Propionibacterium freudenreich	0					0																			41	0			0		
Salmonella enterica	0					36																				149	0		0		
Staphylococcus aureus	3																										88	0	0		
Streptococcus pneumoniae	0	0	0	0			0		0		0					0		0		0		0			0	0		85	0		
Streptomyces ambofaciens	0	0	0	0	0	0	1	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	147		- 0
																				633											

Predicted label

52

Some are relatively easy



Correct index:390 species:Streptomyces ambofaciens



Misspredict index:750 species:Chlamydomonas reinhardtii
Some are group ambiguous



Misspredict index:259 species:Mycobacterium chelonae

Probably a Mycobacterium?

Others are just a toss-up



Ambiguity in classification



Predicted label

57

Meta-models

- Efflux Pump
- Gene Cluster

