

bgwas3: a pipeline for kmer based association testing in bacteria

Gregory Leeman

Contents

1	Abstract	1
2	Acknowledgements	2
3	Abbreviations	2
4	Introduction	2
4.1	Pangenome wide association studies	2
4.2	The Need for Multiple Association Tests - Metabolomics	3
5	Scope of work	3
6	Methods	4
6.1	CGATcore	4
6.2	Genome Annotation	6
6.3	Kmer minig and counting	6
6.4	Phylogeny estimation and covariate estimation	6
6.5	Kmer association testing	7
6.6	Bonferoni correction	7
6.7	Significant Kmer mapping	7
6.8	Pathway analysis	9

6.9	Visualisation	9
6.10	Scientif computing practices	9
7	Results	9
7.1	Processing of phenotypes	14
7.2	Annotation	14
7.3	Kmer mining	14
7.4	Phenotype	14
8	Discussion	14
8.1	Genome annotation	14
8.2	Kmer mining	14
	References	14

1 Abstract

Genome-wide association studies (GWAS) are now applied more often to bacterial populations.

Specificially, a new method which is alignment free is being deployed.

The studies based on mixed width K-mers as apposed to short nucleotide polymorphisms are being deployed, and packages such as *pyseer* are making it possible. However, there is not at present a succint tool which pipelines the necessary intermediry data processing and analysis steps. Let alone combines all these tests and allows for multiple association studies to be peformed at once, and distrubuted onto multiple nodes in a computer cluster.

which has utility in a time metabolomics data is becoming more common, and so the need to determine the genetic bases of possible hundreds of phenotypes as once is a requireid.

I researched, developed, and wrote and installable pipeline tool, *bgwas* which can test the association of multiple phenotypes at once to genetic loci. The tool integrates various open source software tools for Kmer counting, gene annotation, phylogeny estimation, kmer association testing, kmer-gene mapping and finally visualisation.

In creating the tool, best practices for computational biology were exercised; resulting in a final tool that expresses appropriate testing, documentation and packaging; meaning it can easily be utilised by others in the future.

The tool was also used to test the association of n phenotypes - some relating to antibiotic resistance and others metabolomics measures - and the outputs suggest the tool is usable, but similarly provides insights into future improvements that could be made

2 Acknowledgements

3 Abbreviations

SNPs

4 Introduction

4.1 Pangenome wide association studies

There is a recent wealth of bacterial genomes and phenotype data, and the repositories are rapidly expanding. In conjunction, there are increasingly more studies whose focus is on identifying the genetic basis of traits. (TODO ref) In bacteria, these traits may correspond with virulence, resistance, or some other quality that improves their ability as effective pathogens, and for this reason, discovering the genetic reasoning could lead to identifying targets for pharmaceutical intervention, and may increase the efficacy of treatment

Some studies have taken the approach of finding specific clonal groups as opposed to genetic loci. (ref) Genome wide association studies can be used to identify genetic and phenotype associations in a hypothesis free manner.

They have been used successfully in multiple human based studies.

Commonly, the units of variation that have been tested for association have been single nucleotide polymorphisms (SNPs).

SNPs have many advantages in organisms where recombination is reliable and genomes remain relatively constant among populations.

Because of linkage disequilibrium, identifying a highly associated SNP often means that the true associated loci is nearby

However, the nature of bacteria replication makes the use of SNPs in association studies less useful.

For one, significant regions of the genome can be in linkage disequilibrium: And though some species experience high rates of recombination, the recombination is not as reliable or consistent as that in, say, humans to reduce LD.

Secondly, bacteria experience much higher trans- during recombination. Therefore, bacteria of the same species can vary considerably in both order of gene content, and in the accessory genome - regions of the genome which differ.

Therefore, first finding SNPs becomes a tricky process which involves multiple alignments. Such tools such as mummer (and other?) do meet this problem, but even still, SNPs associated with only the core genome can be retrieved and does not consider the huge variable accessory genome in bacteria/ will not encapsulate the full

As an alternative, some recent studies have instead utilised K-mers.

Kmers are simply variable length sequences of DNA Initially a concept of genome assembly Can capture multiple genetic variations SNPs, longer deletions/ insertions and recombination The size of kmers affects the genetic variation captured Longer more specific Shorter are more sensitive

Getting kmers is alignment free, releasing the burden of full alignment But more specifically utilising kmers means the core and accessory genome can be tested.

Seer (see (ref) and its reimplementation (pyseer) dbpwa

4.2 The Need for Multiple Association Tests - Metabolomics

Metabolomics provides a detailed snapshot of an organism's physiological state through the quantification of hundreds to thousands of small molecules (Zampieri et al., 2017).

Coupling genomics and metabolomics to study *Pseudomonas aeruginosa* in unique clinical samples can allow us to understand how bacteria adapt to the lung environment.

The nature of metabolomics means that multiple molecules abundance may be measured and of interest to investigate.

Determine the genetic basis of metabolomics of *P* in the context of chronic infections

5 Scope of work

The foundation of this project was to create a tool, *bgwas3*, an integration of multiple open source genomics tools and custom scripts written in python, R and bash into a single installable package that facilitates conducting multiple pangenome wide association studies in a single user step.

The tool aims to be comprehensive and robust enough to convert only the most basic required input files into valuable and interpretable results including static and interactive visualisations.

When run on a node within a computer cluster, tasks that are computationally intensive or long are distributed among the cluster, and can run simultaneously.

Though the additional configuration file, the specifics of multiple intermediary steps can be altered, which can have variable effects on the final results.

In light of the necessity to perform multiple.. (1)

6 Methods

6.1 CGATcore

bgwas3 was written primarily in python, utilising the pipelining framework CGATCore, (Cribbs et al. 2019), an update to CGATRuffus (ref).

The framework comes as a dependency free Python package.

The ethos behind the framework is to wrap pipeline steps into discrete Python functions. Through the use of python decorators, these functions (referred to as tasks)

The tool's usefulness of CGATcore is twofold.

First, the tool can make use of/ incorporate a yaml configuration file

The tool determines the dependency of tasks, and completes them in order.

It checks for what tasks need to be run by files last modification date. Therefore if files are removed or added, only tasks that are out of date are re run.

Allows independent tasks to be run in parallel.

The recent addition (CGATCore) was written so as to take advantage of the benefits that computer clusters provide.

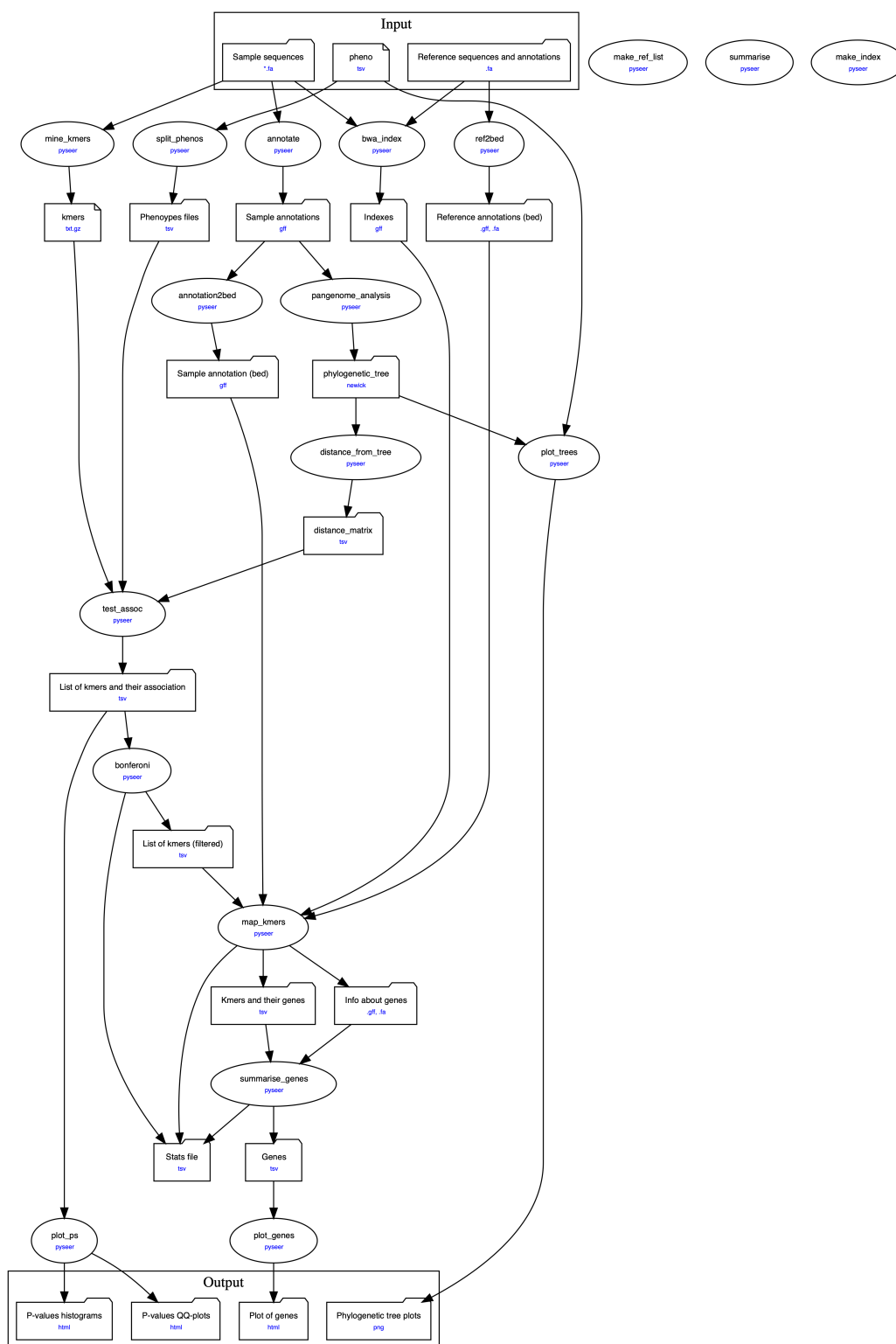


Figure 1: Pipeline of bgwas

It works alongside a Distributed Resource Management Application API (DRMAA) such as PBS-pri/Torque (used by imperial)

In its current iteration, *bgwas3* is comprised of n ‘tasks’, as visualised in figure. Starting with just three file types.

6.2 Genome Annotation

In the *bgwas* task “annotate” all sample sequences are annotated by Prokka (Seemann 2014). The tool functions as follows: It first makes predictions on features using a number of external tools including Prodigal [1] for the annotation of coding region, RNAmmer for ribosomal RNA genes, Aragorn for transfer RNA genes, SignalP for signal leader peptides and Infernal for non-coding RNA. After feature prediction, the speculative features are queried against a number of databases such as UniProt [2], Refseq [3] and Pfam.

The function of annotating input sequences in the *bgwas3* pipeline is twofold. First, gene annotations are later used to estimate the phylogenetic tree of the sample, and secondly, significantly associated Kmers are later mapped to the annotated genomes when identifying the Kmer’s genetic identity.

6.3 Kmer mining and counting

bgwas3 integrates the tool fsm-lite (Välimäki 2018) to count ‘mine’ and count Kmers. of user defined variable length kmers in the all samples. fsm-lite benefits first by being able to run on a single core, but also because of its ability to read and count Kmers not of a single length, but of a range of lengths, unlike similar tools such as DSK. As mentioned before, the length of Kmers affects the specificity and ? of the association testing. Therefore *bgwas3* allows the Kmers length to be changed by editing the configuration yml file ?

6.4 Phylogeny estimation and covariate estimation

bgwas3 currently takes a pangenomic approach to phylogeny estimation; as in it considers the relative presence and absence of genes in the samples. An estimated phylogenetic tree is estimated with *bgwas3* primarily utilising the tool Roary (Page et al. 2015). In summary, Roary determines genes which fall within the core genome, and then performs clustering of isolates are clustered based on gene presence in the accessory genome. A newick tree output of roary, is then converted into a distance matrix using a python script from pyseer.

The single reasoning for predicting phylogeny in the *bgwas3* pipeline is to use the distances defined in the distance matrix as covariates in the association testing.

This is different to the standard method of multi-dimensional scaling often used in human gwas and as part of the standard seer workflow. In general, if a phylogeny is accurate, th

6.5 Kmer association testing

For linear mixed model association testing, a *bgwas3* constructs a distance matrix based on a phylogenetic tree. Association between kmers and phenotypes is implemented with *pyseer* (Lees et al. 2018).

Specifically, *pyseer* is used to perform a linear mixed model using the Fast-LMM algorithm (Lippert et al. 2011).

LMM tackle confounders in association tests by using a measure of similarity (in this case a distance defined from an estimated phylogenetic tree) as a random effect in a linear model.

TODO equation

Given a means of determining the hierarchical relatedness between samples, the mixed model is generally preferred, and has been shown in past studies to control the inflation of p-values better (Lees et al. 2017).

6.6 Bonferroni correction

The output of *pyseer*, a list of all kmers and statistics relating to their association to the given phenotype, are then filtered by their p-value through Bonferroni correction. The Bonferroni correction is a multiple-comparison correction used when several dependent or independent statistical tests, utilised to limit the number of spurious positive tests. If multiple hypotheses are tested, the chance of a rare event increases, and therefore, the likelihood of incorrectly rejecting a null hypothesis (i.e., making a Type I error) increases. An R script **bonf.R** calculates a threshold p-value of which to accept an associated kmer as significant based on its p-value

6.7 Significant Kmer mapping

A new script, written in R **map_kmers.R** takes a list of kmers

fig>

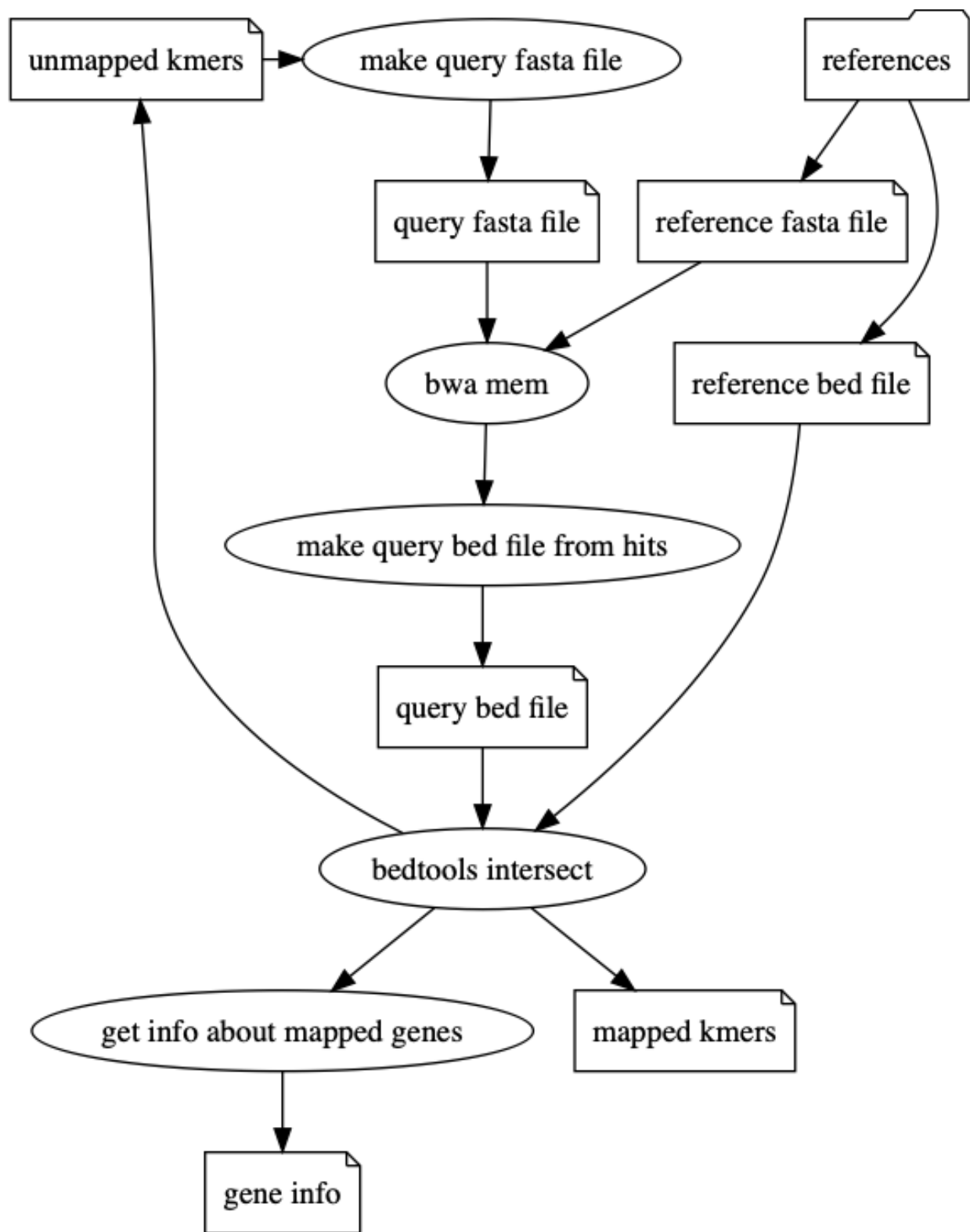


Figure 2: Graphical representation of algorithm to map Kmers to reference genes

6.8 Pathway analysis

6.9 Visualisation

6.10 Scientific computing practices

A significant goal of *bgwas3* was to implement a useful and reusable tool that may be used outside the scope of this singular project.

As such, various software development principles and concepts were applied to the project.

6.10.1 Testing with PyTest and Travis

Pytest and

6.10.2 Documentation with Sphinx

6.10.3 Packaging with Conda

7 Results

For the development, testing, and evaluation of *bgwas3* a dataset of

Cystic fibrosis is an autosomal recessive disorder that, due to a single gene mutation, limits the ability to *Pseudomonas aeruginosa* is one of the primary bacteria present in the lungs of late stage and terminal patients with cystic fibrosis. ref(Schaedel et al., 2002; Hauser, 2011) (becomes chronic). when the lung begins to exhibit decreased function and signs of failure For this reason, understanding the genetic adaptations *pseudomonas* experience when in chronic state is significantly important in treating this disease

This dataset involves *n* strains collected over a period of 24 years from 18 patients. ref

In paper () it was shown that there was ‘clear evidence’ of the metabolic lung environment surrounding the bacteria. Specifically acetate production was negatively associated with length of infection and others?

To test the usability and application of *bgwas3*,

The strong LD caused by the clonal reproduction of bacterial populations means that non-causal k-mers may also appear to be associated.?

Confirmation of known resistance

Name	Description
Tobromycin	Resistance to inhalant antibiotic Tobromycin
Imipenem	Resistance to intravenous antibiotic Imipenem
Aztreonam	Resistance to intravenous/intramuscular antibiotic Aztreonam
Ciprofloxacin	Resistance to oral antibiotic Ciprofloxacin
Colistin	Resistance to ‘last-resort’ antibiotic Colistin
Swim	Measure of cell surface bacteria movement by flagella
Swarm	Measure of rapid surface movement by multiple bacteria with rotating flagella
Twitch	Measure of slow bacteria movement powered by pili

Chemical

Hydrogen Cyanide

Cyanide

2-Furoate

3-Hydroxyisovalerate

3-Methylthiopropionic acid

Anthranilate

Betaine

Cystine

Formate

Fumarate

Histidine

Isoleucine

Leucine

Methanol

Methionine

Tryptophan

Uracil

Valine

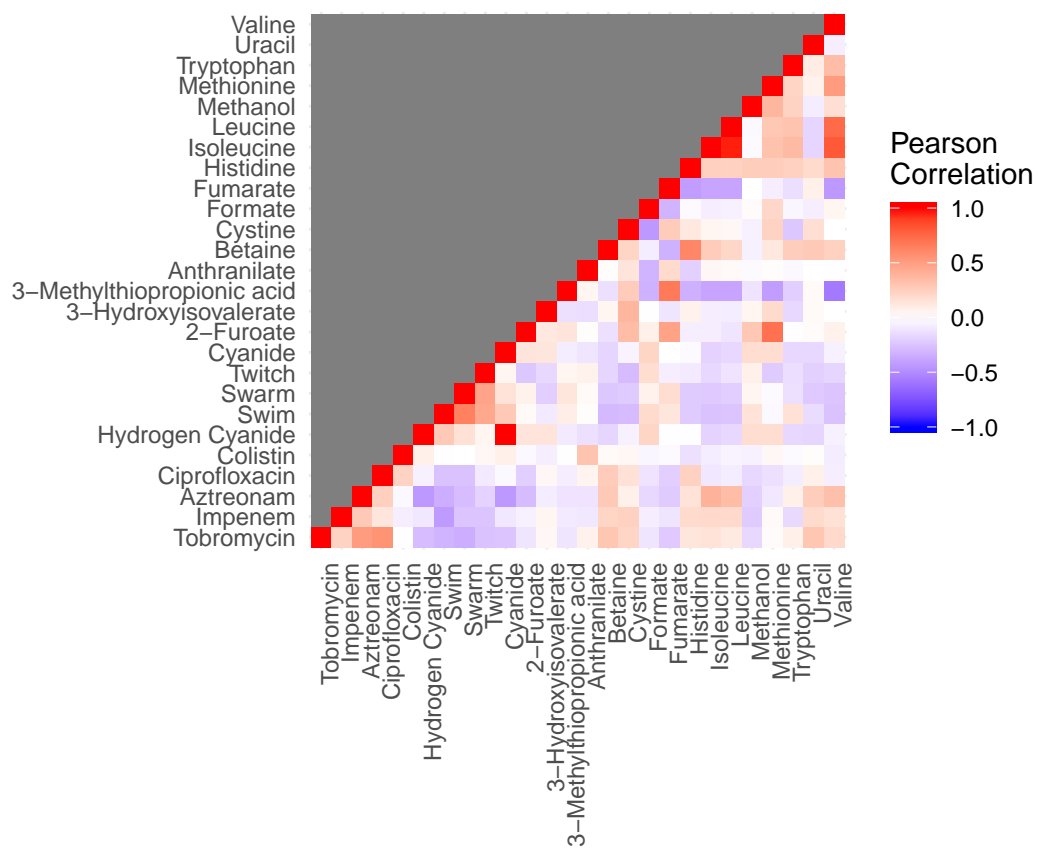


Figure 3: Correlation matrix of phenotypes

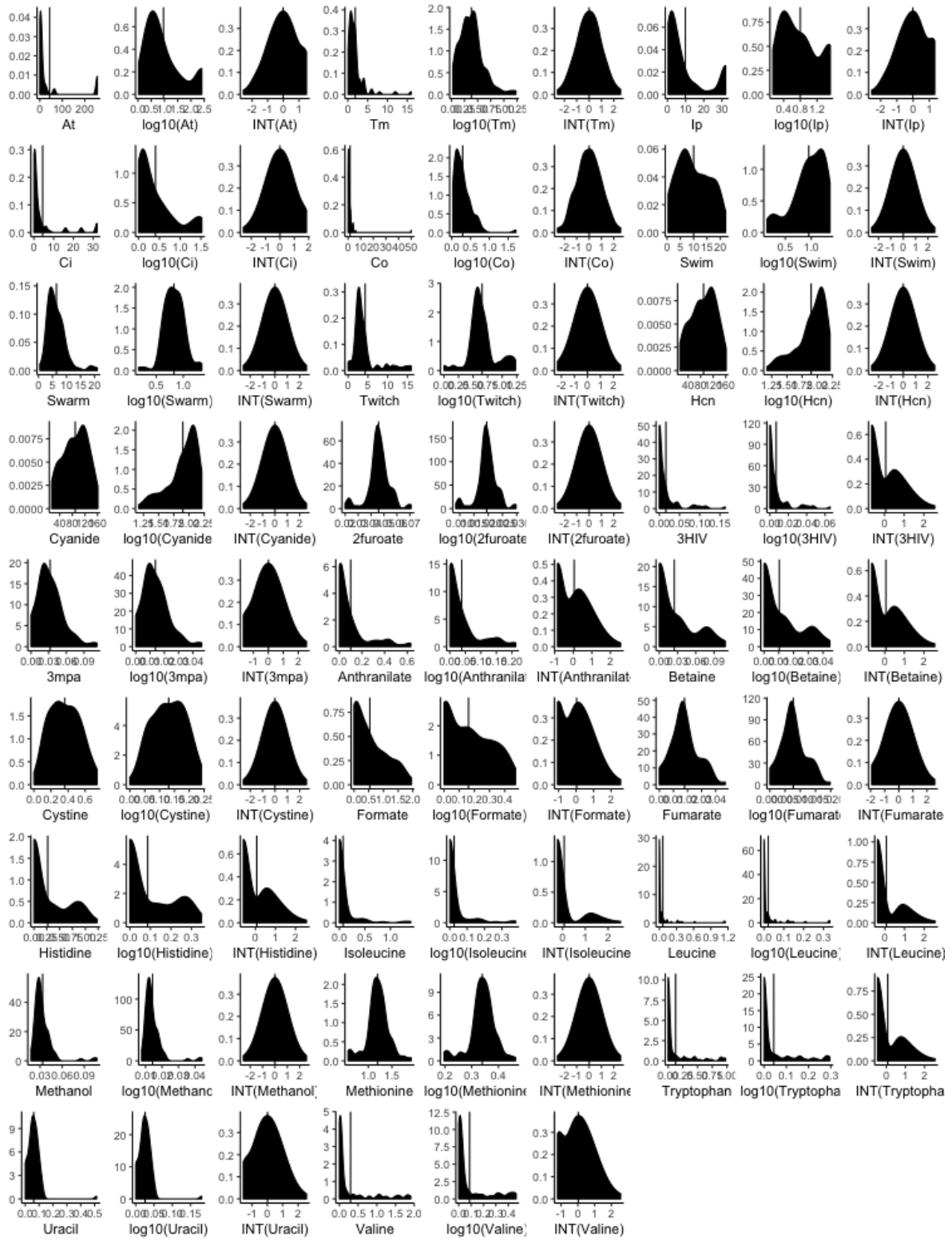


Figure 4: Density plots

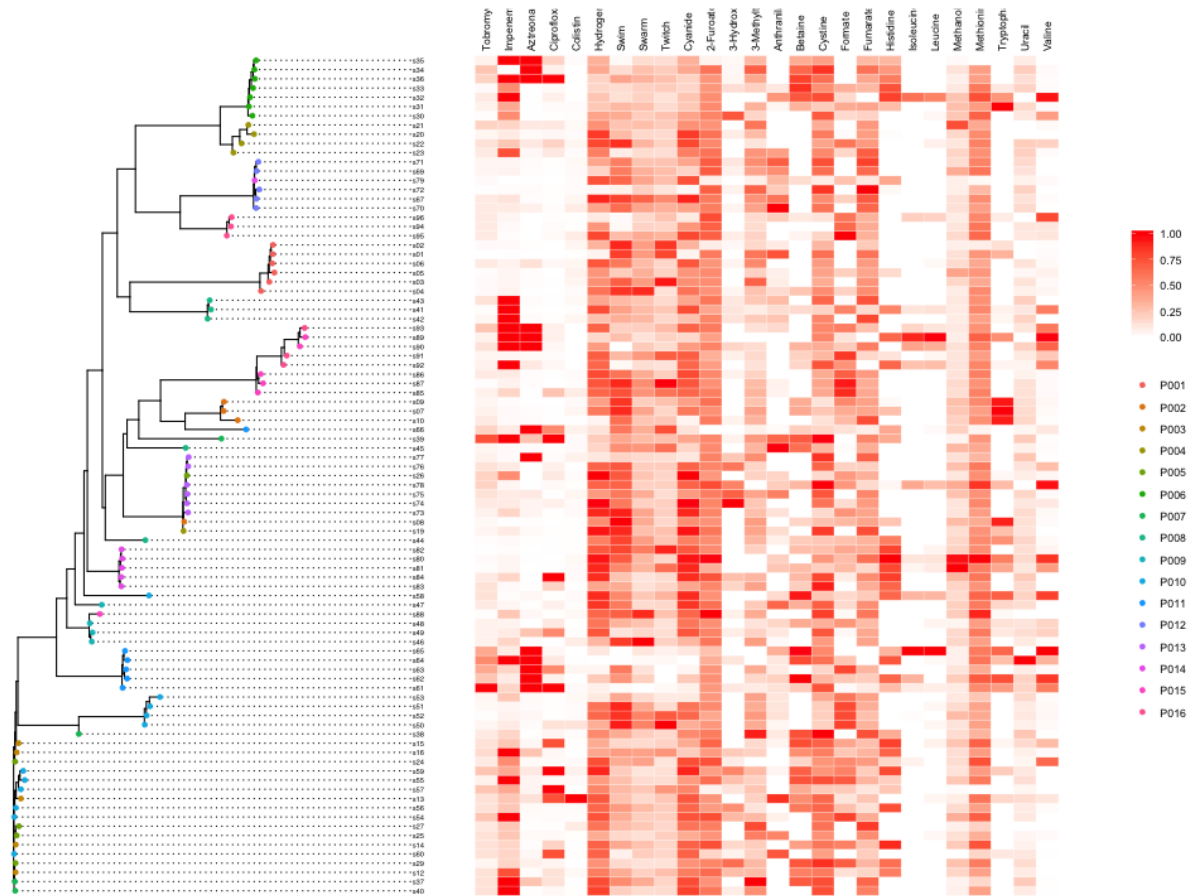


Figure 5: Density plots

7.1 Processing of phenotypes

Prior to running *bgwas*, the *n* phenotypes were transformed in two ways. Log transformed and INT INT is becoming a common method of normalising data The problem of normalising data

fig density plots

discuss achievement on normalising

7.2 Annotation

All 91 genomes were annotated and 14643 unique genes were identified using Prokka.

7.3 Kmer mining

Number etc

7.4 Phenotype

8 Discussion

8.1 Genome annotation

8.2 Kmer mining

The tool fsm

Prokka is a good tool ?Settings ## Phylogeny prediction Currently, *bgwas3* implements only a pangenomic approach approach of distance estimation. There are other tools which involve alignment of the core genome and snps... May or may not be better Reintroduce the problem of a large multiple alignment.

References

Cribbs, Adam P., Sebastian Luna-Valero, Charlotte George, Ian M. Sudbery, Antonio J. Berlanga-Taylor, Stephen N. Sansom, Tom Smith, et al. 2019. “CGAT-Core: A Python Framework for Building Scalable,

Reproducible Computational Biology Workflows.” *F1000Research* 8 (April): 377. <https://doi.org/10.12688/f1000research.18674.1>.

Lees, John A, Nicholas J Croucher, David Goldblatt, François Nosten, Julian Parkhill, Claudia Turner, Paul Turner, and Stephen D Bentley. 2017. “Genome-Wide Identification of Lineage and Locus Specific Variation Associated with Pneumococcal Carriage Duration.” Edited by Sarah Cobey. *eLife* 6 (July): e26255. <https://doi.org/10.7554/eLife.26255>.

Lees, John A, Marco Galardini, Stephen D Bentley, Jeffrey N Weiser, and Jukka Corander. 2018. “Py-seer: A Comprehensive Tool for Microbial Pangenome-Wide Association Studies.” Edited by Oliver Stegle. *Bioinformatics* 34 (24): 4310–2. <https://doi.org/10.1093/bioinformatics/bty539>.

Lippert, Christoph, Jennifer Listgarten, Ying Liu, Carl M Kadie, Robert I Davidson, and David Heckerman. 2011. “FaST Linear Mixed Models for Genome-Wide Association Studies.” *Nature Methods* 8 (10): 833–35. <https://doi.org/10.1038/nmeth.1681>.

Page, Andrew J., Carla A. Cummins, Martin Hunt, Vanessa K. Wong, Sandra Reuter, Matthew T. G. Holden, Maria Fookes, Daniel Falush, Jacqueline A. Keane, and Julian Parkhill. 2015. “Roary: Rapid Large-Scale Prokaryote Pan Genome Analysis.” *Bioinformatics* 31 (22): 3691–3. <https://doi.org/10.1093/bioinformatics/btv421>.

Seemann, Torsten. 2014. “Prokka: Rapid Prokaryotic Genome Annotation.” *Bioinformatics* 30 (14): 2068–9. <https://doi.org/10.1093/bioinformatics/btu153>.

Välimäki, Niko. 2018. “Fsm-Lite.” <https://github.com/nvalimak/fsm-lite>.