

## Introduction and Purpose

The purpose of this project was to develop pedigree graphs for use in pedigree-linkage analysis projects. More specifically, this project sought to:

- find a source of pedigree and quantitative trait loci (QTL) data that is suitable (i.e.: large enough) for building arbitrary pedigrees
- use the aforementioned data to build bayesian network representations of arbitrary pedigrees
  - for this project, segregation networks were used (TODO)
- convert arbitrary bayesian network pedigrees into a standardized format (UAI)
- extract meaningful subsets of variables from the aforementioned bayesian network pedigrees for use in queries

## Generating Pedigree/QTL Data with QMSim

Because of an ostensible scarcity of reliable pedigree/QTL data, I used a pedigree/QTL data simulator (QMSim) to generate all the data necessary for this project.

“QTL and Marker Simulator”, or QMSim, is lightweight software that generates pedigree, quantitative trait loci, and genetic marker data for simulated livestock populations. QMSim was developed at the University of Guelph. To read more about QMSim, see Sargolzaei’s and Schenkel’s paper “QMSim: A Large Scale Genome Simulator for Live Stock.”

QMSim is a fairly complex piece of software with myriad uses, input parameters, and output parameters, and so can be somewhat daunting to use. Because of the small scope of this project, however, there is a fairly small subset of usage parameters that are useful for this project. Below, I will describe how to use QMSim in the context of this project, which parameters should be considered in the context of this project, and the effects each parameters has on networks built from data generated by QMSim.

## How to Use

Although QMSim comes with a large user manual that contains an exhaustive list of parameter descriptions, many of the use cases and parameters described in the user manual are not useful in the context of this project. The full user manual can be found at: [http://www.aps.uoguelph.ca/~msargol/qmsim/QMSim\\_documentation.pdf](http://www.aps.uoguelph.ca/~msargol/qmsim/QMSim_documentation.pdf).

## Parameters and Parameter Files

### Parameter Files

QMSim takes as input parameter files which have the file extension “.prm”. There are several mandatory parameters that all QMSim parameter files must contain, but to avoid repeating what is already present in the QMSim instruction manual, I’ll only discuss the parameters that were of interest to me in generating data for pedigrees. There are several valid **prm** files in the **QMSim/prm** directory that came with this report, most of which were provided as samples with QMSim and modified as necessary. When generating new pedigrees, I recommend using any of these files as templates and modifying them as necessary.

### Parameters

#### Historical Population Parameters

Historical population parameters affect the simulation of the populations that precede the population for which data is generated by QMSim. In other words, if QMSim generates data for generations 0-10, historical population parameters affect the simulation process for all generations  $< 0$ . It’s especially important to pay attention to population data when generating small pedigrees: if the historical population isn’t sufficiently large, QMSim will throw QTL variance errors.

Historical population parameters are specified between **hp** tags:

```
begin_hp;  
    ...  
end_hp;
```

**hg\_size**

The **hg\_size** is an optional parameter that affects the number of individuals in each historical generation, as well as the number of historical generations. To use:

```
hg_size = v1 [v2] ... v3 [v4];
```

Where **v1** is the number of individuals in generation **v2** (which should always start at 0, and **v3** is the number of individuals in generation **v4**. Arbitrary size/generation values can be placed after between the values for the first and last generations.

Although these parameters don’t have a direct effect on networks generated from this data, they are nonetheless very important for correctly generating

diverse populations with sufficiently diverse QTLs. The QMSim manual states that legal values for the populations size of the historical generations (**v1**) can be as small as 2, and that the number of historical generations (**v4**) can be as small as 1, these values, in my experience, must be significantly higher to produce populations with desirable genetic diversity and with sufficiently balanced gender distributions. Anecdotally, population sizes should be  $> 10$  and the number of generations should be  $> 20$ . Of course, these number all correlate, so some experimentation may be required depending on the situation. For all of my experiments, I used 200 generations with population sizes of 500. In general, more and larger historical generations produces more genetically diverse reported populations (depending on a few other paramters; see genome parameters section).

### Population Parameters

Population parameters affect the simulation of the populations that QMSim actually reports on.

Population parameters are specified between **pop** tags:

```
begin_pop = "str";
...
end_pop;
```

where **str** is a string name for the population.

**male/female**

The **male** and **female** founder parameters control various aspects of the male and female founder populations, and are primarily used, along with **ls** and **ng**, to control the size of the generated pedigree. Although there are several ways in which QMSim allows its users to configure the male/female founder populations, the most “interesting” use for these parameters in the context of this project is controlling the male/female ratio in the founding generation, which in turn affects the sizes of subsequent generations, and therefore the number of variables in a segregation network build from generated data.

For example, a founder population of 10 males and 10 females produced a segregation network with 840 variables. A founder population consisting of 2 males and 18 females produced a segregation network with 1480 variables, albeit with more inbreeding. I did find, also, that a larger male/female ration increases the induced width of the corresponding segregation graph. For example, using a  $\frac{1}{5}$  male/female ( $\frac{5}{25}$ ) generated a graph with an induced width of  $\approx 85$  (calculated using an average of 5 different networks), whereas using a ration of  $\frac{1}{1}$  ( $\frac{15}{15}$ ) resulted in networks with induced widths of  $\approx 110$  (also calculated using an average of 5 different networks). To use:

```

begin_founder
    male  [n = int, pop = "str"];
    female[n = int, pop = "str"];
end_founder

```

where **int** is an integer in the range 1-50000 and **str** is the string name of the population for which data is being generated. The male/female ration can be at most 1.

**ng**

This parameter controls the number of generations in the population. To use:

```
ng = int;
```

where **int** is an int in the range 0-2000. **ng**, as well as **ls** and **male/female**, is used to control the size of the generated pedigree.

Based on my observations, **ng** has no measurable effect on network complexity outside of affecting changes in network size. A network created using 20 male founder and 20 female founders, with an **ng** of 10, contained 840 variables , had a pseudotree depth of 80, and an induced width of 54. Doubling the number of generations to 20 almost doubles the number of variables and the pseudotree depth (to 1640 and 144, respectively), but only increased the width to 63. Doubling **ng** once more to 40 about doubled the number of nodes and tree depth once again (to 3240 and 277, respectively), but again caused a modes increase in tree width to from 63 to 72.

**ls**

This parameters controls the maximum number of offspring produced by each female. **ls**, along with **ng** and **male/female** are used to control the size of the generated pedigree.

## Population Output Parameters

Population output paramters allow users to specify the data that QMSim outputs.

Population output parameters are specified between **popoutput** tags, and are also nested between **pop** tags:

```

begin_pop;
    ...
    begin_popoutput;
        ...
    end_popoutput;
end_pop;

```

### **data**

Specifying the **data** output parameter causes QMSim to output family/pedigree data for the specified generations. To use:

```
data \gen0 ... \gen_n;
```

where **\gen\_n** refers to an arbitrary generation. If no generation is specified, data for all generations will be reported. See the ‘Output’ section for a more detailed description of the format of files generated by the **data** tag.

Note: the files generated by **data** are necessary to build segregations networks.

### **allele\_freq**

Specifying the **allele\_freq** output parameters causes QMSim to output allele frequency statistics. To use:

```
allele_freq \gen0 ... \gen_n
```

where **\gen\_n** refers to an arbitrary generation. If no generation is specified, data for all generations will be reported. See the ‘Output’ section for a more detailed description of the format of files generated by the **allele\_freq** tag.

Note: the files generated by **allele\_freq** are necessary to build segregations networks.

### **genotype**

Specifying the **genotype** output parameters causes QMSim to output allele assignments for individuals in the pedigree. To use:

```
genotype \gen0 ... \gen_n
```

where **\gen\_n** refers to an arbitrary generation. If no generation is specified, data for all generations will be reported. See the ‘Output’ section for a more detailed description of the format of files generated by the **genotype** tag.

Note: the files generated by the **genotype** are necessary to build segregations networks.

## **Genome Parameters**

Genome parameters control the way QMSim simulates genotype/QTL/marker data for the generated pedigrees. Genome parameters are specified between **genome** tags. Additionally, all genome parameters relevant to this project are specified between **chr** (chromosome) tags:

```
begin_genome;
  begin_chr = int;
  ...
  end_chr;
end_genome;
```

where `int` is a integer that specifies the number of chromosomes for which to generate data.

QMSim produces both genetic marker and QTL data, but this project uses only QTL data.

`nqloci`

This parameter affects the number of quantitative trait loci (QTL) on each chromosome. To use:

```
nqloci = int;
```

where `int` is a integer in the range 0-50,000. If `nqloci` is 10, and the number of chromosomes being simulated is 20, then QMSim will generate data for  $10 \cdot 20 = 200$  QTL.

`nqloci` affects the complexity of segregation network by allowing users to incorporate more alleles into a network. For examples, creating a pedigree using `male = female = 10`, `ls = 10`, and `ng = 10` produces a network with 840 variables if only one QTL is used to build the network. Adding another QTL to the network increases its size of the network to 1680, and adding another QTL to the network increases its size to 2520. It is easy to see that linearly increasing the number of QTLs incorporated into a network results in a linear increase in that network's size (measured in number of variables). It is not quite so easy to describe the increase in network complexity that occurs as a result of the incorporation of more QTLs: the increase in network complexity caused by adding more QTL is largely dependent on the number of allelic types that occur at that QTL and the underlying distribution of those types. See TODO for more information.

`qpos`

This parameter affects the positions of QTLs on their respective chromosomes. This parameter has several potential values. See the QMSim manual for a full list of possible values. To use:

```
qpos = even;
```

Causes QTLs to be spaced evenly along their respective chromosomes.

```
qpos = rnd;
```

Causes QTLs to be spaced randomly along their respective chromosomes.

`qpos = pd float_1 float_2 ... float_n`

Causes QTLs to be positioned at predefined locations specified by the float values `float_1` to `float_n`, where all values  $v$  are  $0 \leq v \leq 100$ . If the `pd` option is used, the number of float position values specified must be the same as the value of `nqloci`.

`qpos` affects the complexity of networks in a manner that is difficult to measure: QTL positions are used to build factor tables for cliques in segregation graphs that consists of multiple segregation nodes. See TODO for more information.

`nqa`