# Bioinformatics workshop

# 3$^{rd}$ SoyaGen Annual Meeting

# 7 December 2018

# Genetic Mapping workshop

Manel Fallah, François Belzile
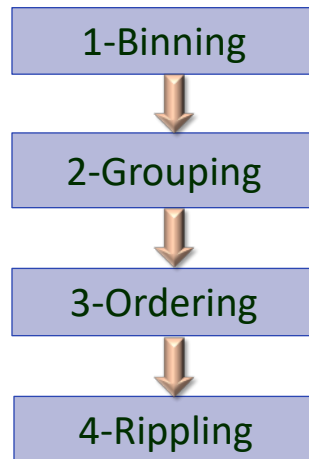
3rd Annual SoyaGen meeting

07-12-2018

---

## The steps for genetic mapping:

A. Optimization of the data :

1. Filtering to remove excess heterozygosity
   a. per line
   b. per marker

2. Minor allele frequency (expect 0.5:0.5 for two alleles)

3. Conversion script (nucleotides to numerical genotype)

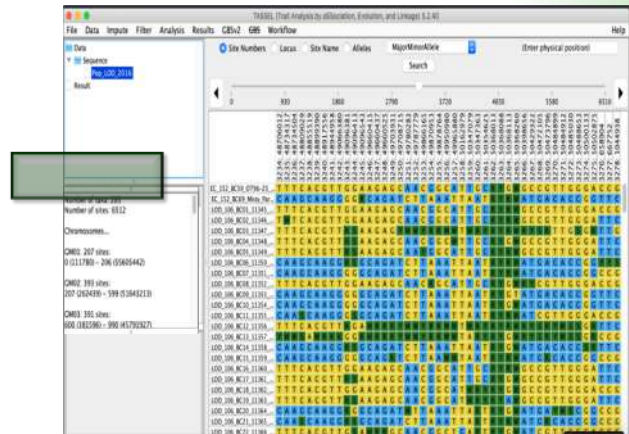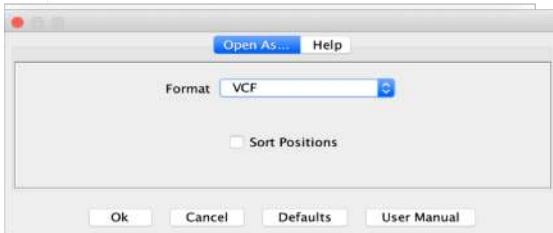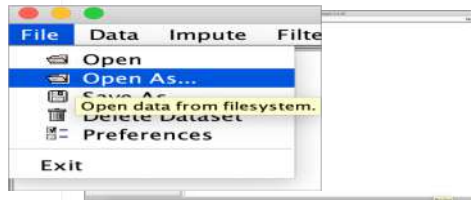## B. Construction of a genetic linkage map

1-Binning

⬇

2-Grouping

⬇

3-Ordering

⬇

4-Rippling

---

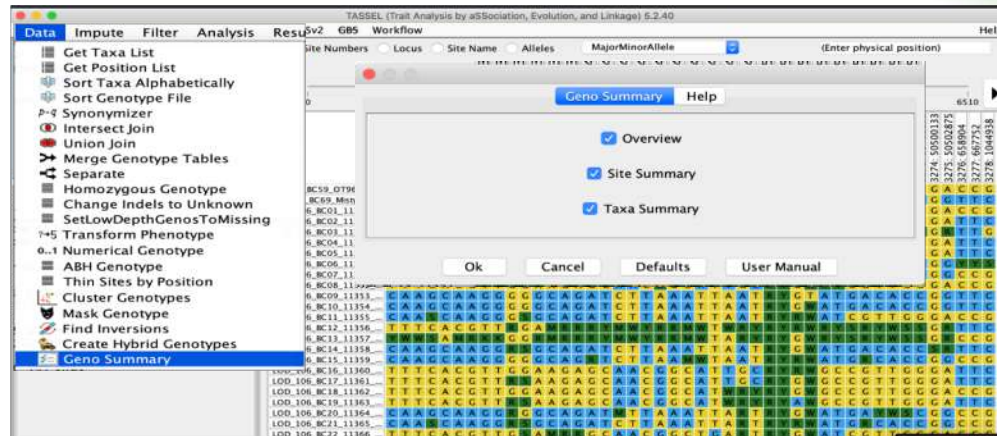A-Optimization of the data

# The heterozygosity filter

➢ The point of heterozygosity filters is to eliminate lines and markers that have an excessive amount of heterozygosity.
➢ There will always be some heterozygosity (biological and technical causes)
➢ The amount of "normal" heterozygosity depends on the type of lines (e.g. F5 vs RILs)

Apply a Genosummary

Calculate the threshold for heterozygosity
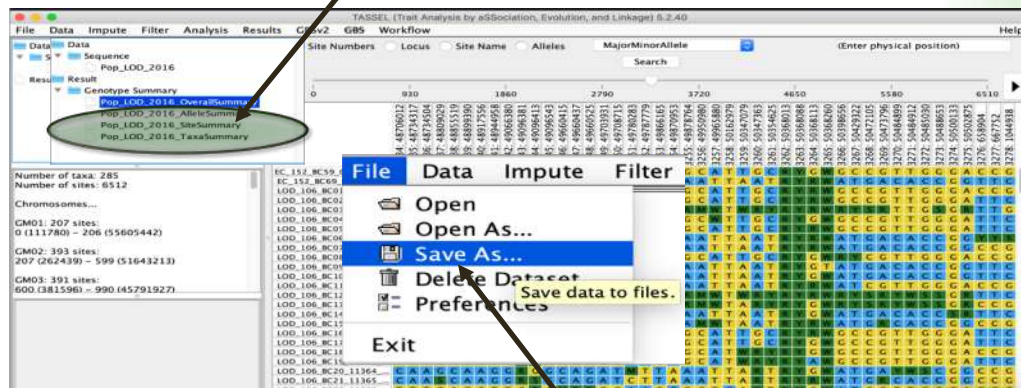
Filter the data

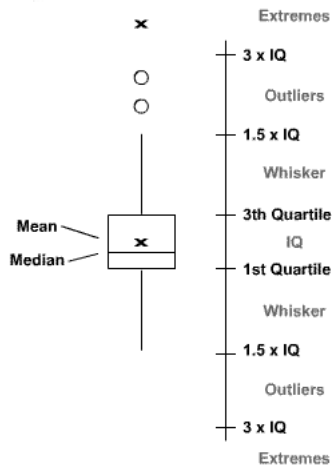The steps of the Genosummary: **Data** then **Geno Summary**

## The Genosummary results

The two Result files to export to Excel to determine an appropriate threshold (i.e. how much heterozygosity is "too much")
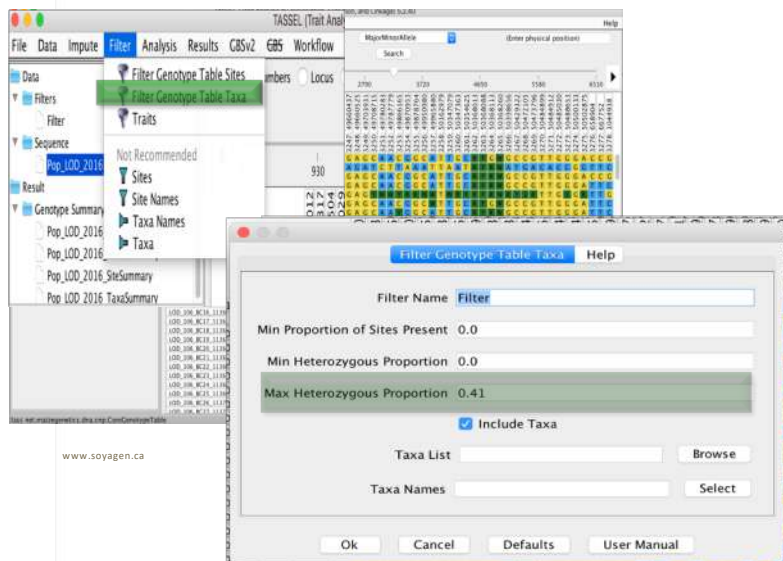
Save as "Table"

## Extreme heterozygosity values are determined using an inter-quartile-range (IQR) approach



- The distribution of heterozygosity (found both in lines and in markers) is plotted in the form of a box plot
- Extreme outliers are lines/markers whose proportion of heterozygosity exceeds a value corresponding to >1.5 x IQR
- The critical value used to filter the dataset (both for markers and lines) is that which corresponds to 1.5 x IQR
- This analysis is carried out in Excel

www.soyagen.ca

9

---

## a-Filtering the lines that have "too much" heterozygosity



✓ Verify if the filter was correctly done

www.soyagen.ca

10

b-Filtering the markers that have "too much" heterozygosity


2-Filtering markers for Minor Allele Frequency
(expectation = 0.5; tolerate down to 0.3)

## Heterozygosity distribution before and after filtration



The excessive heterozygosity is eliminated

---

## 3-Conversion script

- SNP catalogues provide the actual nucleotide present at a given locus
- Genetic mapping softwares:
  - do not care about this information
  - want to know which allele came from which parent
- We need to convert each genotype into a number
  - For example:
    - AA genotype present in parent 1 is converted to a 0
    - GG genotype present in parent 2 is converted to a 2
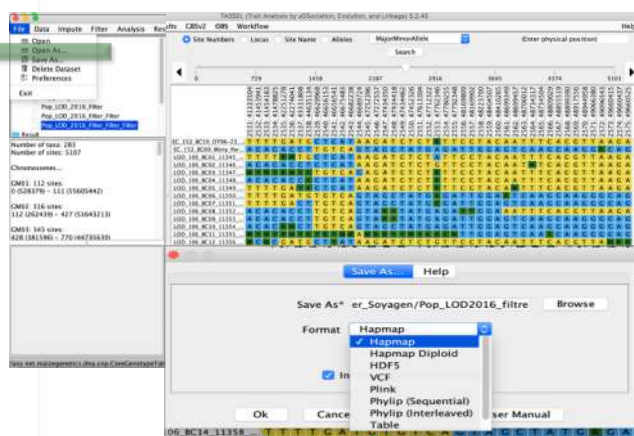    - AG genotype (present in a heterozygote) is converted to 1

## Prepare the file for the conversion

1-Open the hapmap file (.hmp.txt) in Excel



2-Copy and transpose the rs# column in a new Excel sheet, delete the rs# cell and save as a text (.txt) file

| SNP1 | SNP2 | SNP3 | SNP4 | SNP5 | SNP6 | SNP7 | SNP8 | SNP9 |
|------|------|------|------|------|------|------|------|------|
| Parent1 | C | C | C | T | C | T | C | A |
| Parent2 | T | T | T | A | G | C | G | G |
| Indiv1 | T | T | T | A | G | C | G | G |
| Indiv2 | T | T | T | A | G | C | G | G |
| Indiv3 | C | C | C | T | C | T | C | A |

(Parental lines **must** be the first two individuals in this file)

17

---

## Run the script conversion

Run the bash command line :   ./ped2carto.py Pop_atelier.txt

Script name          File Name

⚠ There must be a space between these



An output file with a new extension (_carto.txt) is going to be produced
(in the case illustrated, it will be called "Pop_atelier_carto.txt")

| locus_name | Indiv1 | Indiv2 | Indiv3 | Indiv4 | Indiv5 | Indiv6 | Indiv7 |
|------------|--------|--------|--------|--------|--------|--------|--------|
| SNP1 | 0 | 0 | 2 | 0 | 2 | 0 | |
| SNP2 | 0 | 0 | 2 | 0 | 2 | 0 | |
| SNP3 | 0 | 0 | 2 | 0 | 2 | 0 | |
| SNP4 | 0 | 0 | 2 | 0 | 2 | 0 | |
| SNP5 | 0 | 0 | 2 | 0 | 2 | 0 | |
| SNP6 | 0 | 0 | 2 | 0 | 1 | 0 | |

After running the script (rows = SNP loci, parents are eliminated and all genotypes are now 0, 1 or 2)

18

9

# Elimination of "double recombinants"

When long stretches of one parental allele are interrupted by one or a few alternate genotype calls, these are most often genotyping errors (called "double recombinants"

These must be eliminated and there is a tool for that!

Suspected errors are replaced with missing data (coded as "-1")

---

# B. Construction of a genetic linkage map

1- Binning: Redundant markers are markers that are identical and do not provide additional information. The point is to remove redundant markers prior to generating the map file.

2-Grouping: Linkage groups are assembled based on anchor info (physical position).

3-Ordering: After the groups are correctly formed, the ordering is done "by input".

4-Rippling: After ordering, the position of each marker needs to be "fine tuned"

5-Outputting : Output the Results file.

# 1-Binning

**a. Preparation of the bin file**

Population type →
Unit (cM) →
Mapping function →
Numbers of lines →
Numbers of markers →

| | |
|---|---|
| 8 | |
| 1 | |
| 2 | |
| 281 | |
| 6512 | |

⚠ No Lines list

This is done by:

1) Copying the list of loci (column A)

2) Pasting the additional copy of loci at the bottom of the genotype table

3) Indicating in column B the appropriate chromosome number

⚠ 3-Convert from (.txt) to (.bin)

---

# QTLIciMapping
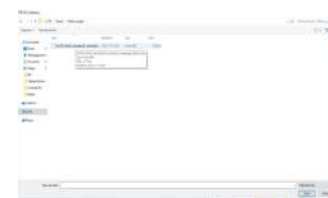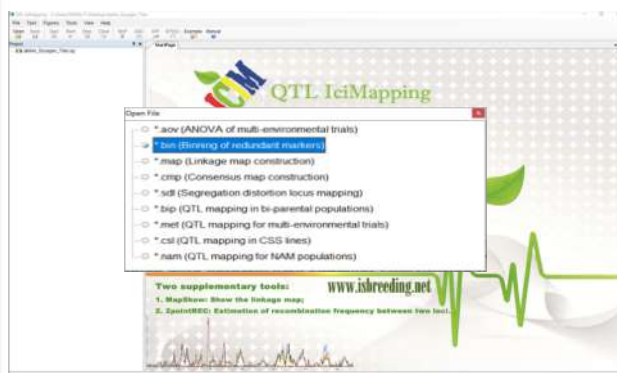
# 1-Create a new project

1-file
2-create a new project

✓ The project created

23



# 2-Load the bin file

✓Open then choose the file type (here we choose bin)
✓Go to the working directory to choose the file
✓Select "Open" or "ouvrir"

24

12

# 3- Binning

✓The file is loaded

**Project**
- Atelier_Soagen.ipj
- BIN
  - Pop_Atelier_carto_binfile.bin

| StartPage | Pop_Atelier_carto_binfile.bin | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| **Marker Information** | | | | | | | | |
| MarkerID | MarkerName | AnchorInfo | Missing(%) | ChiSquare | DistortP | BinID | Deleted |
| 1 | SNP1 | 1 | 0 | | | | |
| 2 | SNP2 | 1 | 1,0601 | | | | |
| 3 | SNP3 | 1 | 1,7668 | | | | |

**Parameters**

| Delete markers | Anchor Information | Missing values | Delete redundancy |
|---|---|---|---|
| By missing rate (%) 10.00 | ☐ Consider anchor info | ☐ Consider missing values | ◉ By Missing Rate (%) |
| By distortion P value 0,0000 | | | ○ By Random |
| Non-polymorphism, and markers with higher missing rate or lower P value will be deleted. For those markers, BinID = -1. | If selected, redundant markers in same anchor group will be assigned to one BIN group.If not, redundancy is the only factor considered in BINNING. | If selected, missing values are used, resulting in two markers at same position in map construction. If not selected, non-redundant markers may be in one bin. | For non-redundancy, BinID = 0. One is retained in each bin. |
| | | | Binning |

---

# The result file for the binning

✓ Two columns that are important: "**BinID**" and "**Deleted**"
✓ **BinID**:
    * 0 : the marker is unique and it doesn't belong to any linkage group
    * Every number >0 indicates a group identification or BinID (markers with the same BinID belong to the same linkage group
✓ **Deleted**: 0 the marker was not deleted (saved for the map file)
          1 the marker has been deleted

| StartPage | Pop_Atelier_carto_binfile.bin | | | | | | |
|---|---|---|---|---|---|---|---|
| **Marker Information** | | | | | | | |
| MarkerID | MarkerName | AnchorInfo | Missing(%) | ChiSquare | DistortP | BinID | Deleted |
| 1 | SNP1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 2 | SNP2 | 1 | 1,0601 | 0 | 0 | 1 | 0 |
| 3 | SNP3 | 1 | 1,7668 | 0 | 0 | 1 | 1 |
| 4 | SNP4 | 1 | 2,8269 | 0 | 0 | 1 | 1 |
| 5 | SNP5 | 1 | 1,4134 | 0 | 0 | 2 | 1 |
| 6 | SNP6 | 1 | 1,0601 | 0 | 0 | 2 | 1 |
| 7 | SNP7 | 1 | 0,3534 | 0 | 0 | 2 | 0 |

**4-Load the Map file**

✓ The same steps to load the bin file but we choose map



**5-The steps for mapping**



1-Create the linkage groups

2-Order the markers in the linkage groups

3-Optimize the marker order

4-Export the results

# 6-The Output files

# 7-Draw the linkage map