



ISMU 2.0: A Multi-Algorithm Pipeline for Genomic Selection

5th International Conference on
Next Generation Genomics and Integrated Breeding
for Crop Improvement

Wednesday, February 18, 2015

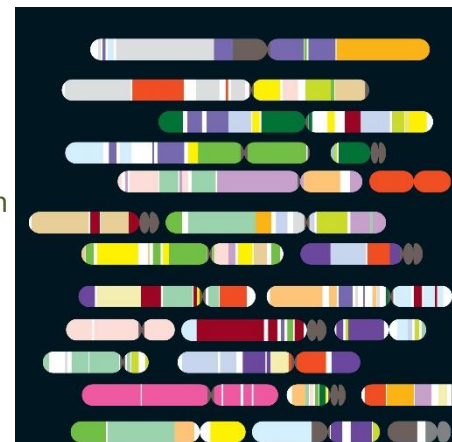
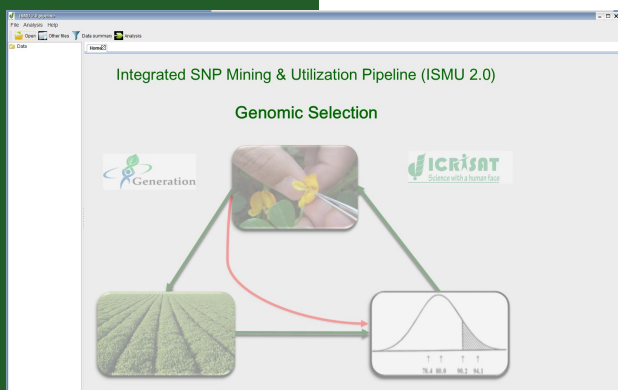
Abhishek Rathore¹, Roma R. Das¹, Manish Roorkiwal¹,
Dadakhalar Doddamani¹, Mohan Telluri¹, David
Edwards², Mark E Sorrells³, Janez Jenko⁴, John Hickey⁴,
Jean-Luc Jannink³ and Rajeev K. Varshney¹

¹ ICRISAT, Hyderabad, India

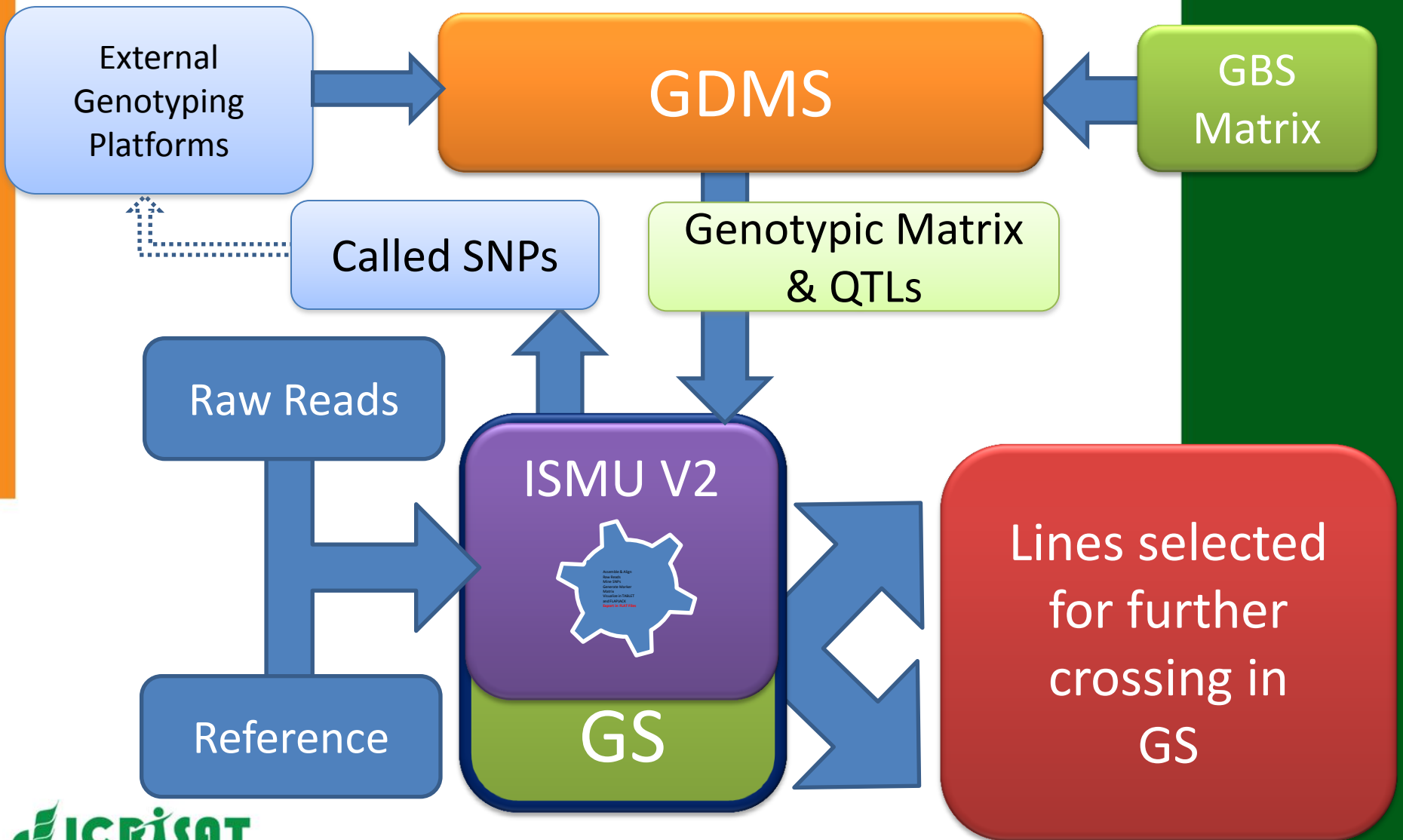
² University of Queensland, Brisbane, Australia

³ Cornell University, Ithaca, NY

⁴ The University of Edinburgh, Scotland, United Kingdom



ISMU V2.0





Genomic Selection (GS)

Genomic tool to accelerate breeding cycle

- Increases **genetic gain** per cycle through early selection
- Very useful for **complex traits** (Difficult/expensive/takes long time to phenotype, etc.)
- Breeding values are predicted on the basis of genome wide markers, called Genomic Estimated Breeding Values (**GEBVs**)
- Several analytical approaches / GS models have been proposed for prediction of GEBVs





GS Approaches / Models?

- To meet the challenges, statistical methods that can handle high-dimensional data developed
- Respective properties are still not fully understood
- Causing considerable uncertainty about the choice of models for genomic prediction
- Factors affecting GS are also not very clear





Factors Affecting GS-Models?

- Marker density, genome size and structure?
- Size of the training population?
- Historical effective population size?
- Trait heritability?
- Relationship between training population & selection candidates?
- Number of genes and distribution of their effects?
- Method used for the estimation of marker effects?
- GxE?

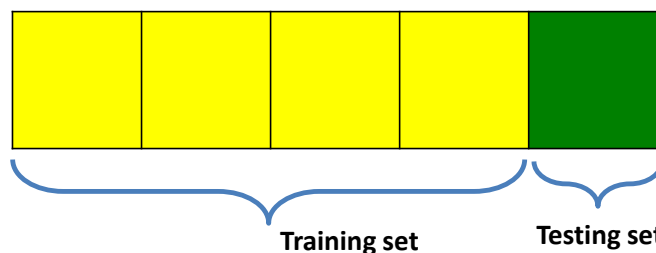




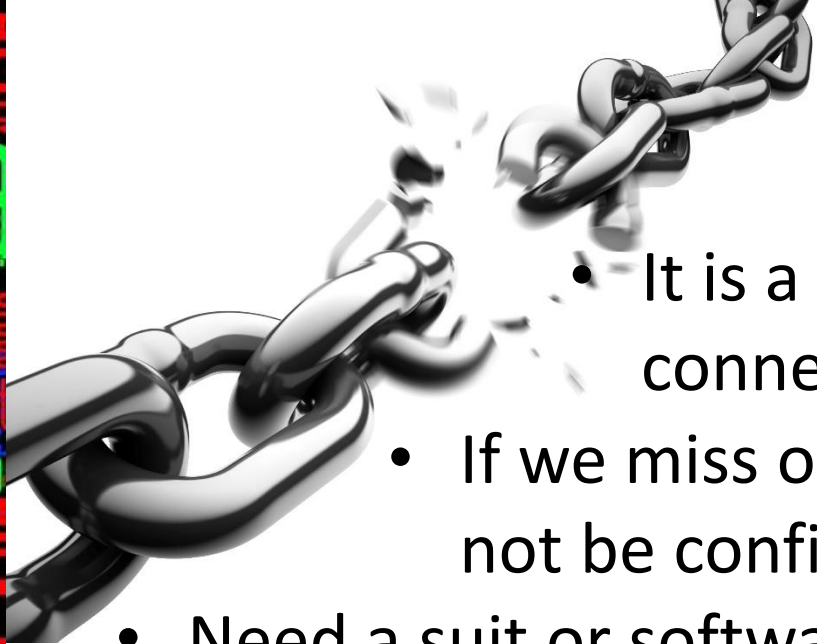
Many Steps in Genomic Selection...

- ✓ Get Training Population (Marker & Phenotype)
- ✓ Quality control / data filtering
- ✓ Model Population Structure / Covariates
- ✓ Fit available models
- ✓ Perform Cross Validation
- ✓ Prepare matrix of scores
- ✓ Select final method
- ✓ Get Testing Population, Predict GEBVs
- ✓ Make Selection based on GEBVs
- ✓ Add new data & rebuild model

Cross Validation $K(=5)$ - fold cross-validation



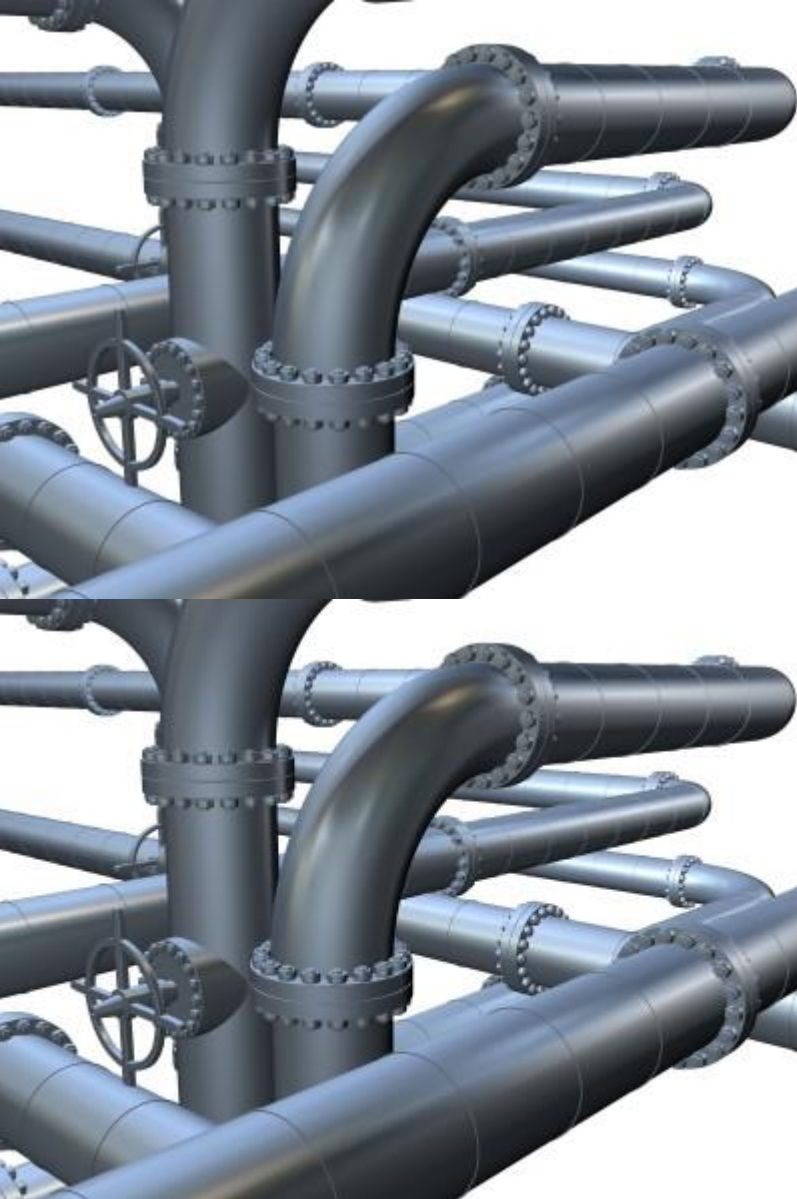
Genomic
Popular
Final matrix
Validation matrix
Marker Prepare
Perform
model data
repeat
Genomic
Final matrix
Validation matrix
Marker Prepare
Perform
model data
repeat
Genomic
Final matrix
Validation matrix
Marker Prepare
Perform
model data
repeat
Genomic
Final matrix
Validation matrix
Marker Prepare
Perform
model data
repeat



Difficulties in GS Application

- It is a whole chain of inter-connected tasks
- If we miss one link, predictions will not be confident
- Need a suit or software pipeline to deal with all steps with ease and confidence

ISMU 2.0



ISMU 2.0 Pipeline

- GUI for Genomic Selection
- Multicore Support
- R and Fortran Libraries for GS
- Project Mode Development
- IDE Supports
- Multiple Method & Traits at once
- Platform Support
 - Windows x64
 - Windows x32
 - CentOS x64
 - Ubuntu x64

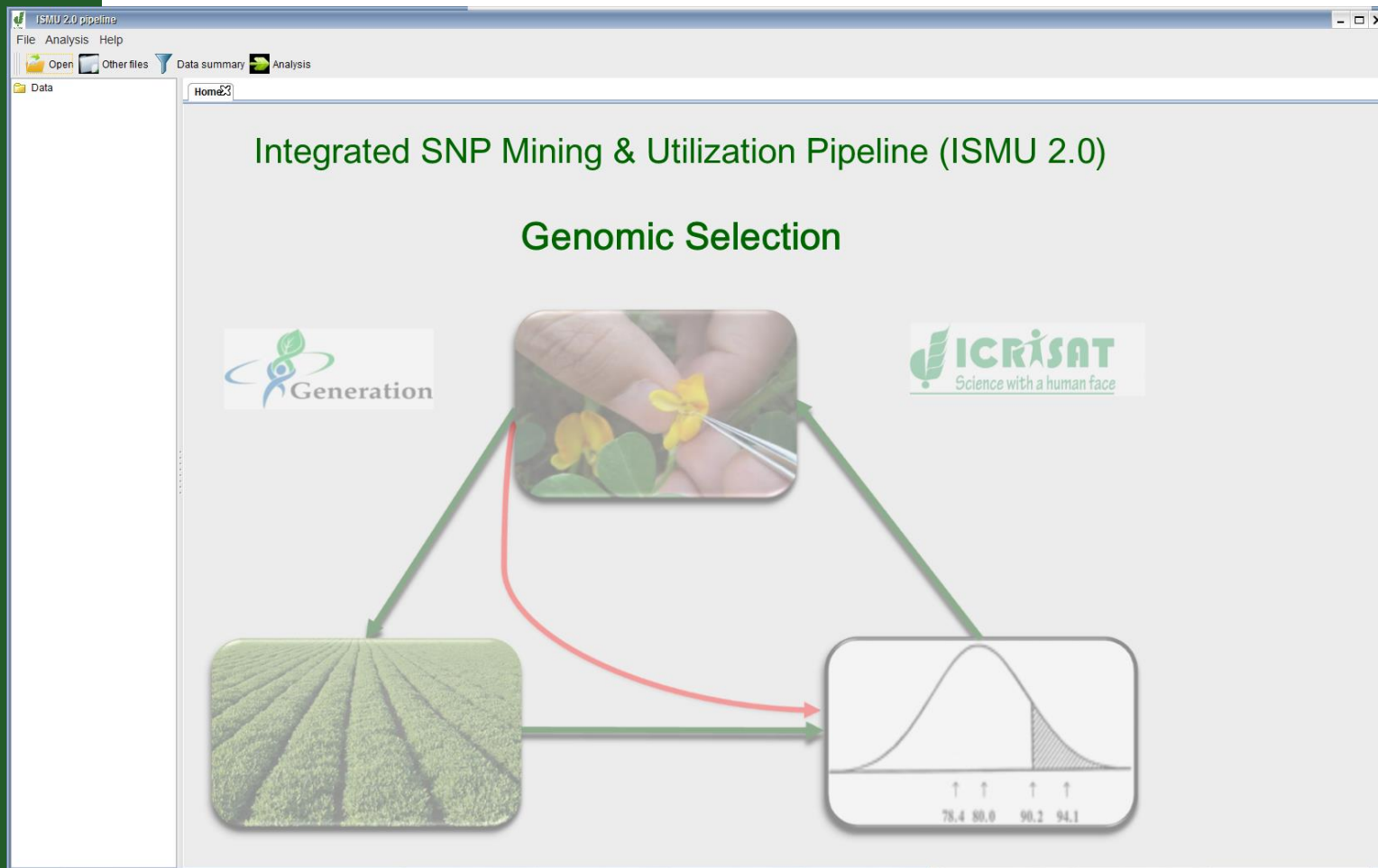


ISMU 2.0 Pipeline

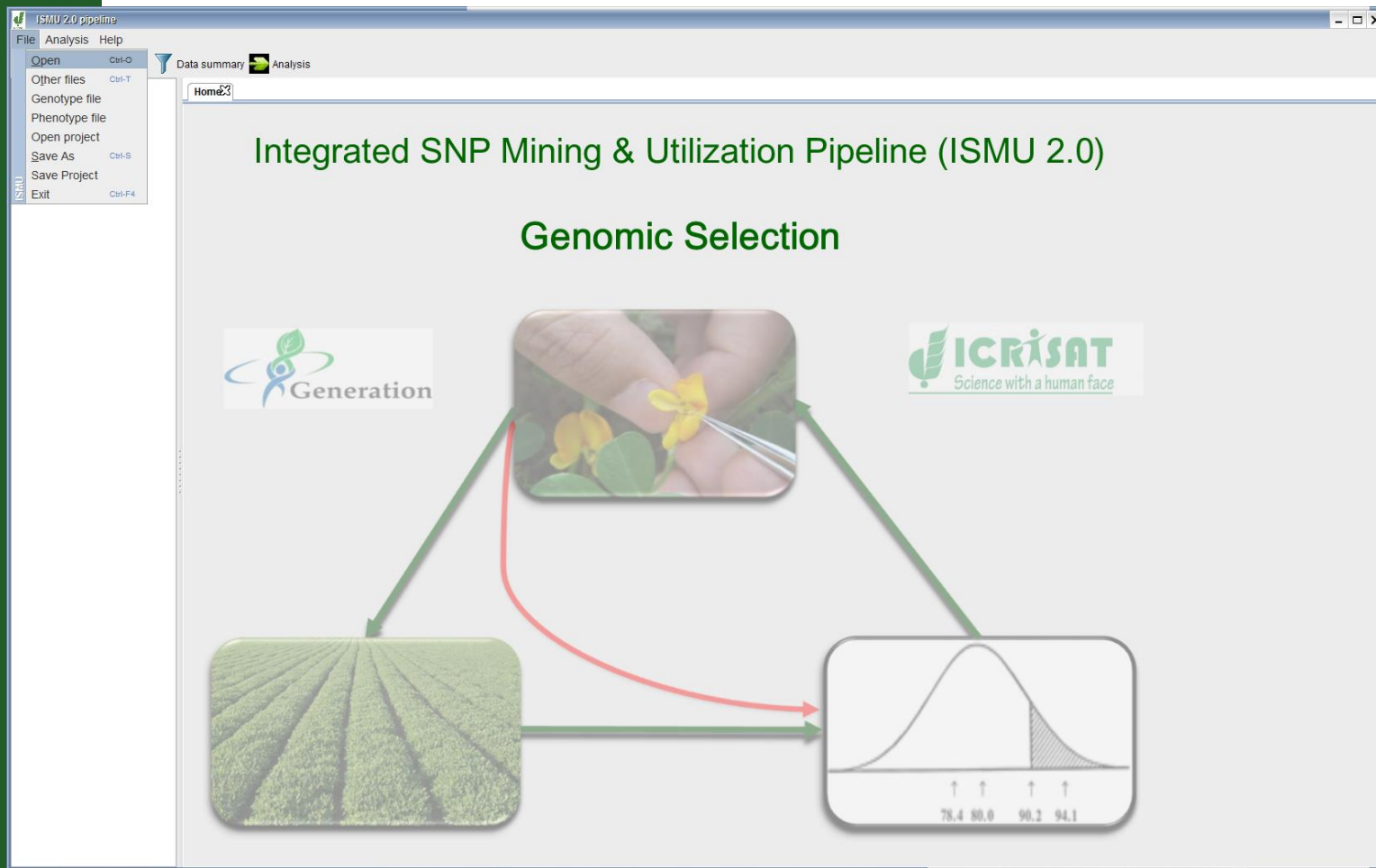
- Data Diagnostics
 - Graphical Summary
 - Tabular Summary
- Subset Data
 - Missing %
 - MAF
 - PIC
- Genomic Selection
 - RR-BLUP
 - Kinship Gauss
 - Bayesian LASSO
 - BayesA, BayesB and BayesC π
 - Random Forest Regression (RFR)
- Excel, HTML & PDF Output



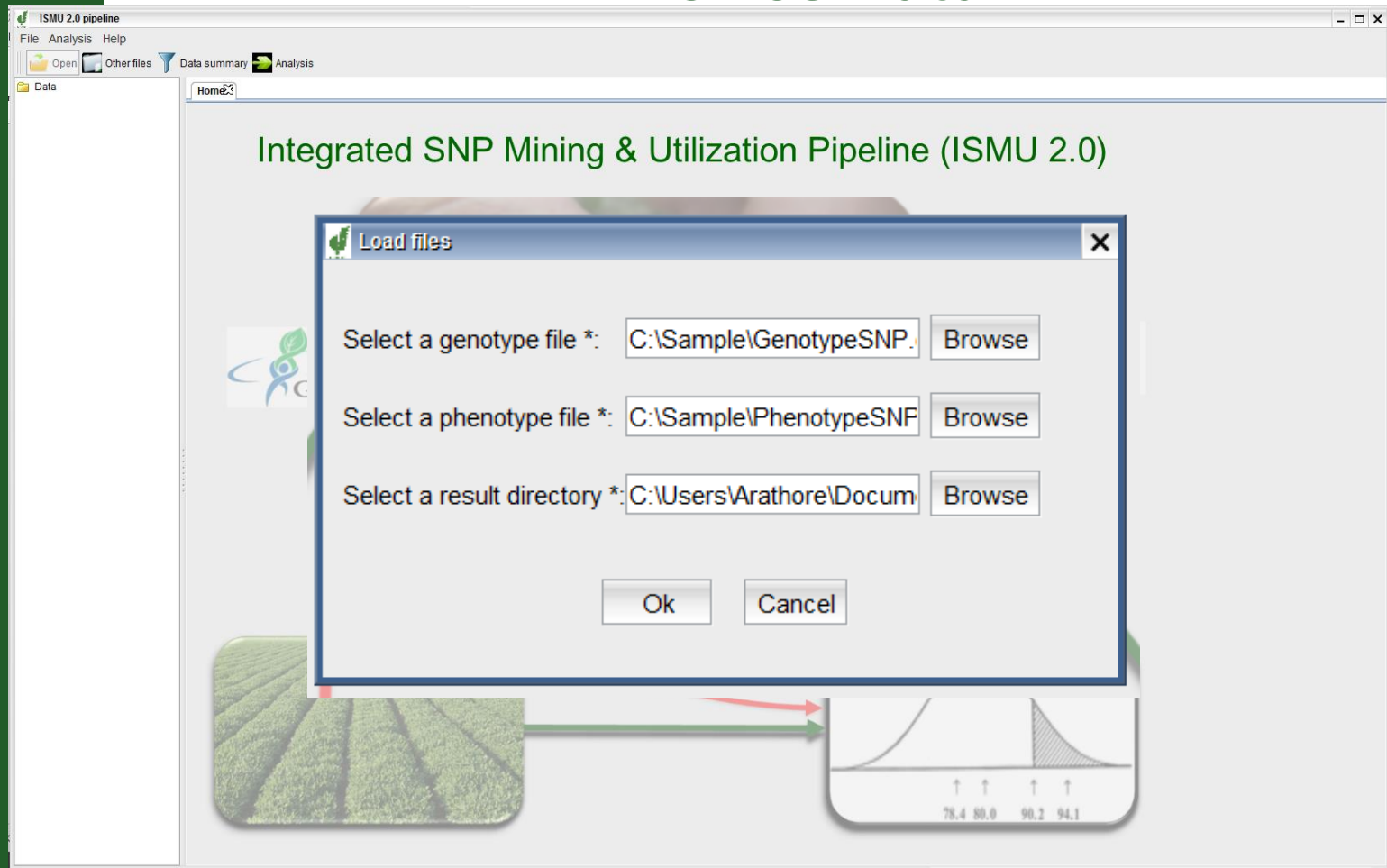
ISMU 2.0



ISMU 2.0



Browse Data



Calculation of Marker Summary

ISMU 2.0 pipeline

File Analysis Help

Open Other files Data summary Analysis

Data

- Genotype
 - GenotypeSNP.csv
- Phenotype
 - PhenotypeSNP.csv
- Relationship matrix
- Pedigree data
- Population structure
- Result Directory
- Log

Marker	G1	G2	G3	G4	G5	G6	G7	G8	G9	G10	G11	G12	G13	G14	G15	G16	G17	G18
M1	NA	C/C	C/C	T/T	C/C	T/T	C/C	NA	T/T	C/C	C/C	C/C	C/C	C/C	NA	C/C	T/T	
M2	G/G	A/A	A/A	A/A	A/A	A/A	A/A	A/A	A/A	A/A	A/A	A/A	G/G	G/G	NA	A/A	A/A	NA
M3	C/C	NA	C/C	C/C	C/C	C/C	NA	C/C	C/C	NA	C/C	NA	C/C	C/C	NA	NA	C/C	C/C
M4	T/T	T/T	T/T	T/T	G/G	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T
M5	G/G	A/A	NA	G/G	G/G	G/G	G/G	G/G	G/G	NA	G/G	NA	G/G	G/G	NA	NA	NA	G/G
M6	C/C	T/T	T/T	T/T	C/C	T/T	C/C	C/C	T/T	NA	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T
M7	T/T	G/G	T/T	G/G	T/T	G/G	NA	T/T	G/G	NA	T/T	NA	NA	T/T	T/T	T/T	NA	T/T
M8	G/G	G/G	G/G	G/G	G/G	G/G	G/G	NA	G/G	G/G	G/G	NA	G/G	NA	NA	NA	NA	G/G
M9	C/C	C/C	C/C	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	NA	T/T	NA	T/T	C/C
M10	A/A	A/A	A/A	A/A	A/A	A/A	NA	A/A	A/A	A/A	G/G	A/A	A/A	A/A	A/A	A/A	A/A	NA
M11																A/A	A/A	A/A
M12																NA	NA	T/T
M13																T/T	T/T	T/T
M14																NA	G/G	G/G
M15																A/A	A/A	A/A
M16																NA	NA	C/C
M17																NA	T/T	T/T
M18																NA	G/G	T/T
M19																NA	A/A	G/G
M20																G/G	A/A	G/G
M21																NA	NA	A/A
M22																NA	NA	C/C
M23																T/T	T/T	T/T
M24																A/A	NA	A/A
M25																NA	C/C	NA
M26																NA	NA	A/A
M27																G/G	G/G	G/G
M28																NA	NA	G/G
M29																C/C	C/C	C/C
M30																NA	NA	G/G
M31	C/C	C/C	C/C	C/C	C/C	C/C	NA	T/T	C/C	NA	C/C	NA	C/C	NA	NA	NA	NA	T/T
M32	T/T	G/G	G/G	G/G	G/G	G/G	G/G	NA	G/G	NA	G/G	G/G	G/G	T/T	NA	NA	NA	G/G
M33	T/T	C/C	C/C	C/C	C/C	C/C	C/C	NA	C/C	NA	C/C	C/C	C/C	T/T	NA	NA	NA	C/C
M34	A/A	G/G	A/A	G/G	G/G	G/G	NA	G/G	G/G	NA	A/A	NA	A/A	G/G	NA	G/G	A/A	G/G
M35	A/A	C/C	A/A	C/C	C/C	C/C	NA	NA	C/C	NA	A/A	NA	NA	A/A	C/C	C/C	C/C	C/C
M36	T/T	C/C	T/T	T/T	T/T	T/T	T/T	T/T	T/T	NA	T/T	NA	T/T	T/T	T/T	T/T	T/T	T/T
M37	A/A	G/G	A/A	G/G	A/A	A/A	NA	A/A	G/G	NA	A/A	G/G	G/G	A/A	A/A	A/A	A/A	A/A
M38	C/C	C/C	C/C	C/C	C/C	A/A	C/C	C/C	C/C	NA	C/C	C/C	C/C	A/A	C/C	NA	C/C	C/C
M39	A/A	G/G	G/G	G/G	A/A	A/A	A/A	A/A	G/G	A/A	A/A	G/G	G/G	A/A	A/A	NA	NA	NA
M40	T/T	C/C	C/C	T/T	T/T	T/T	T/T	T/T	T/T	NA	NA	C/C	NA	T/T	T/T	T/T	NA	T/T
M41	NA	C/C	NA	C/C	NA	C/C	C/C	NA	C/C	T/T	NA	C/C	NA	NA	T/T	NA	C/C	C/C
M42	T/T	G/G	G/G	G/G	T/T	T/T	T/T	T/T	G/G	NA	T/T	NA	NA	NA	T/T	T/T	G/G	NA
M43	NA	T/T	T/T	T/T	T/T	T/T	NA	T/T	T/T	T/T	C/C	T/T	T/T	T/T	T/T	NA	NA	T/T
M44	NA	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	T/T	NA	NA	T/T

Data summary

Select a genotype file Genotype2450.csv

Select a phenotype file Phenotype2450.csv

☒ Calculate % missing

☒ Calculate Data Summary

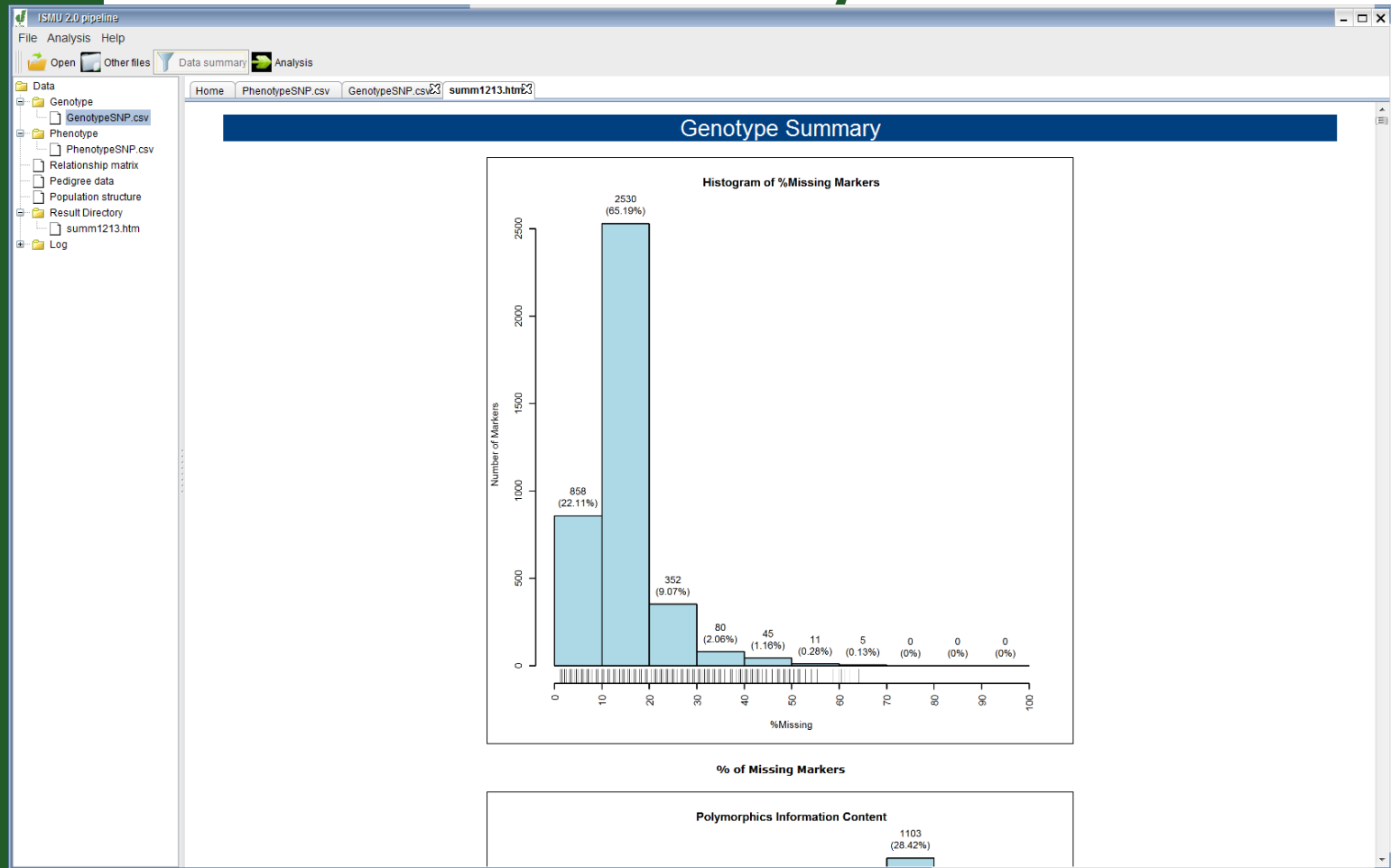
☒ Calculate PIC value

☒ Calculate MAF

Start Cancel



Summary Plots



Various Statistics

ISMU 2.0 pipeline

File Analysis Help

Open Other files Data summary Analysis

Data

- Genotype
 - GenotypeSNP.csv
- Phenotype
 - PhenotypeSNP.csv
- Relationship matrix
- Pedigree data
- Population structure
- Result Directory
 - summ1213.htm
 - resControl_Yd1620.htm
- Log

Home PhenotypeSNP.csv GenotypeSNP.csv summ1213.htm resControl_Yd1620.htm

Marker Statistics

	Marker	AlleleNumber	MisPercent	PIC	MAF
1	M1	2	13.3	0.206	0.135
2	M2	2	17.9	0.327	0.289
3	M3	2	32.5	0.036	0.019
4	M4	2	2.9	0.086	0.047
5	M5	2	11.2	0.278	0.211
6	M6	2	12.1	0.331	0.299
7	M7	2	15.0	0.266	0.196
8	M8	2	14.2	0.334	0.306
9	M9	2	20.8	0.171	0.105
10	M10	2	10.0	0.128	0.074
11	M11	2	11.2	0.174	0.108
12	M12	2	27.9	0.184	0.116
13	M13	2	9.6	0.141	0.083
14	M14	2	8.8	0.223	0.151
15	M15	2	10.1	0.219	0.177



Export to MS-Excel (Windows)

The screenshot illustrates the process of exporting data from a bioinformatics pipeline to an Excel spreadsheet. The 'ISMU 2D pipeline' window on the left shows a file tree with 'GenotypeSNP.csv' selected. The Excel window in the foreground displays a spreadsheet with the following data:

	A	B	C	D	E	F	G	H
1		Marker	AlleleNumber	MisPercent	PIC	MAF		
2	1	M1	2	13.3	0.206	0.135		
3	2	M2	2	17.9	0.327	0.289		
4	3	M3	2	32.5	0.036	0.019		
5	4	M4	2	2.9	0.086	0.047		
6	5	M5	2	11.2	0.278	0.211		
7	6	M6	2	12.1	0.331	0.299		
8	7	M7	2	15	0.266	0.196		
9	8	M8	2	14.2	0.334	0.306		
10	9	M9	2	20.8	0.171	0.105		
11	10	M10	2	10	0.128	0.074		
12	11	M11	2	11.2	0.174	0.108		
13	12	M12	2	27.9	0.184	0.116		
14	13	M13	2	9.6	0.141	0.083		
15	14	M14	2	8.8	0.223	0.151		
16	15	M15	2	10.4	0.249	0.177		
17	16	M16	2	19.2	0.375	0.5		
18	17	M17	2	16.2	0.276	0.209		
19	18	M18	2	14.2	0.283	0.218		
20	19	M19	2	13.8	0.271	0.203		
21	20	M20	2	23.3	0.359	0.375		
22	21	M21	2	45.8	0.328	0.292		

Analysis

Select file names from combo box

Genotype file name : Genotype2648.csv

Phenotype file name : Phenotype2648.csv

Covariate file name :

Select a method(s) to start analysis

R **Fortran**

☒ Ridge Regression BLUP

☒ Bayes Cpi

Select the trait

Trait2

Trait3

Bayes

Rounds 1,000

Burnin 100

Thinning 20

Forests 10

Processor

Cores 1

Replication 1

Fold 2

Subset

Missing markers 10

MAF (MAF) 0

Start

Cancel

Analysis

Performing analysis of trait : Trait1

Present running method : BayesLasso...

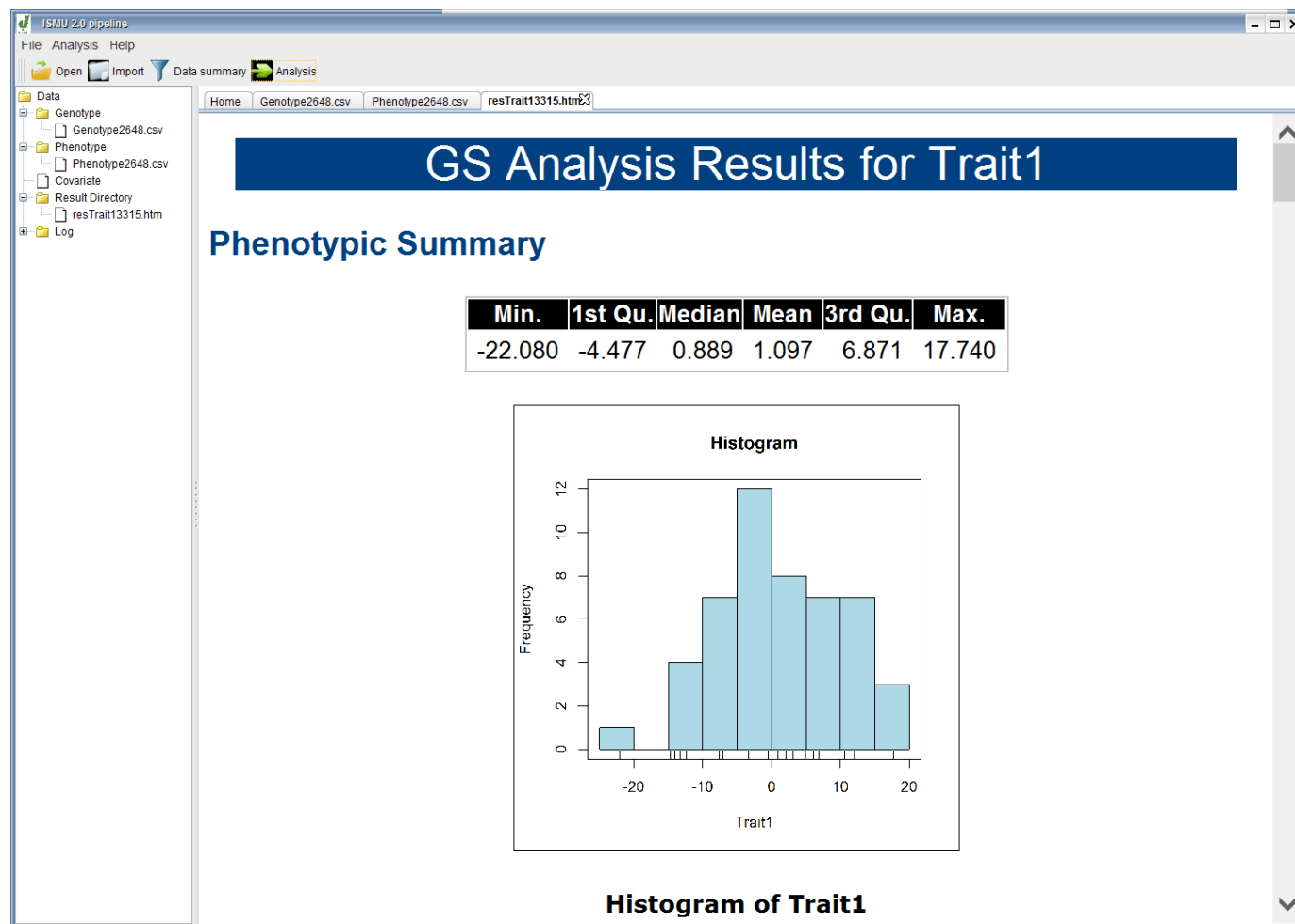
☒ RidgeRegression

☒ BayesCpi

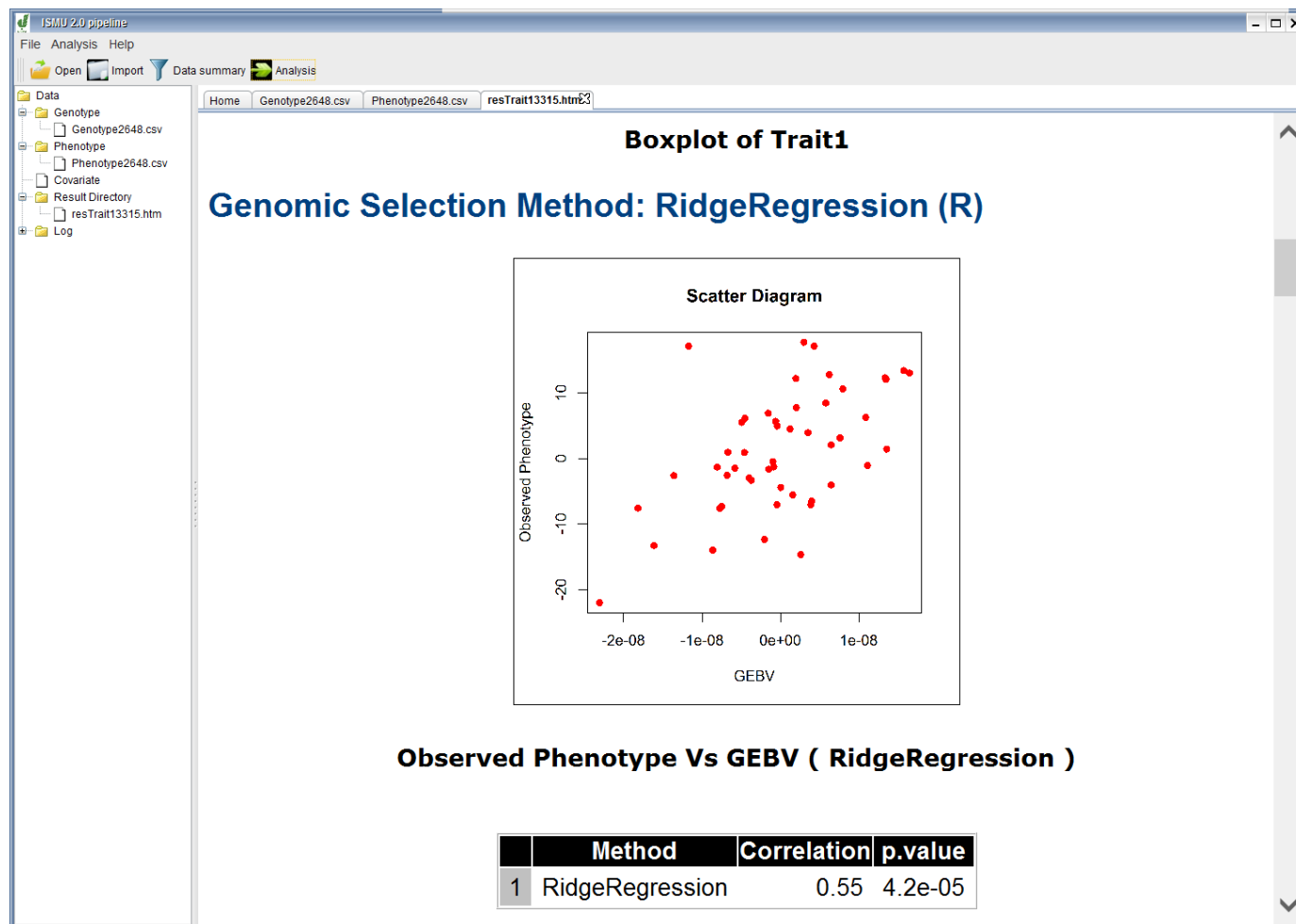
☐ BayesLasso

Cancel

GS Results

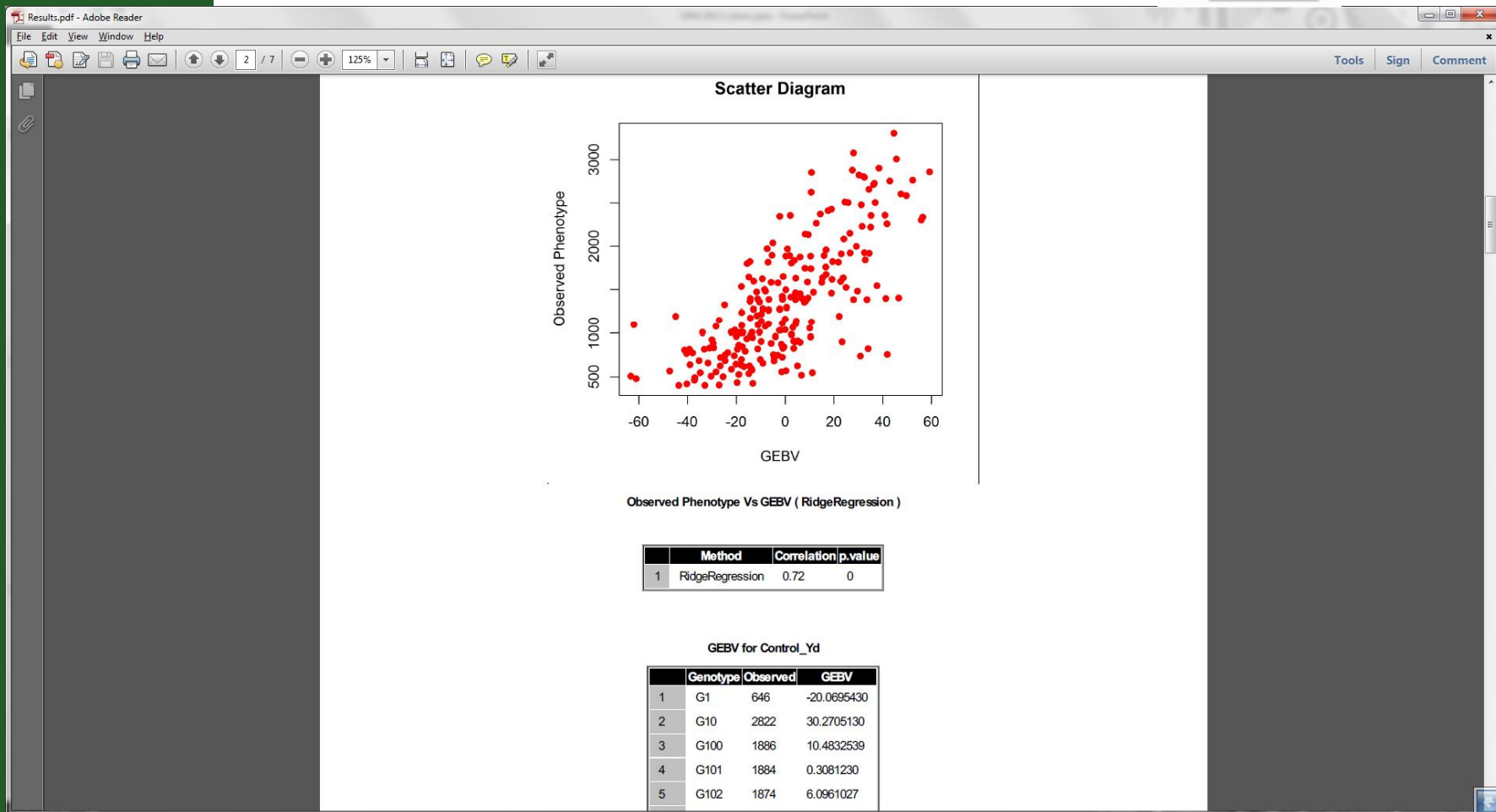


GS Results

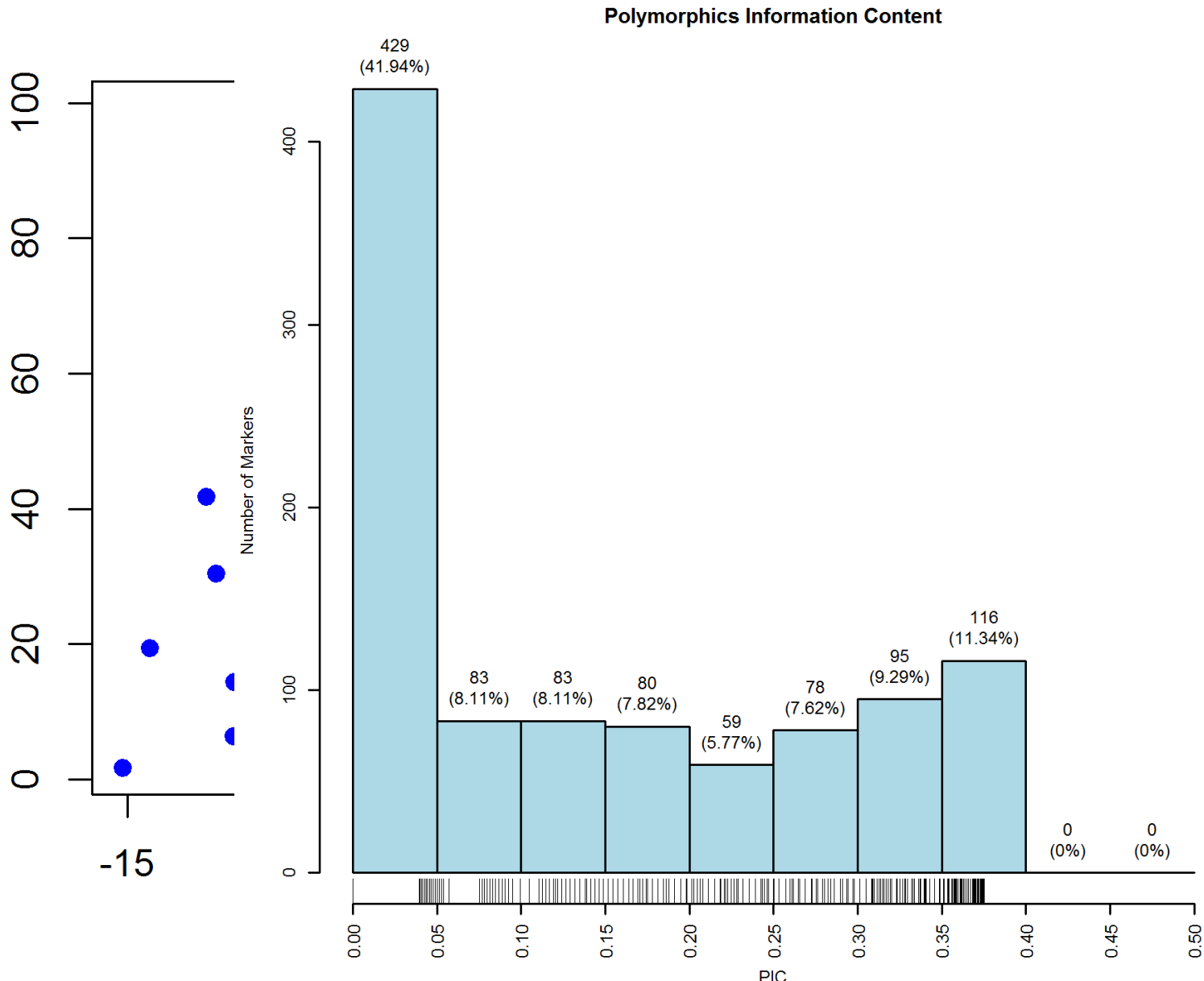


Export to PDF

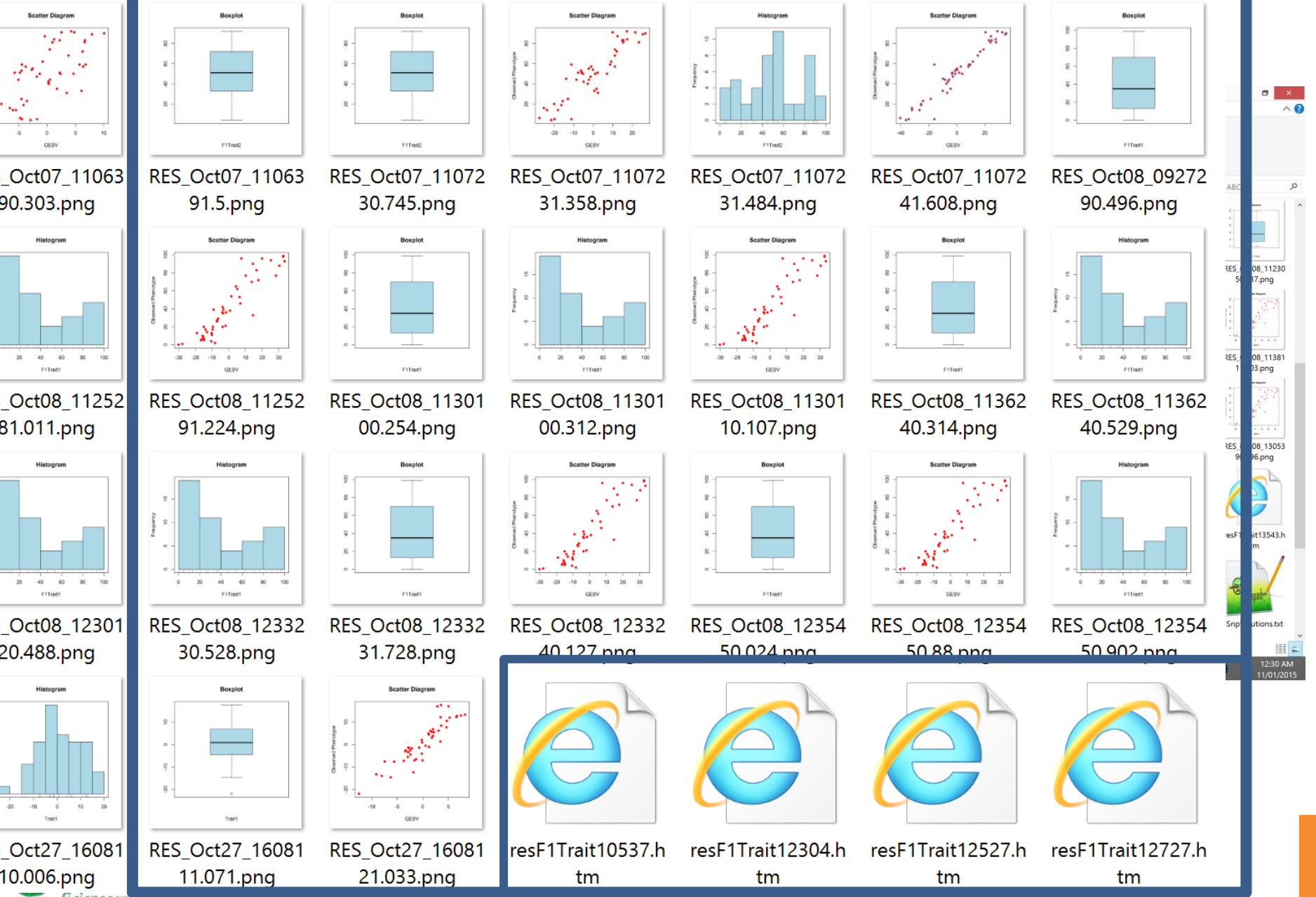
PDF



Export to High Quality Graphics 300DPI



Graphs & HTML Reports saved



Support Large Data Sets : R & F Cocktail

- R is relatively slow when apply GS on large data sets
 - 1500 Individuals and 50 K Markers?
 - Or Even 5000 Individuals and 50 K Marker?
- A cocktail of Native FORTRAN binaries and R was used as a solution
 - 5-6 times faster
- FORTRAN was used for data processing and fitting GS Models
- R was used to compile generated results and produce high quality graphics and dynamic reports



ISMU 2.0 pipeline

File Analysis

Select file names from combo box

Genotype file name : Genotype0849.csv

Phenotype file name : Phenotype0849.csv

Covariate file name : Select

Select a method(s) to start analysis

R Fortran

☐ Ridge Regression

☐ BayesA

Select the trait(s) for analysis

Trait1
Trait2
Trait3

>
>>
<<
<

Data Subset

Percentage(%) of missing markers 10

PIC value 0

Minor allele frequency (MAF) 0

Additional Parameters

Bayes

Rounds 1,000

Burnin 100

Thinning 20

Random Forest

Forests 10

Processor

Cores 1

Cross Validation

Replication 1

Fold 2

Start

Cancel

Genomic Selection Method RidgeRegression (Fortran)

Index	Marker	Value 1	Value 2
37	G42	-5.578000	-5.0972290
38	G43	5.480000	3.7197820
39	G44	13.436000	11.0970100
40	G45	13.019000	10.1830100
41	G46	5.650000	2.8407890
42	G47	-1.692000	-3.4900590
43	G48	1.051000	1.3180110

Summary of Selected GS Methods

	Method	Correlation	Prob.t
1	RidgeRegression	0.627	0.000
2	CrossValidation(R=1,F=2)	0.513	0.0025
3	BayesA	0.57	0.000
4	CrossValidation(R=1,F=2)	0.45	0.002

Results Generated by ISMU 2.0 : Mon Jan 12 10:30:30 AM 2015 (Total Time: 00 Hr : 00 Min : 04 Sec)





Plans

- Make online version
- Support import of various popular formats
- VCF, PED, hapmap and etc
- Integration of newer methods
- Multiple trait GS
- GxE



Acknowledgements



Cornell University



THE SAINSBURY LABORATORY



RESEARCH
PROGRAM ON
Grain Legumes



RESEARCH
PROGRAM ON
Dryland Cereals

Thanks...