

Bioinformatics workshop

3rd SoyaGen Annual Meeting

7 December 2018





SOYA ON

Basic data handling, diversity analysis and phylogenetic trees, parental line selection and confirmation of hybrid status with TASSEL

Martine Jean 3rd SoyaGen Annual Meeting, December 7th, 2018 Université Laval, Québec

www.soyagen.ca























































Taxa Summary								
Taxa Name N SoyaGenDemo_p01c01_Sample2.sort SoyaGenDemo_p01c01_Sample26.sort SoyaGenDemo_p01c01_AAC-Mandor.sort	umber of Gam 18803 0 18803 0 18803 0	Prop Nu 0 0	umber He Proporti 2645 2610 2455	on Heterozygous 0.14067 0.13881 0.13056	Iter Analysis Results GBSv2 Filter Cenotype Table Sites Filter Cenotype Table Taxa Traits	Warning: if this option is missing, you are in a		
SoyaGenDemo_p01c01_Sample3.sort SoyaGenDemo_p01c01_Sample4.sort SoyaGenDemo_p01c01_Sample7.sort SoyaGenDemo_p01c01_S03-W4.sort SoyaGenDemo_p01c01_Sample41.sort	18803 0 18803 0 18803 0 18803 0 18803 0 18803 0	000000	2408 1995 1993 1831 1686	0.12806 0.1061 0.10599 0.09738 0.08967	Not Recommended ¥ Sites ¥ Site Names Taxa Filter N	Heip I		
 Unexpected excess of usually indicative of a 	f heterozyg ı technical amination	;ous g proble	enotypes ir em	n a cultivar is	Max Heterozygous Propor Max Heterozygous Propor Taxa List	arcel entropy liser Manual		
Accidental outcro	SS			·				





	0. 0
Renaming file in Excel:	Step 2
0	•
Dataset with original sample names	
🛛 🕼 🕼 🕫 🕐 👌 🔹 🔥 SoyaGendierne, pöllcitt järpuned, 91st88003.hmp 🛛 🖓 - Nachenster dava is in 🧐 -	
Accesi Insertion More explage Formulaes Dennées Bérlelon Africhage	Using the "replace" function from
Caster (Corpe) - 12 - A - A - E = T - Standard - Markan - Standard - Standard - Markan - Standard - Standard - Markan - Standard - Standard-	
under ge auf 2 an	Excel to modify sample names
AF A F F F F F F F F F F F F F F F F F	
na alietes dreim pes strait assentate centri pretS0 auxqSD pareES0 (Cocee Septemberg.g000), jampel and Septemberg.g0000, jampel and Septemberg.g0000, jampel and	Romplacer
542_05053 1/C I 3505/1+ MA NA NA MA NA NA T T T Y	Bechester
1932_193864 C/T I 235064 NA NA NA NA NA NA TC C T 052_293524 G/L 2351624 NA NA NA NA NA KA G G A	SoyaGenDerno, pD1x01,
1952,293864 1/6 i 299964 NA NA NA NA NA NA T T G	Dans : Feulite B Respector ta casee
511_0550 A/6 I 25007- NA NA NA NA NA NA A A	Sens) Par Sgre
1921_3782.7 L/G 1_2762.07 + 18.4 N/A N/A N/A N/A N/A N/A K A C 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5 5	
g > StyleDemo p01c01 inputed 91s1 +	Remplacer par
ha 📰 🗉 🖉 + 50 %	
	Remplacer Remplacer Ituit Fermer Elizabet
Dataset with modified sample names	The second se
🙃 💿 💼 Sayalambana, joloot jaquad (httilliot have 🔹 💽 - tamaaanaa dalaa ta ay	Berglinae.
Assuall Inversion Mise en page Formules Données Révision Affichage 🌲 Portager A	Protocolo 1
Celibri (Corps) + 12 + A+ A+ Standard - B: Maximi torm conditionals - Thesis - ∑ A → .	100 CH
Caller at 6 7 5 *	Dane : Eeulle C Respecter la casse
1 X V & A	Sans Dar ligner 0
net alleles them pos exact exactly entry postSP enapSP predSP SampleS SampleS SampleS SampleS	
302_230531 T/C 1 130032 • NA NA NA NA NA NA A A T T Y 502_230531 T/C 1 135592 • NA NA NA NA NA A A C C T	Remplecer par :
501_259644 C/T 1 259644 - NA NA NA NA NA A C C T	An and the second se
542,59535 W/A 1 295354 W/A 1 295354 W/A 1 295354 W/A 1 A 1 A 1 A 1 A 1 A 1 A 1 A 1 A 1 A 1	Remplacer Remplacer tour Fermer Sulvant
150_25425 0/A 1 254224 NA NA NA NA NA NA G 0 A 150_25425 0/A 1 254224 NA NA NA NA NA NA G G G	
500_270027 C/G 1 270227 - NA NA NA NA NA E C S	
S02_SMBEX A/T 1_384882 + AA NA NA NA NA NA A A A W	



























el
Ne l
BO Ker - F
4,3 Mu
14,8 Mu 4
Ouver
<i>K</i> 2
e en page Formulies Donnáes Rávi
) • 12 • A• A• = ⊗
 ⊥·▲·▲· = = = =
252
D E F G Sample3 Sample4 Sample5 AAC-Mandor
07 0.3389103 0.2768973 0.31689093 0.37153646 0.2365318 0.28202945 0.2542945 0.32401744
0.0 0.26341543 0.25259268 0.28056683 M5 0.26341543 0.0 0.2793969 0.34061053
0 0.25259268 0.2795969 0.0 0.2479575 144 0.28056693 0.34061053 0.24796575 0.0
NAT 0.37392968 0.29359677 0.3378716 0.36225602 R01 0.30120035 0.26225517 0.3770941 0.4043356
91.0.301/0003.9.20134317 0.3/70341 0.4042709
1 3 1 5 9 9





SOYA GEN Using a distance matrix to compare lines and select parental combination that maximize allele diversity Use the color-code the distance matrix to visualize the level of similarity between the lines • Use the min and max values per line to identify the lines that are the most similar or different to your line of interest • Unexpected very high similarity between lines may indicate sample mislabelling or breeding Highly similar Highly different problems lines lines • Use the max value for the entire dataset to identify the lines that would make the most divergent parental combination www.soyagen.ca



















Using haplotypes to predict allelic states at important agronomic genes

Objectives

SOYA GEN

 Select lines with a desired trait using haplotype information around the gene of interest

Protocol

- Step 1. Get the genomic position of your gene of interest from a database
- Step 2. Select SNPs surrounding your gene of interest in your dataset
- Step 3. Use the "Linkage Disequilibrium" and "LD plot' functions to create a LD plot
- Step 4. Identify SNPs that always co-segregate (LD=1) with your gene of interest
- Step 5. Create haplotypes with the "Cluster genotype" function
- Step 6. Create a tree to visualize lines belonging to each cluster
- Step 7. Use reference lines to identify the allele present in uncharacterized lines

www.soyagen.ca



GEN	Late maturi cultivar	ty Early m cult	aturity ivar	 An example: finding haplotypes for the E1 gene The E1 gene is one of the most important genes controlling maturity in soybean Among the 136 Canadian lines characterized by Tardival et al. (2108) most sarry either the s1 m/ 		
Gene	Gene model	Allele E1	Count %	or <i>e1-as</i> alleles.		
E1	Glyma.06g207800	e1-ni e1-as e1-fs Tardivel et al. 20	39 58.2 25 37.3 2 3.0	 KASP assays could be used to determine the allele present at the E1 gene in multiple lines Costly and specific to a single allele of a single gene 		
	SNP 1	A — T allele1 🗖 all	ele2	 SNPs from GBS dataset can often be grouped to create haplotypes specific for each allele of a gene of interest 		
ww	SNP 2 - C - G			Haplotypes for multiple genes of interest can be assessed in the same GBS dataset		



















