

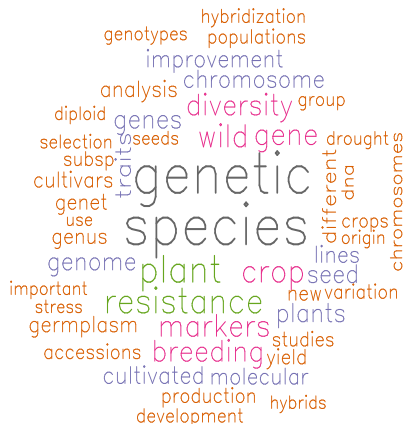
# Measuring Intraspecific Genetic Diversity

Fernando Henrique RB Toledo\*

September 29-30, 2016



<f.toledo@cgiar.org>




## 💡 Outline:


- 1 Introduction
  - 1.1 Frequency of Alleles & Genotypes
- 2 Diversity
  - 2.1 Distances & Diversity Indices
- 3 Intraspecific
  - 3.1 Wright Statistics ( $F$ )
- 4 Cluster to identify "groups"
- 5 BIO-R to make the things easy...

# 1 Introduction and some Definitions

- **Genotype:**

Genotype is ...  \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

- **Markers:**

Genetic Markers are ...  \_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_  
\_\_\_\_\_

- **Population:** ... [CONT]

## Populations

- **Population Genetics:** Studies the heredity's mechanisms at the population level.
- **Population:** Set of conspecifics individuals in the *same place and time*, which have the ability to mate (exchange alleles).



In a population, the individual has a transitory importance, what matters are the alleles it has, which will be transmitted to subsequent generations

## Reproductive systems

- **Allogamous:** frequency of cross pollination is  $\geq 95\%$  *e.g.*, maize
- **Autogamous:** frequency of cross pollination is  $\leq 5\%$  *e.g.*, wheat.

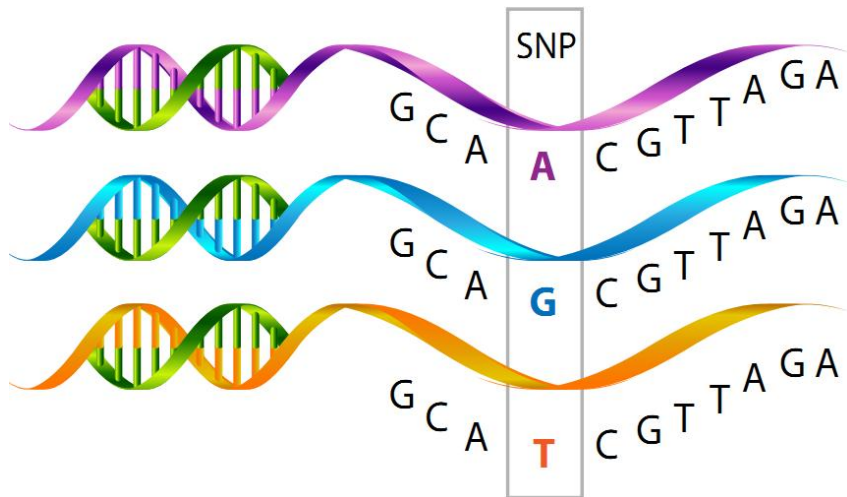
## 1.1 Frequency of Alleles &amp; Genotypes

Genotype		Phenotype	Observed	Frequency
<i>BB</i>	→	White	100	0.05
<i>Bb</i>	→	Yellow	1000	0.50
<i>bb</i>	→	Red	900	0.45
<b>Total:</b>			2000	1.00

$$f(B) = p = 0.05 + \frac{0.50}{2} = 0.30$$

$$\begin{aligned} f(b) &= q = 0.45 + \frac{0.50}{2} = 0.70 \\ &= (1 - p) \end{aligned}$$

## Genetic Markers



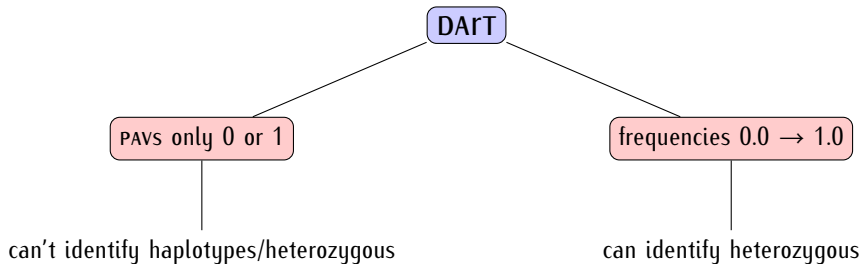
## Examples

Marker	Amount	Expression	Polymorphic degree	Specific by locus
Isoenzyme	< 100	codominant	low	✓
RFLP	$\infty$	codominant	medium	✓
RAPD	$\infty$	dominant	medium	–
SSR	$\infty$	codominant	very high	✓
SCAR	$\infty^*$	both	low	✓
AFLP	$\infty$	dominant	high	–
SNP	$\infty$	codominant	very high	✓
<b>DART</b>	$\infty$	both	very high	✓**

\* After the deployment of other kind of markers.

\*\* Clone aren't but markers are.

## DART flexibility...





## Maize Bulks w/ frequencies

Marker	Allele	Bulks					
		1	2	3	4	...	$g$
1	1	0.67	NA*	0.57	NA	...	0.36
1	2	0.33	NA	0.43	NA	...	0.64
2	1	NA	1.00	1.00	1.00	...	1.00
2	2	NA	0.00	0.00	0.00	...	0.00
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
$m$	1	0.00	NA	1.00	0.00	...	1.00
$m$	2	1.00	NA	0.00	1.00	...	0.00

\* NA means *Not Available* or missing.

Frequencies, so:

$$p + q = 0.67 + 0.33 = 1.00$$

## Maize Genotypes w/ PAV – (presence/absence)

Marker	Allele	Bulks					
		1	2	3	4	...	$g$
1	1	1	NA	1	NA	...	1
2	1	NA	1	1	1	...	1
⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
$m$	1	1	NA	1	1	...	1

\* NA means *Not Available* or missing.

Allele is always 1... the reference allele.



When Genotypes belong to the same bulk, the summary of their PAV generates frequencies

## 2 Genetic Diversity

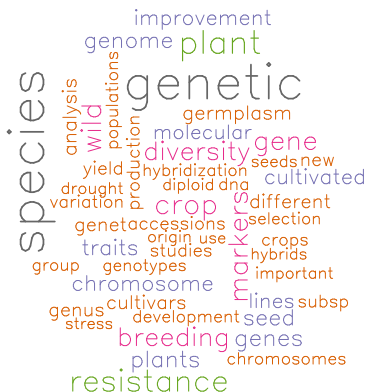
### 2.1 Genetic Distances

#### Markers w/ allelic information

For example: SSR, SNP and DART.

- Euclidean;
- Modified Rogers; and
- Cavalli-Sforza & Edwards.

References: [1, 2, 3]



## Euclidean:

$$Ed_{[x,y]} = \sqrt{\sum_l^L \sum_a^A (\hat{p}_{xla} - \hat{p}_{yla})^2}, \quad (0.0 \leq Ed_{[x,y]} \leq \sqrt{2L}^*)$$

where:

$\hat{p}_{xla}$  is the allele frequency for allele  $a$  at locus  $l$  in the genotypes  $x$ ;

$\hat{p}_{yla}$  is the same as above for genotype  $y$ ;

$L$  is the number of locus; and

$A$  is the number of alleles at locus  $l$ .

\*This is the estimator and its respective domain as shown by [3].

## Modified Rogers:

$$Rd_{[x,y]} = \frac{1}{\sqrt{2L}} \sqrt{\sum_l^L \sum_a^A (\hat{p}_{xla} - \hat{p}_{yla})^2}, \quad (0.0 \leq Rd_{[x,y]} \leq 1.0)$$

where:

$\hat{p}_{xla}$  is the allele frequency for allele  $a$  at locus  $l$  in the genotypes  $x$ ;

$\hat{p}_{yla}$  is the same as above for genotype  $y$ ;

$L$  is the number of locus; and

$A$  is the number of alleles at locus  $l$ .

For non-diploid species, replace  $2L$  by the number of copies.

## Cavalli-Sforza &amp; Edwards:

$$Cd_{[x,y]} = \sqrt{\frac{1}{L} \sum_l^L \left( 1 - \sum_a^A \sqrt{\hat{p}_{xla} \times \hat{p}_{yla}} \right)}, \quad (0.0 \leq Cd_{[x,y]} \leq 1.0)$$

where:

$\hat{p}_{xla}$  is the allele frequency for allele  $a$  at locus  $l$  in the genotypes  $x$ ;

$\hat{p}_{yla}$  is the same as above for genotype  $y$ ;

$L$  is the number of locus; and

$A$  is the number of alleles at locus  $l$ .

## Mean, Variance and Standard Error

$d_{[x,y]}$  is the distance between genotypes  $x$  and  $y$ , thus:

- **Mean:**  $\mu_{d_{[x,y]}} = \frac{1}{\frac{A(A-1)}{2}} \sum_{x < y} d_{[x,y]}$ ;
- **Variance:**  $\sigma_{d_{x,y}}^2 = \frac{1}{\frac{A(A-1)}{2} - 1} \sum_{x < y} (d_{[x,y]} - \mu_{d_{[x,y]}})^2$ ; and
- **Standard Error:**  $S_{\bar{d}_{[x,y]}} = \sqrt{\frac{\sigma_{d_{[x,y]}}^2}{\frac{A(A-1)}{2}}}$ .



$\frac{A(A-1)}{2}$  is the  $2 \times 2$  combination of the alleles i.e., pairs in the distances estimators

## Markers w/o allelic information – PAVS

		Genotype $x$	
		1	0
Genotype $y$	1	$a$	$b$
	0	$c$	$d$

- **Simple Matching:**  $d_{[x,y]} = \frac{b+c}{a+b+c+d}$ ;
- **Jaccard:**  $d_{[x,y]} = \frac{b+c}{a+b+c}$ ; and
- **Nei and Li (or Dice):**  $d_{[x,y]} = \frac{2(b+c)}{2a+b+c}$ .

The same references for the previous cases i.e., [1, 2, 3].



## Diversity Indices



A locus is polymorphic if, and only if, the frequency of one of its alleles is  $\leq 0.95$  or  $0.99$

### Raw statistics

- Polymorphic proportion:  $P = \frac{n_{poly}}{n_{total}}$ ; and
- Mean allele number by locus:  $n_a = \frac{1}{L} \sum_l^L n_l$

### Others

- Observed Heterozygosity ( $H_o$ );
- Expected Heterozygosity ( $H_e$ );
- Effective Number of Alleles ( $A_e$ ); and
- Shannon Index ( $SH$ ).

**Observed Heterozygosity:**

*It is defined as the percentage of heterozygous loci per individual or the number of heterozygous individuals per locus.*

**Expected Heterozygosity:**

$$He = \frac{1}{L} \sum_l^L \left( 1 - \sum_a^A \hat{p}_{la}^2 \right), (0.0 \leq He \leq 1.0)$$

where:

$\hat{p}_{la}$  is the estimated frequency of the allele  $a$  at locus  $l$ ; and the other quantities ( $L$  and  $A$ ) were already defined.



**NOTE:**

$$\sum_a^A \hat{p}_{la} = 1.0$$

## Effective Number of Alleles:

$$Ae_l = \frac{1}{\sum_a^A \hat{p}_a^2}$$
$$Ae = \frac{1}{L} \sum_l^L Ae_l$$

**NOTE:**

$$Ae_l = \frac{1}{1 - He_l}$$

## Shannon Index:

- Total Frequencies:







$$SH_{\text{Total}} = - \sum_a^A \hat{p}_a \log_{10} (\hat{p}_a) , \sum_a^A \hat{p}_a = 1.0$$

- By Locus Frequencies:

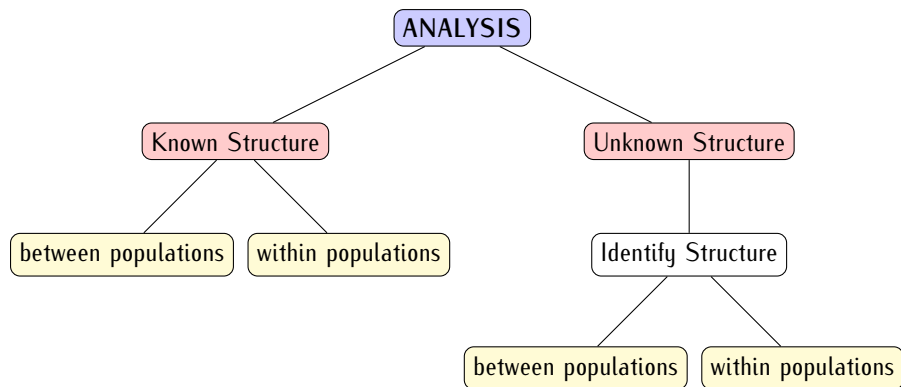
$$SH_{\text{Locus}} = - \sum_a^A \hat{p}_a \log_2 (\hat{p}_a) , \sum_a^A \hat{p}_a = L \therefore (0.0 \leq SH_{\text{Locus}} \leq L)$$

## Diversity Indices for PAVs

Statistics ...:

- Polymorphic proportion \_\_\_\_\_ 
- Mean number of alleles by locus \_\_\_\_\_ 
- Observed Heterozygosity ( $H_o$ ) \_\_\_\_\_ 
- Expected Heterozygosity ( $H_e$ ) \_\_\_\_\_ 
- Effective Number of Alleles ( $A_e$ ) \_\_\_\_\_ 
- Shannon Index ( $SH$ ) \_\_\_\_\_ 

### 3 Intraspecific Diversity (*1<sup>st</sup> branch*)



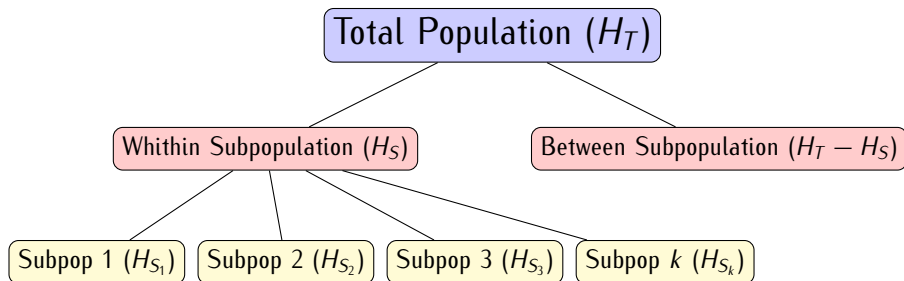


## Recapping ...

Parameter	Estimator
Population (by locus)	$He_l(H_{T_l}) = 1 - \sum_a \hat{p}_{a_l}^2$
Population (average)	$He(H_T) = \frac{1}{L} \sum_l He_l$
Subpopulation <sub><i>i</i></sub> (by locus)	$He_{i_l} = 1 - \sum_a \hat{p}_{i_{a_l}}^2$
Subpopulation <sub><i>i</i></sub> (average)	$He_i = \frac{1}{L} \sum_l He_{i_l}$
<b>Within Subpopulation</b>	$H_S = \frac{1}{J} \sum_i He_i$
<b>Between Subpopulation</b>	$D_{ST} = H_T - H_S$



Total = Between + Within



### 3.1 Wright Statistics ( $F$ )

$F_{IS}$  Heterozygosity proportional deviations within subpopulations:

$$F_{IS} = \frac{H_S - H_o}{H_S} [-1; 1]$$

$F_{IT}$  Overall heterozygosity proportional deviation (inbreeding coefficient):

$$F_{IT} = \frac{H_T - H_o}{H_T} [-1, 1]$$

$F_{ST}$  Heterozygosity proportional deviation between subpopulations:

$$F_{ST} = \frac{H_T - H_S}{H_T} [0; 1]$$

Nei, 1987 [4] ...:

1. Obtain  $H_T$ :

$$H_T = \frac{1}{L} \sum_l^L He_l$$

2. Obtain  $H_S$  (by locus):

$$H_S = \frac{1}{L} \sum_l^L He_{S_l}$$

3. Thus, obtain  $G_{ST}$ :

$$G_{ST} = \frac{D_{ST}}{HT} = \frac{H_T - H_S}{H_T} = 1 - \frac{H_S}{H_T}$$

**Berg, 1997 [1] ...:**

1. Obtain  $G_{ST}$  by locus:

$$G_{ST} = \frac{D_{ST}}{HT} = \frac{H_T - H_S}{H_T} = 1 - \frac{H_S}{H_T}$$

2. Summarize as average:

$$\bar{G}_{ST} = \frac{1}{n_a} \sum_n^{n_a} G_{ST_n}$$

3. That is the same as:

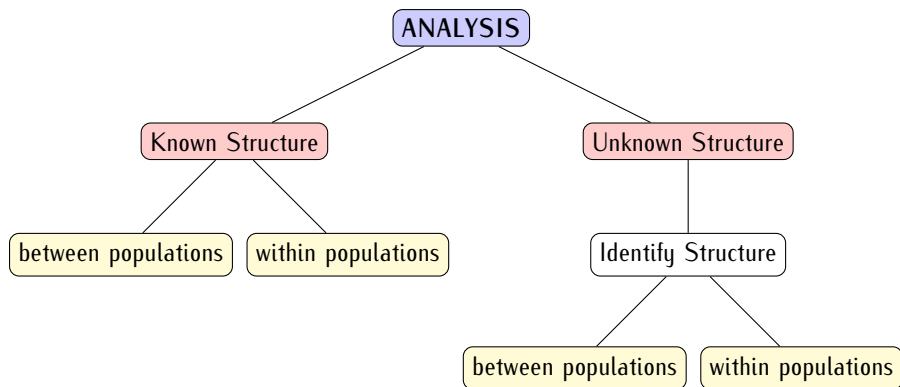
$$\bar{G}_{ST} = 1 - \frac{\sum_n^{n_a} \left( \frac{H_{S_n}}{H_{T_n}} \right)}{n_a}$$

Be aware ... 

- When we have just one (1) allele:  $G_{ST} = F_{ST}$ ;
- $G_{ST}$  generalizes  $F_{ST}$ ;
- $G_{ST}$  is the proportion of total diversity that is between subpopulations; and
- Thus,  $(1 - G_{ST})$  is the proportion of total diversity within subpopulations.
- The meaning of diversity according to  $F_{ST}$ :

$F_{ST}$	means
$[0; 0.05)$	small
$[0.05; 0.15)$	medium
$[0.15; 0.25)$	high
$\geq 0.25$	very high

## ... Intraspecific Diversity (2<sup>nd</sup> branch)



## 4 Identifying Groups

**OBJECTIVE:**

Minimize within variability → Maximize between variability

References can be consulted for deep understanding:

- Foundation book ...[5];
- Foundation paper in genetics (Nei) ...[4]; and
- More modern reference ...[6].

$$\mathbf{Y}_{n \times p} = \begin{bmatrix} y_{1,1} & y_{1,2} & y_{1,3} & \cdots & y_{1,p} \\ y_{2,1} & y_{2,2} & y_{2,3} & \cdots & y_{2,p} \\ y_{3,1} & y_{3,2} & y_{3,3} & \cdots & y_{3,p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ y_{n,1} & y_{n,2} & y_{n,3} & \cdots & y_{n,p} \end{bmatrix}$$

**NOTE:**

Everything starts from obtaining the distance matrix between all  $2 \times 2$  pairs ...

$$\frac{n(n-1)}{2}$$



$$D = \begin{bmatrix} 0 & d_{1,2} & d_{1,3} & \cdots & d_{1,j} & \cdots & d_{1,n} \\ & 0 & d_{2,3} & \cdots & d_{2,j} & \cdots & d_{2,n} \\ & & 0 & \cdots & d_{3,j} & \cdots & d_{3,n} \\ & & & \ddots & \vdots & & \vdots \\ & & & & 0 & \cdots & d_{j,n} \\ & & & & & \ddots & \vdots \\ & & & & & & 0 \end{bmatrix}$$

$d_{i,i} = 0.0$  and  $D$  is symmetric  $\therefore d_{i,j} = d_{j,i}$

## Clustering Methods

- **Geometrical:**
  - Hierarchical;
  - Neighbor Joining;
  - *k-means* (density search); and others
- **MANOVA:**

$$\text{Total} = \text{Between} + \text{Within}$$

- **Statistical:**
  - Mixture and Bayesian (STRUCTURE)

## Hierarchical:

1. All distances;
2. Find the most close individuals (smaller distance); and
3. Iterate over that.
  - UPGMA: merge groups with smaller distance;
  - WARD: merge groups that generate a new one with smaller  $S.S.W$ ; and
  - NJ: merge groups if: (i) smaller distance & (ii) higher distances in comparison to the others.

## Bayesian:

The model account for the presence of HW or LD, introduces population structure... find population structure (groups) with the smaller possible disequilibrium [6].

## Graphical Representation:

### Multidimensional Scaling (MDS)

Two concepts:

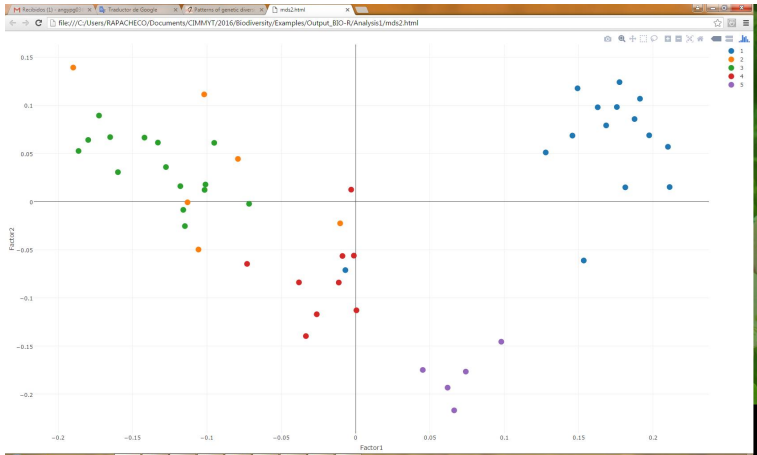
1. **Similarity** ( $s_{[x,y]}$ ):  $d_{[x,y]} = \sqrt{2 \times (1 - s_{[x,y]})}$ ;

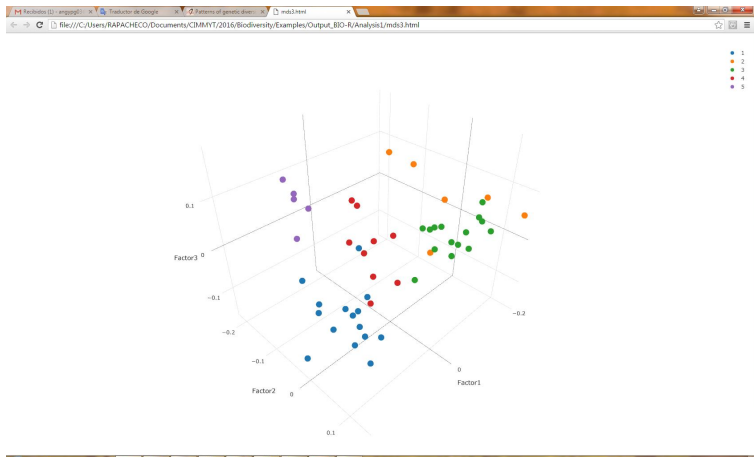
2. **stress** ( $S$ ): 
$$S = \sqrt{\frac{\sum_x \sum_y (d_{[x,y]} - \mu_{d_{[x,y]}})^2}{\frac{n(n-1)}{2}}}$$

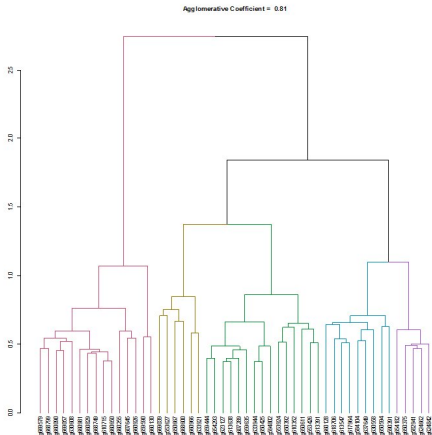


#### NOTE:

Search for the reduced representation (2 or 3 dimensions) that minimizes the  $S.S._B$  in the  $p$  dimensional space and with minimum the estimated distance.





Trees *a.k.a.* Dendograms:

## 5 Software BIO-R and family...

- Phenotypic data

- *Experimental Design*: ADEL, AUDE & STAD;
- *Individual & Multi-Env*: META (> 1500), AGD (> 1300);
- *Interaction/Stability*: GEA (> 1800); and
- *Spatial Analysis*: SPAT.

- Genotypic data

- *kinship*: BROWSE & COP; and
- *Diversity*: BIO.


- Fusion data

- *Relationship/Interaction*: GEA;
- *Genome prediction*: BGLR (\*R package & application tool); and
- *Selection Index*: SI (> 300), RINDSEL.

- **Decision**: Eval  $L \times T$ .



## Download

- BIO-R is coming soon... go to: <http://data.cimmyt.org>;
- Java interface/application for  embedding R [7] scripts to perform Diversity analysis; and
- heterozygosity, diversity B & W groups, shannon index, number of effective allele, % of polymorphic loci, Rogers & Nei distance, clusters and multidimensional scaling (2 & 3d plots).



### CREDITS: BSU/CIMMYT



Angela Pacheco	<R.A.Pacheco@cgiar.org>
Francisco Rodriguez	<F.R.Huerta@cgiar.org>
Gregorio Alvarado	<G.Alvarado@cgiar.org>
Juan Burgueño	<J.Burgueno@cgiar.org>

BIO-R (Biodiversity Analysis with R for Windows)

Open File Help

BSU  
Biodiversity & Conservation

CIMMYT  
International Center for Tropical Agriculture

Welcome to BIO-R (Biodiversity Analysis with R for Windows), Version 1.0 (2016.05.21)  
Copyright © 2016 Centro Internacional de Mejoramiento de Maíz y Trigo (CIMMYT).

Authors:  
Angela Pacheco  
Gregorio Alvarado  
Francisco Rodríguez  
José Chessa  
Juan Burgueño

This program is based in some components from Java, developed by ORACLE AMERICA, INC. and R, developed by R Core Team. Any Java component of this program is hereby licensed under the Oracle Binary Code License Agreement for the Java SE Platform Products and JavaFX made available by Oracle (available at <http://www.oracle.com/technetwork/java/javase/terms/license/index.html>). Any R component of this program as well as the program as a whole developed by CIMMYT are hereby licensed as per the terms of the GNU General Public License version 3 (available at <http://www.gnu.org/licenses/gpl-3.0.html>), as specified below.

This program is free software; you can redistribute it and/or modify it under the terms of the GNU General Public License as published by the Free Software Foundation, either version 3 of the License, or any later.

This program is distributed in the hope that it will be useful, but WITHOUT ANY WARRANTY, without even the implied warranty of MERCHANTABILITY or FITNESS FOR A PARTICULAR PURPOSE. See the GNU General Public License for more details.

You should have received a copy of the GNU General Public License along with this program; if not, write to the Free Software Foundation, Inc., 51 Franklin Street, Fifth Floor, Boston, MA 02110-1301, USA.

For further information, please contact CIMMYT at [CIMMYT-Knowledge-Center@cgiar.org](mailto:CIMMYT-Knowledge-Center@cgiar.org) or at Km. 45 Carretera México-Veracruz, El Estero, Texcoco, Estado de México, México, C.P. 56237.

We have invested a lot of time and effort in creating BIO-R, please cite it when using it for data analysis.

OS: Windows 7 amd64  
Working Directory: C:\Users\apacheco\Documents\CIMMYT\BIO-R\Biodiversity  
R version: R\_3.2.3

 **Input data** MyData.csv

<i>mark</i>	<i>g1</i>	<i>g2</i>	<i>g3</i>	<i>g4</i>	<i>...</i>	<i>gX</i>
1	0.67	NA	0.57	NA	...	1
2	NA	1	1	1	...	NA
3	1	1	0.52	0.50	...	0.80
4	1	NA	0.50	1	...	NA
5	1	NA	1	NA	...	1
6	0.67	0	0.71	1	...	0.33
7	1	NA	0	1	...	1
8	1	NA	1	1	...	1
9	NA	1	0.60	0.22	...	0.71
⋮	⋮	⋮	⋮	⋮	⋮	⋮
<i>N</i>	1	1	1	1	...	1

## Setting parameters and analysing...











1. *Markers* – selects the column that identify the **markers**;
2. *# Clusters* – type the number of groups to split the population;
3. *Output folder* – type the path of the output folder where results will be saved;
  - It will be created inside the bio-R's Output folder;
  - You can change the name for different sets; and
  - It is necessary to change the name for each analysis.
4. *Genotypes* – selects the columns that identify the **genotypes**;
5. *Distance* – selects the method to calculate **distances**; and
6. *ColorMDSPlot* – can specify \*.csv file containing additional information for colors in MDS plot (see manual).

## Results...



### Output

#### Analysis 1

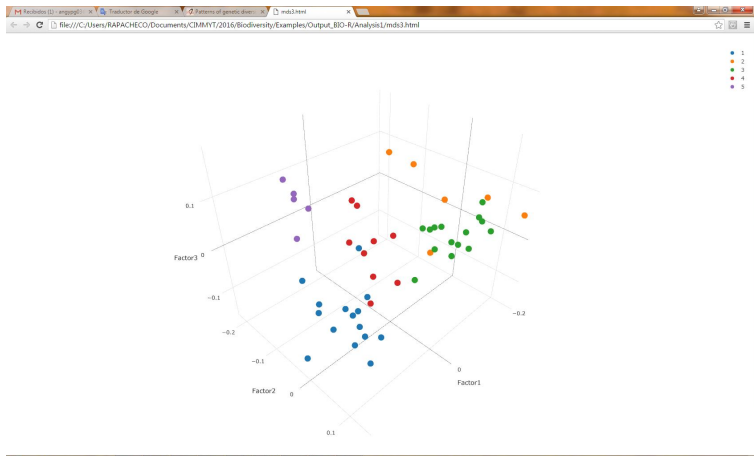
-  CalculusPerGenotype.csv
-  CalculusPerLocus.csv
-  Dendogram.wmf
-  mds2.html
-  mds3.html
-  MDStable.csv
-  RogersDistances.csv
-  SummaryDiversityAnalysis.csv
-  mds2\_files
-  mds3\_files

Output : CalculusPerLocus.csv

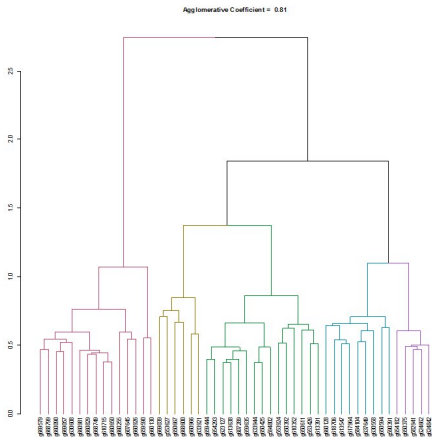
<i>Marker</i>	<i>He</i>	<i>Ho</i>	<i>Ae</i>	<i>Shannon</i>	<i>%NA</i>
1	0.50	0.29	1.99	0.95	0.22
2	0.01	-0.05	1.01	0.02	0.08
3	0.45	0.38	1.82	0.93	0.16
4	0.41	0.42	1.69	0.87	0.24
5	0.03	-0.08	1.03	0.11	0.10
6	0.40	0.43	1.68	0.86	0.20
7	0.50	0.31	2.00	0.99	0.24
8	0.49	0.47	1.94	0.98	0.18
9	0.45	0.47	1.82	0.93	0.14
:	:	:	:	:	:
<i>N</i>	0.08	-0.03	1.09	0.26	0.12

Output : CalculusPerGenotype.csv

<i>Genotype</i>	<i>He</i>	<i>Ho</i>	<i>Ae</i>	<i>Shannon</i>	<i>%NA</i>	<i>clusterGroup</i>
1	0.39	0.06	1.63	0.83	0.11	1
2	0.38	-1.46	1.62	0.82	0.37	2
3	0.40	0.51	1.67	0.85	0.03	3
4	0.39	0.03	1.64	0.83	0.19	4
5	0.37	0.39	1.60	0.81	0.03	3
6	0.41	-3.56	1.69	0.86	0.42	1
7	0.39	0.34	1.65	0.84	0.04	3
8	0.37	-0.28	1.54	0.81	0.24	5
9	0.40	-0.70	1.61	0.82	0.32	1
:	:	:	:	:	:	:
X	0.38	0.33	1.66	0.84	0.04	3





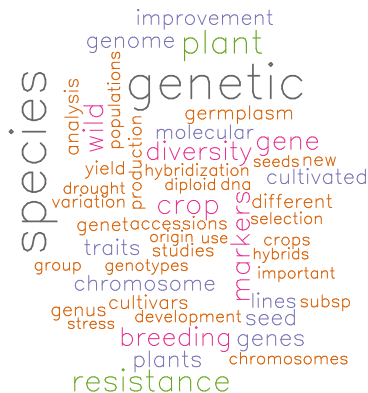


**Summary** : SummaryDiversityAnalysis.csv

- % of polymorphic loci: 0.94;
- Exp. Heterozygosity: 0.30;
- Std. Dev. of *He*: 0.01;
- Obs. Heterozygosity: 0.22;
- Std. Dev. of *Ho*: 0.01;
- Number of effective alleles: 1.55;
- Std. Dev. of *Ae*: 0.02;
- Shannon Index: 0.63;
- Std. Dev. Shannon: 0.02 ...

## 📣 Remarks

- Frequency of Alleles & Genotypes;  
(Population Genetics) The fundamental concept!
- Distances & Diversity Indices;  
(Genetic Diversity) Several ways and flavors.
- Wright Statistics;  
(Intraspecific Diversity) Between/Within variability
- BIO-R; and  
(Software) it has been built considering everything you may want



🗨️ Questions/Suggestions?? 🗨️

 **For Further Studies:**

- [1] E. E. Berg and J. L. Hamrick, "Quantification of genetic diversity at allozyme loci," *Canadian Journal of Forest Research*, vol. 27, pp. 415–429, 1997.
- [2] S. A. Mohammadi and B. M. Prasana, "Analysis of genetic diversity in crop plants – salient statistical tools and considerations," *Crop Science*, vol. 43, pp. 1235–1248, 2003.
- [3] J. C. Reif, A. E. Melchinger, and M. Frish, "Genetical and mathematical properties of similarity and dissimilarity coefficients applied in plant breeding and seed bank management," *Crop Science*, vol. 45, pp. 1–7, 2005.
- [4] N. Saitou and M. Nei, "The neighbour-joining method: A new method for reconstructing phylogenetic trees," *Molecular Biology and Evolution*, vol. 4, no. 4, pp. 406–425, 1987.
- [5] L. Kaufman and P. J. Rousseeuw, *Finding groups in data: an introduction to cluster analysis*. New York: John Wiley, 1990.
- [6] J. K. Pritchard, M. Stephens, and P. Donnelly, "Inference of population structure using multilocus genotype data," *Genetics*, pp. 945–959, 2000.
- [7] R Core Team, *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria, 2016.

That's all folks!

---

✉ <[f.toledo@cgiar.org](mailto:f.toledo@cgiar.org)>

