

caña

CIDCA

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

Mejoramiento Genético de la caña de azúcar mediante el uso de marcadores moleculares

CIDCA, A.C.

Noviembre, 2017.

caña

CIDCA

Datos
genotípicos

Datos fenotípicos

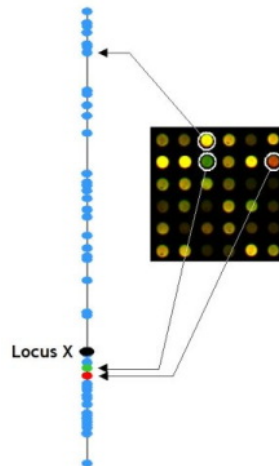
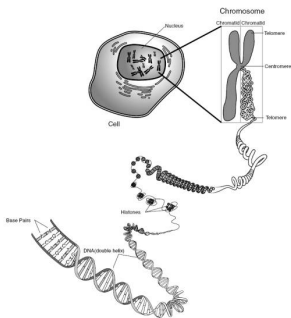
Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis



- Marcadores del tipo “presencia” (1) y ausencia (0).
- **52,828** marcadores obtenidos.

caña

CIDCA

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

R Data Editor

row.names	col1	col2	col3	col4	col5	col6	col7	col8	col9
ATEMex 96-40	0	1	1	1	1	0	0	0	
B 35187	1	0	1	1	1	0	1	0	
B 45181	0	1	0	1	1	0	1	1	
C 323-68	1	0	0	1	1	0	1	1	
CB 36-14	1	0	1	0	0	1	0	0	
CB 40-77	1	1	1	1	0	0	1	1	
CC 83-29	0	0	0	0	1	0	0	0	
CC 84-59	1	0	0	0	0	1	0	0	
CC 85-92	0	1	1	0	0	1	0	0	
CC 86-29	0	1	1	0	0	0	0	0	
CC 87-505	0	0	1	1	1	1	0	0	
CC 92-2198	0	0	0	1	1	1	1	0	
CC 92-2358	0	0	1	0	0	1	1	1	
CC 93-3423	1	0	0	0	1	0	1	0	
CC 93-3458	0	1	0	0	0	1	0	1	

Figura 4: Sub conjunto de marcadores.



caña

CIDCA

Contenido

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

1 Datos genotípicos**2** Datos fenotípicos**3** **Análisis de datos**
Control de calidad
Diversidad
Predicción**4** Otros análisis



caña

CIDCA

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

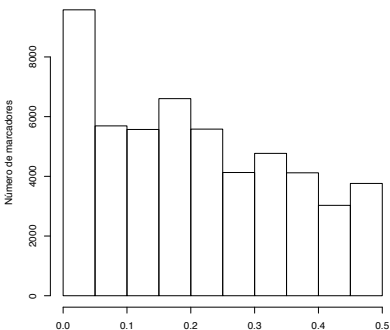


Figura 7: Distribución de frecuencias del alelo menor (MAF).

\hat{p}_j , proporción de 1's para marcador $j = 1, 2, \dots$,

$$maf_j = \begin{cases} \hat{p}_j & \text{si } \hat{p}_j < 0.5 \\ 1 - \hat{p}_j & \text{en caso contrario.} \end{cases}$$

caña

CIDCA

Continuar...

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

En R:

```
Z=scale(X,center=TRUE,scale=TRUE)
```

```
G=tcrossprod(Z)/ncol(Z)
```

```
heatmap(G,cexCol=0.5,cexRow=0.5)
```

Distancias

Una vez que se tienen los marcadores, es posible calcular diferentes tipos de distancias, por ejemplo Euclideana, Rogers, etc.

En el caso de distancias Euclidianas para dos individuos cualesquiera, i y $j = 1, \dots, 93$. con $x_{ij} \in \{0, 1\}$,

$$d_{ij} = \sqrt{\|\mathbf{x}_i - \mathbf{x}_j\|^2} = \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2}$$

para Rogers,

$$d_{ij} = \frac{1}{p} \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2 / 2}$$

caña

CIDCA

Notas para distancia Euclideana

- El valor mínimo de d_{ij} es 0.
- El valor máximo de d_{ij} es \sqrt{p} , donde p es el número de marcadores.

The screenshot shows an Excel spreadsheet with the following data and formulas:

	A	B	C	D	E	F	G
1							
2							
3	Individuo/Marcador	M1	M2	M3	M4	M5	
4		1	1	1	0	0	1
5		2	1	1	0	0	1
6		d12=sqrt((1-1)^2+(1-1)^2+(0-0)^2+(0-0)^2+(1-1)^2)					
7		d12=sqrt(0+0+0+0)=sqrt(0)=0					
8							
9	Individuo/Marcador	M1	M2	M3	M4	M5	
10		1	1	1	0	0	0
11		2	0	0	1	1	1
12		d12=sqrt((1-0)^2+(1-0)^2+(0-1)^2+(0-1)^2+(0-1)^2)					
13		d12=sqrt(1+1+1+1)=sqrt(5)					
14							
15							

Figura 9: Ejemplo de valores máximos y mínimos de d_{ij} .

caña

CIDCA

Continuar...

En R:

```
D=dist(X,method = "euclidean")
D=as.matrix(D)
write.csv(D,file="Euclideanas.csv")
```

	A	B	C	D	E	F	G
1		ATEMex 96-40	B 35187	B 45181	C 323-68	CB 36-14	CB 40-77
2	ATEMex 96-40	0	121.9918	126.05158	121.11978	120.46576	112.97787
3	B 35187	121.991803	0	123.268	122.87392	121.80312	106.32027
4	B 45181	126.0515767	123.268	0	123.51113	121.35485	122.41323
5	C 323-68	121.1197754	122.87392	123.51113	0	120.87183	122.27837
6	CB 36-14	120.4657628	121.80312	121.35485	120.87183	0	123.22337
7	CB 40-77	112.9778739	106.32027	122.41323	122.27837	123.22337	0

Figura 10: Distancias Euclidianas para 6 variedades.

Agrupamiento jerárquico

La idea es combinar las observaciones (es decir, las filas de X) en grupos relativamente homogéneos llamados conglomerados (clusters). Las observaciones del mismo grupo estarán, en algún sentido, “cerca”.

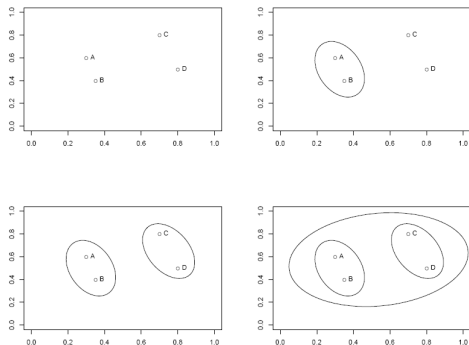


Figura 11: Ejemplo de cluster jerárquico.

Continuar...

En R:

```
D=dist(X)
out_hclust=hclust(D)
plot(out_hclust)
```

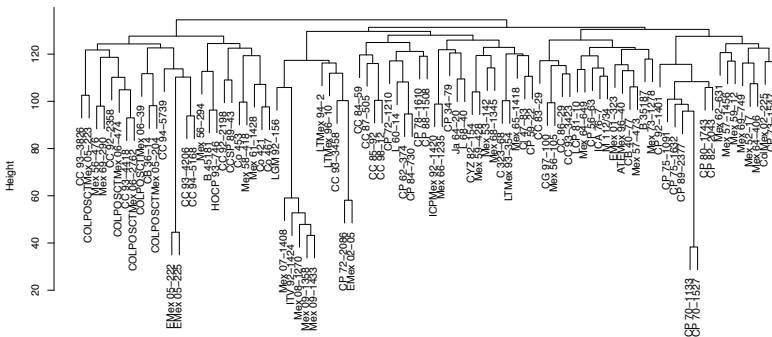


Figura 12: Dendrograma basado en distancias Euclidianas con encadenamiento completo.

caña

CIDCA

Contenido

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

1 Datos genotípicos

2 Datos fenotípicos

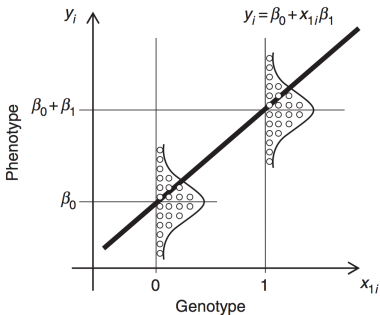
3 **Análisis de datos**
Control de calidad
Diversidad
Predicción

4 Otros análisis

Continuar...

En Selección genómica, los valores genéticos son aproximados utilizando regresión lineal (Meuwissen et al., 2001), esto es:

$$y_i = g_i + e_i = \mu + \sum_{j=1}^p x_{ij}\beta_j + e_i \quad (1)$$



Relación entre marcadores (x_{1i} : 0 and 1) and fenotipos (y_i) para un conjunto de individuos (círculos sin relleno) en una población de entrenamiento (Nakaya e Isobe, 2012).

caña

CIDCA

Continuar...

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

Los valores de cría (GEBVs) se obtienen como una combinación lineal de los marcadores con los pesos dados por los efectos de los mismos, es decir:

$$GEBV_i = \sum_{j=1}^p x_{ij} \hat{\beta}_j \quad (2)$$

GBLUP-RR

El modelo más simple usado en selección genómica:

$$y_i = g_i + e_i = \mu + \sum_{j=1}^p x_{ij}\beta_j + e_i$$

Los efectos de marcadores pueden obtenerse resolviendo el siguiente problema de optimización,

$$\min_{\beta, \lambda} \left\{ (y - \sum X_j\beta_j)'(y - \sum X_j\beta_j) + \lambda \sum \beta_j^2 \right\}, \quad (3)$$

donde $\lambda > 0$ es un parámetro de regularización.

caña

CIDCA

Continuar...

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad


Predicción

Otros análisis

[Theoretical and Applied Genetics](#)October 2013, Volume 126, [Issue 10](#), pp 2575–2586

Experimental assessment of the accuracy of genomic selection in sugarcane

[Authors](#)[Authors and affiliations](#)

M. Gouy, Y. Rousselle, D. Bastianelli, P. Lecomte, L. Bonnal, D. Roques, J.-C. Efile, S. Rocher, J. Daugrois, L. Toubi, S. Nabeneza, C. Hervouet, H. Telismart, M. Denis, A. Thong-Chane, J. C. Glaszmann, J.-Y Hoarau, S. Nibouche, L. Costet , [show less](#)

caña

CIDCA

Continue...

Datos genotípicos

Datos fenotípicos

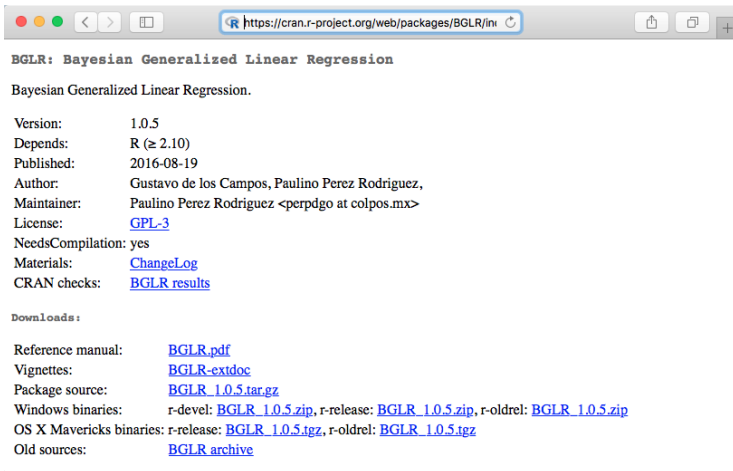
Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis



The screenshot shows a web browser window with the address bar containing the URL <https://cran.r-project.org/web/packages/BGLR/index.html>. The page content is as follows:

BGLR: Bayesian Generalized Linear Regression

Bayesian Generalized Linear Regression.

Version: 1.0.5
Depends: R (≥ 2.10)
Published: 2016-08-19
Author: Gustavo de los Campos, Paulino Perez Rodriguez,
Maintainer: Paulino Perez Rodriguez <perpdgo at colpos.mx>
License: [GPL-3](#)
NeedsCompilation: yes
Materials: [ChangeLog](#)
CRAN checks: [BGLR results](#)

Downloads:

Reference manual: [BGLR.pdf](#)
Vignettes: [BGLR-extdoc](#)
Package source: [BGLR_1.0.5.tar.gz](#)
Windows binaries: r-devel: [BGLR_1.0.5.zip](#), r-release: [BGLR_1.0.5.zip](#), r-oldrel: [BGLR_1.0.5.zip](#)
OS X Mavericks binaries: r-release: [BGLR_1.0.5.tgz](#), r-oldrel: [BGLR_1.0.5.tgz](#)
Old sources: [BGLR archive](#)

BGLR

- A novel software for whole genomic regression and prediction for continuous, discrete traits, censored and uncensored.
- Suitable for big p and small n problems.
- Many non-parametric and parametric models implemented in a consistent manner.
- Large collection of Bayesian models included:
 - Bayesian ridge regression.
 - Bayesian LASSO.
 - BayesA, BayesB, BayesC- π .
 - Reproducing Kernel Hilbert Spaces.
 - Reproducing Kernel Hilbert Spaces with Kernel-Averaging.

caña

CIDCA

En R

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

```
#Clean workspace
rm(list=ls())

library(BGLR)
#Load data
load("data/Pheno_Geno.RData")

#Compute the Genomic Relationship matrix
Z=scale(X,center=TRUE,scale=TRUE)
G=tcrossprod(X)/ncol(Z)

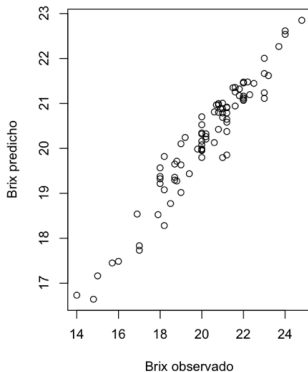
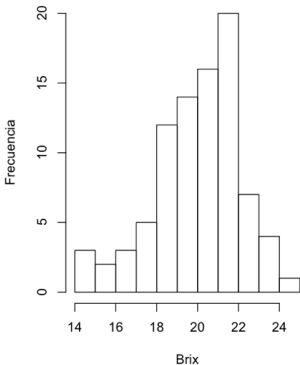
#Exploratory data analysis
hist(Pheno$Brix)

#Regression model
ETA0=list(list(X=Z,model="BRR"))
fm0=BGLR(y=Pheno$Brix,ETA=ETA0,nIter=10000)
```

caña

CIDCA

Predicción de grados Brix



caña

CIDCA

Heredabilidad genómica

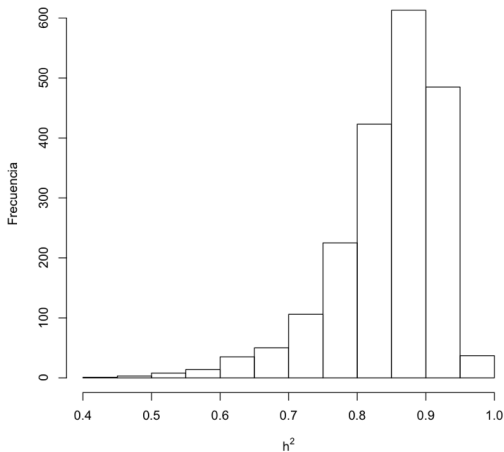


Figura 14: Distribución de $h^2 = \frac{\sigma_a^2}{\sigma_a^2 + \sigma_e^2}$

Predicción valores fenotípicos de cruzas simples

- Predecir valores fenotípicos de cruzas usando información genotípica de padres.
- Datos:
 - Información fenotípica de algunas cruzas.
 - Información genotípica de padres y madres.
- El modelo debe predecir cruzas que no se han hecho y no se han probado en campo.

caña

CIDCA

Lo que se hace en trigo y maíz...

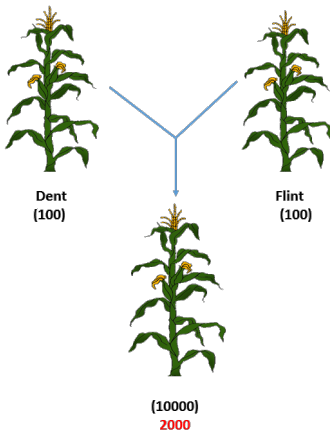


Figura 15: Cruza simple

caña

CIDCA

Modelos

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

HIGHLIGHTED ARTICLE
GENOMIC SELECTION

Genome Properties and Prospects of Genomic Prediction of Hybrid Performance in a Breeding Program of Maize

Frank Technow,* Tobias A. Schrag,* Wolfgang Schipprack,* Eva Bauer,¹
Henner Simianer,² and Albrecht E. Melchinger*¹

*Institute of Plant Breeding, Seed Sciences, and Population Genetics, University of Hohenheim, 70599 Stuttgart, Germany, ¹Plant Breeding, Technische Universität München, 85354 Freising, Germany, and ²Department of Animal Sciences, Georg-August-University Goettingen, 37075 Goettingen, Germany

caña

CIDCA

Continuar...

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

$$\mathbf{y} = \mathbf{Z}_E \boldsymbol{\beta}_E + \mathbf{Z}_D \mathbf{g}_D + \mathbf{Z}_F \mathbf{g}_F + \mathbf{Z}_H \mathbf{h} + \mathbf{e}$$

- \mathbf{Z}_E matriz de incidencia para ambientes, $\boldsymbol{\beta}_E$ el efecto de ambientes.
- \mathbf{Z}_D matriz de incidencia para machos.
- \mathbf{g}_D vector de efectos aleatorios para la aptitud combinatoria general de machos, $\mathbf{g}_D \sim N(\mathbf{0}, \sigma_D^2 \mathbf{G}_D)$.
- \mathbf{Z}_F matriz de incidencias para hembras.
- \mathbf{g}_F vector de efectos aleatorios para aptitud combinatoria general para hembras, $\mathbf{g}_F \sim N(\mathbf{0}, \sigma_F^2 \mathbf{G}_F)$.



caña

CIDCA

Continue...

Datos
genotípicos

Datos fenotípicos

Análisis de datos

Control de calidad

Diversidad

Predicción

Otros análisis

- \mathbf{Z}_H matriz de incidencias para híbridos.
- \mathbf{h}_F vector de efectos aleatorios asociada a aptitud combinatoria específica de las cruzas, $\mathbf{h} \sim N(\mathbf{0}, \sigma_H^2 \mathbf{H})$.
- $\mathbf{H} = \mathbf{G}_D \otimes \mathbf{G}_H$.

Modelo de umbrales para resistencia...

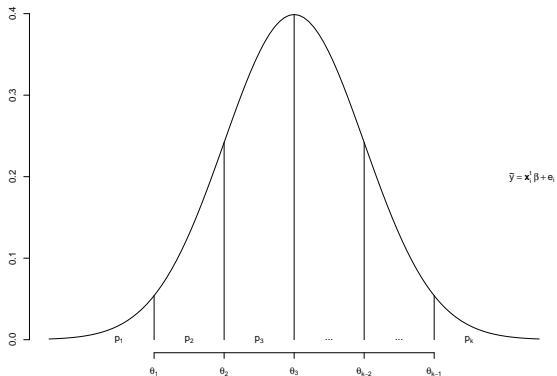


Figura 16: Función de densidad de la variable latente

