



FACULTAD DE CIENCIAS AGRARIAS
UNIVERSIDAD NACIONAL DE ROSARIO

Paquete de R y aplicación Web para el análisis de datos
provenientes de ensayos multiambientales

JULIA ANGELINI

TRABAJO FINAL PARA OPTAR AL TITULO DE ESPECIALISTA EN
BIOINFORMÁTICA

DIRECTOR: Gerardo Cervigni
CO-DIRECTOR: Marcos Prunello

AÑO: 2019

Paquete de R y aplicación Web para el análisis de datos provenientes de ensayos multiambientales

Julia Angelini

Licenciada en Estadística – Universidad Nacional de Rosario

Este Trabajo Final es presentado como parte de los requisitos para optar al grado académico de Especialista en **Bioinformática**, de la Universidad Nacional de Rosario y no ha sido previamente presentada para la obtención de otro título en ésta u otra Universidad. El mismo contiene los resultados obtenidos en investigaciones llevadas a cabo en el **Centro de Estudios Fotosintéticos y Bioquímicos (CEFOBI)**, durante el período comprendido entre los años **2017 y 2019**, bajo la dirección del **Dr. Gerardo Cervigni** y **Mgs. Marcos Prunello**.

Nombre y firma del autor

Nombre y firma del Director

Nombre y firma del Co - Director

Defendida: _____ de 20____.

Agradecimientos

En este trabajo final, directa o indirectamente, participaron muchas personas a las que les quiero agradecer.

En primer lugar al Dr. Gerardo Cervigni por confiar en mí y permitirme explorar el mundo de la Bioinformática durante mi tesis doctoral, para que hoy sea parte de mis conocimientos. Al Mgs. Marcos Prunello por acompañarme en el desarrollo del trabajo final, por su dedicación y sus consejos.

Todo esto nunca hubiera sido posible sin el apoyo y el cariño de mis padres, de mi hermano, de Segundo y Kalita. Siempre estuvieron a mi lado, las palabras nunca serán suficientes para agradecerles.

A mis compañeras Jor y Lu, por su ayuda y por compartir excelentes momentos.

A Gaby y Euge mis compañeras de CEFOBI, gracias a ustedes este camino ha sido mas fácil!

A mis amigos, por estar siempre presentes.

Muchas gracias a todos!

Abreviaturas y Símbolos

EMA: ensayos multiambientales **IGA:** interacción genotipo ambiente **NCOI:** interacción sin cambio de rango, del inglés *no crossover interaction* **COI:** interacción con cambio de rango, del inglés *crossover interaction* **ANOVA:** análisis de la variancia, del inglés *analysis of variance* **AMMI:** modelo de los efectos principales aditivos y interacción multiplicativa, del inglés *Additive Main effects and Multiplicative Interaction* **ACP:** análisis de componentes principales **SREG:** modelo de regresión por sitio, del inglés *Site Regression model* **DVS:** descomposición de valores singulares **GNU:** *General Public Licence* **CRAN:** *Comprehensive R Archve Network* **EM:** maximización de la esperanza, del inglés *Expectation-Maximization*

Resumen

Los programas informáticos se han convertido, hoy en día, en una herramienta básica utilizada por el análisis estadístico como apoyo fundamental a la hora de realizar diferentes operaciones y para facilitar una mayor comodidad a los usuarios. Actualmente, R es uno de los programas más utilizados debido a su potencia y a su distribución como software libre.

Palabras Clave:

Abstract

Keywords:

Índice general

Capítulos	Página
1. Introducción	1
2. Objetivos	7
2.1. Objetivo general	7
2.2. Objetivos específicos	7
3. Métodos	8
3.1. Métodos estadísticos	8
3.1.1. Modelo AMMI	8
3.1.2. Modelo SREG	11
3.1.3. Métodos de imputación	17
3.2. Paquete de R	18
3.2.1. Esqueleto y estructura del paquete	19
3.2.2. Creación de funciones y conjuntos de datos	20
3.2.3. Documentación	21
3.2.4. Pruebas del flujo de trabajo	22
3.2.5. Compilación e instalación	23
3.2.6. Publicación	24
3.3. Shiny APP	24
3.3.1. Flujo de trabajo	26
3.3.2. Compartiendo una Shiny Web App	29
4. Resultados	30
4.1. Paquete de R <i>geneticae</i>	30
4.1.1. Conjuntos de datos en <i>geneticae</i>	30
4.1.2. Funciones en <i>geneticae</i>	31
4.2. <i>Geneticae</i> Shiny Web App	48
4.2.1. Los datos	49

4.2.2.	Análisis descriptivo	50
4.2.2.1.	<i>Boxplot</i>	50
4.2.2.2.	Gráfico de correlación	52
4.2.2.3.	Matriz de correlación	54
4.2.2.4.	Gráfico de interacción	56
4.2.3.	Análisis de la variancia	58
4.2.4.	Biplot GGE	58
4.2.5.	Biplot GE	58
4.2.6.	Ayuda	58
5.	Conclusiones	59
A.	Hoja de referencia Shiny	60
B.	Guías para usuario de Geneticae APP	63
C.	Código R de Geneticae APP	64
	Bibliografía	65

Índice de figuras

1.1. Representación gráfica de tipos de IGA: (A)IGA no crossover, (B) IGA crossover y (C) no IGA	3
3.1. Ejemplo de un biplot GE	10
3.2. Ranking de genotipos en el ambiente D a través del biplot GGE	13
3.3. Ambientes favorables y desfavorables para el genotipo 2 en el biplot GGE .	13
3.4. Comparación de los genotipos 6 y 8 en el biplot GGE	14
3.5. Biplot GGE con el polígono envolvente y las perpendiculares a sus lados .	15
3.6. Eje de coordenadas de ambiente medio para un mega-ambiente en el biplot GGE	17
3.7. Esquema interno de la aplicación.	25
4.1. Biplot básico obtenido de la función <code>GGEPlot()</code>	33
4.2. Ranking de cultivares para un ambiente determinado obtenido de la función <code>GGEPlot()</code>	33
4.3. Ranking de ambientes para cultivar determinado obtenido de la función <code>GGEPlot()</code>	34
4.4. Relación entre ambientes obtenido de la función <code>GGEPlot()</code>	35
4.5. Comparación entre dos genotipos obtenido de la función <code>GGEPlot()</code>	35
4.6. Identificación del mejor cultivar en cada ambiente a partir de la función <code>GGEPlot()</code>	36
4.7. Evaluación de los ambientes basados tanto en la capacidad de discriminación y representatividad a partir de la función <code>GGEPlot()</code>	37
4.8. Clasificación de ambientes con respecto al ambiente ideal a partir de la función <code>GGEPlot()</code>	37
4.9. Clasificación de genotipos con respecto al genotipo ideal a partir de la función <code>GGEPlot()</code>	38
4.10. Evaluación de los cultivares con base en el rendimiento promedio y la estabilidad a partir de la función <code>GGEPlot()</code>	39
4.11. Biplot GE obtenido del modelo clasico AMMI	40

4.12. Biplot GE obtenido del modelo robusto rAMMI	41
4.13. Biplot GE obtenido del modelo robusto hAMMI	41
4.14. Biplot GE obtenido del modelo robusto gAMMI	42
4.15. Biplot GE obtenido del modelo robusto lAMMI	43
4.16. Biplot GE obtenido del modelo robusto ppAMMI	43
4.17. yan.winterwheat dataset disponible en Shiny Web App	49
4.18. plrv dataset disponible en Shiny Web App	50
4.19. Boxplot de ambientes a través de los genotipos para el conjunto de datos Plrv	51
4.20. Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv	52
4.21. Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv	53
4.22. Boxplot de ambientes a través de los genotipos para el conjunto de datos Plrv	54
4.23. Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv	55
4.24. Boxplot de ambientes a través de los genotipos para el conjunto de datos Plrv	56
4.25. Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv	57
4.26. Boxplot de ambientes a través de los genotipos para el conjunto de datos Plrv	58

Índice de tablas

Lista de tareas pendientes

This is a note	see 1.0 at p. 1
This is a note	see 1.0 at p. 1

Capítulo 1

Introducción

A lo largo de la historia de la agricultura, el hombre ha desarrollado el mejoramiento vegetal (o fitomejoramiento) en forma sistemática y lo ha convertido en un instrumento esencial para la mejora de la producción agrícola en términos de cantidad, calidad y diversidad.

This is a note

La humanidad depende, directa o indirectamente, de las plantas. Para la alimentación, ya que todos sus alimentos son vegetales o se derivan de éstos por ejemplo: carne, huevos y productos lácteos. De las plantas se deriva también la mayoría de las fibras textiles, fármacos, combustibles, lubricantes y materiales de construcción.

This is a note

El fitomejoramiento, en un sentido amplio, es el arte y la ciencia de alterar o modificar la herencia de las plantas para obtener cultivares (variedades o híbridos) mejorados genéticamente, adaptados a condiciones específicas, de mayores rendimientos económicos y de mejor calidad que las variedades nativas o criollas (Allard, 1967). En otras palabras, el fitomejoramiento busca crear plantas cuyo patrimonio hereditario esté de acuerdo con las condiciones, necesidades y recursos de los productores rurales, de la industria y de los consumidores, o sea de todos aquellos que producen, transforman y consumen productos vegetales.

Las variedades mejoradas son el resultado del trabajo de desarrollo genético llevado a cabo en los programas de fitomejoramiento, los cuales se extienden a lo largo de varios años y requieren cuantiosas inversiones. La vigencia comercial de las variedades puede extenderse durante varias décadas, por lo que su elección es crítica para que el productor evite pérdidas económicas por malas campañas y el suministro al mercado sea constante.

Generalmente, en etapas tempranas de estos programas existen un gran número de genotipos experimentales con pocos antecedentes de evaluación; mientras que en etapas posteriores se trabaja con pocos genotipos altamente selectos.



La utilización adecuada de procedimientos de análisis de datos agronómicos y ambientales es una condición inherente al desarrollo actual y futuro de investigaciones orienta-



das a mejorar los cultivos en forma económica y ambientalmente sustentable. En etapas avanzadas de los programas de mejoramiento, los ensayos multiambientales (EMA) que comprenden experimentos en múltiples ambientes son herramientas fundamentales para incrementar la productividad y rentabilidad de los cultivos. Estos son frecuentes en investigaciones agrícolas de comparación de rendimiento, ya que constituyen una de las principales estrategias de identificación de genotipos vegetales superiores y de ambientes en los cuales estos se expresan de manera diferencial.

Debido a que las regiones de producción de los principales cultivos cubren áreas ecológicas muy extensas, se observan variaciones en las condiciones climáticas y de suelo. Por lo tanto, la aparición de la interacción genotipo ambiente (IGA) es inevitable, provocando respuestas altamente variables en los diferentes ambientes (Crossa et al., 1990; Cruz Medina, 1992; Kang y Magari, 1996). La IGA es considerada casi unánimemente por los fitomejoradores como el principal factor que limita la respuesta a la selección y, en general, la eficiencia de los programas de mejoramiento.

Los investigadores agrícolas han sido conscientes de las diversas implicaciones de IGA en los programas de mejoramiento (Mooers, 1921; Yates y Cochran, 1938). Por ejemplo, la IGA tiene un impacto negativo en la heredabilidad, cuanto menor sea la heredabilidad de un carácter, mayor será la dificultad para mejorar ese carácter mediante la selección. Por lo tanto, información sobre la estructura y la naturaleza de la IGA es particularmente útil para los mejoradores porque puede ayudar a determinar si necesitan desarrollar cultivares para todos los ambientes de interés o si deberían desarrollar cultivares específicos para ambientes específicos (Bridges, 1989). Gauch y Zobel (1996) explicaron la importancia de IGA como: “Si no hubiera interacción, una sola variedad de trigo (*Triticum aestivum* L.) o maíz (*Zea mays* L.) o cualquier otro cultivo rendiría al máximo en todo el mundo, y además la prueba de variedades debería realizarse en un sólo lugar para proporcionar resultados universales. No habría ruido, los resultados experimentales serían exactos, identificando la mejor variedad sin error, y no habría necesidad de replicación. Entonces, una réplica en un lugar identificaría la mejor variedad de trigo que florece en todo el mundo”.

Peto (1982) ha distinguido las interacciones cuantitativas, llamada también sin cambio de rango (NCOI), o *no crossover*, de las interacciones cualitativas, denominada también con cambio de rango (COI) o *crossover* (Cornelius et al., 1996). Cuando dos genotipos X e Y tienen una respuesta diferencial en dos ambientes diferentes, pero su ordenación permanece sin cambios se dice que la IGA es *no crossover* (Figura 1.1(A)). Sin embargo, es de tipo *crossover* cuando hay cambios en el orden de los genotipos (Figura 1.1(B)). Cuando los genotipos responden de manera similar en ambos ambientes (Figura 1.1(C)) no hay IGA.

Distintos conceptos como regiones ecológicas, ecotipos, mega-ambientes, adaptaciones

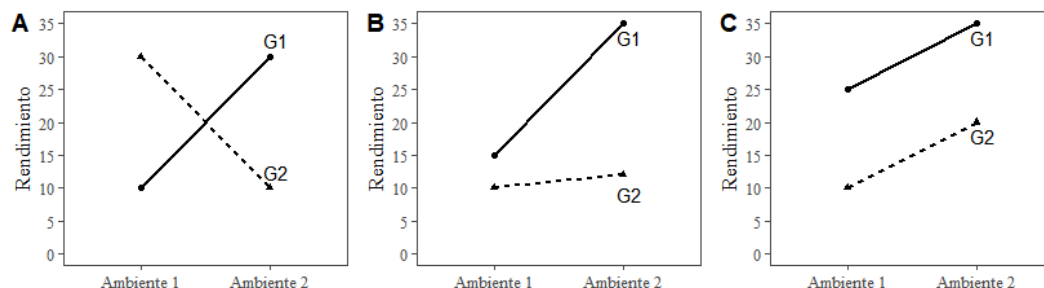


Figura 1.1: Representación gráfica de tipos de IGA: (A) IGA no crossover, (B) IGA crossover y (C) no IGA

de germoplasma tanto en sentido amplio (a través de los ambientes) como específico (para cada ambiente o grupos de ambiente particular) (Kang et al., 2004) se pueden analizar a partir de la **interacción genotipo-ambiente** (Yan y Hunt, 2001).

Un análisis adecuado de la información de los EMA es indispensable para que el programa de mejoramiento de los cultivos sea eficaz. El rendimiento medio en los ambientes es un indicador suficiente del rendimiento genotípico solo en ausencia de IGA (Yan y Kang, 2003). Sin embargo, la aparición de IGA es inevitable y no basta con la comparación de las medias **de los** genotipos, sino que se debe recurrir a una **metodología estadística más apropiada**. Metodología estadística más difundida para analizar los datos provenientes de EMA se basa en modificaciones de los modelos de regresión, análisis de variancia (*Analysis of Variance*, ANOVA) y técnicas de análisis multivariado.

Particularmente, para el estudio de la interacción y los análisis que de ella se derivan, dos modelos multiplicativos han aumentado su popularidad entre los fitomejoradores como una herramienta de análisis gráfico: el modelo de los efectos principales aditivos y **interacción multiplicativa** (*Additive Main effects and Multiplicative Interaction*, AMMI) (Kempton 1984, Gauch, 1988), y el de regresión por sitio (*Site Regression model*, SREG) (Cornelius et al., 1996; Crossa y Cornelius, 1997 y 2002). Estos modelos combinan un ANOVA con la descomposición de valores singulares (DVS) o el análisis de componentes principales (ACP) sobre la matriz residual de ANOVA. En SREG, el ANOVA se realiza sobre el efecto principal de A mientras que en AMMI se considera el efecto de G y A. Mient**a** que a través del modelo AMMI se obtiene el gráfico biplot *Genotipo-Entorno* (**GE**) el cual es usado para explorar patrones puramente atribuibles a los efectos GE, para el modelo SREG, Yan y Hunt (2002), presentaron la técnica GGE biplot usado **para explorar simultáneamente patrones** variación en la suma G+GE.

Una limitación importante de la mayoría de las propuestas de análisis provenientes de EMA es que requieren que el conjunto de datos este completo. Aunque los EMA

están diseñados para que todos los genotipos se evalúen en todos los ambientes, las tablas de datos **genotipo x ambiente** completas son poco frecuentes (no todos los genotipos se encuentran en todos los ambientes). Esto ocurre, **por ejemplo**, debido a errores de medición o causas naturales como, **por ejemplo**, la destrucción de plantas por animales, inundaciones o durante la cosecha, la incorporación de nuevos genotipos y a que otros se descartan por su pobre desempeño (Hill y Rosenberg, 1985). En estos casos, entre las posibles soluciones para tratar con una tabla de datos incompleta es (i) el uso de un subconjunto completo de datos, eliminando aquellos genotipos que tienen valores faltantes (Ceccarelli et al., 2007, Yan et al., 2011), (ii) completar datos faltantes con la media ambiental, o (iii) imputación de datos faltantes con valores estimados utilizando, por ejemplo, un modelo multiplicativo (Kumar et al., 2012).

En este contexto, el análisis de datos provenientes de EMA **requiere metodología estadística sofisticada** cuyas rutinas informáticas se encuentran disponibles en programas desarrollados por diferentes empresas. Esto genera el inconveniente de tener que disponer de todos los programas necesarios para los distintos análisis, atender los requerimientos de formatos de datos usados por cada uno, y comprender los diversos tipos de salidas en las que se ofrecen los resultados obtenidos. Además, algunos procedimientos, especialmente aquellas metodologías recientes, no se encuentran **disponibles**, y los costos de las licencias de dichos programas resultan muy elevados.

El software R se trata de un proyecto de software libre distribuido bajo los términos de la *General Public Licence* (GNU), desarrollado por *The R Foundation for Statistical Computing*. Surge como resultado de la implementación de uno de los lenguajes más utilizados en investigación por la comunidad estadística, el lenguaje S. **Que un software sea libre quiere decir que sus usuarios son libres de usarlo, copiarlo, distribuirlo, editarlo y modificarlo según sus propias inquietudes (FSF 2019)**. A diferencia de los programas estadísticos utilizados frecuentemente, R es un lenguaje de programación y no dispone de una interfaz gráfica en la cual se utilizan **menús** para realizar los distintos análisis de interés, lo cual genera dificultad en su uso **para aquellos que no se encuentran familiarizados con el uso de código computacional**. Sin embargo, brinda mayores posibilidades en cuanto a la manipulación y análisis de los datos **ya que le permite a los usuarios definir sus propias funciones y personalizar el tipo de análisis que desean realizar**. Si bien la versión básica del programa dista mucho de ser amigable, R Studio permite contar con una interacción más fluida con el programa **actuando como una interfaz amigable entre el usuario**. RStudio es un entorno de desarrollo integrado (IDE) gratuito y de código abierto para R.

R forma parte de un proyecto colaborativo ya que promueve el hecho de que los **usuarios creen funciones y las ponga al alcance de toda la comunidad**, es decir que está en continuo

desarrollo y actualización. Sin embargo, como muchas veces no resulta sencillo reutilizar una función creada por algún usuario se ha introducido la posibilidad de crear paquetes (*package*) o librerías. Estas son una colección de objetos creados y organizados siguiendo un protocolo fijo que garantiza un soporte mínimo para el usuario así como la ausencia de errores (de sintaxis) en la programación.

R cuenta con 14 paquetes básicos y 29 recomendados para su funcionamiento instalados automáticamente en él, como por ejemplo, *base* o *stats*. Dado que la comunidad de usuarios que programan en R ha ido creciendo notablemente en los últimos años y que muchos de ellos han ido proporcionando librerías, se cuenta con una gran cantidad de paquetes que extienden las funciones básicas de R. Entre ellos se encuentran, *plyr*, *lubridate*, *reshape2* y *stringr* para la manipulación de los datos; *ggplot2* y *rgl* para la visualización; *knitr* y *xtable* para la presentación de resultados; entre otros. La lista completa de los paquetes oficiales puede consultarse en CRAN¹, se contaba con más de 14.000 paquetes disponibles en CRAN hasta junio de 2019. Esta gran variedad de paquetes es una de las razones por las que R tiene tanto éxito: lo más probable es que alguien ya haya resuelto un problema en el que estás trabajando, y puedes beneficiarte de su trabajo descargando su paquete.

Además de los paquetes oficiales, existen otros que pueden instalarse desde repositorios como, por ejemplo, Github. Sin embargo, no es sencillo encontrar un paquete que puede ser útil para un determinado fin sino que se debe recurrir a varios de ellos para cumplir un determinado objetivo.

Hoy los programas de computación se han convertido en una herramienta básica para el análisis de datos. La decisión de qué software de análisis estadístico utilizar no tiene una respuesta predefinida: la elección dependerá de las necesidades de la investigación. Esto pues los lenguajes de programación son herramientas y el principal criterio para decidir el uso de uno u otro debe efectuarse en función de la particularidad de los objetivos y alcances de la investigación que se busque desarrollar. Actualmente, R es uno de los programas más utilizados debido a su poder y distribución gratuita. R tiene una sintaxis compleja y, por lo tanto, no es del todo amigable para aquellos que no tienen conocimiento del lenguaje de programación R.

Con frecuencia, los mejoradores usan programas que tienen una interfaz para realizar el análisis estadístico deseado y sin necesidad del manejo de un lenguaje de programación. En el año 2012 se creó el paquete *Shiny* de R que permite desarrollar aplicaciones Web utilizando R, acercando la potencia de R a todo tipo de usuarios.

El objetivo del presente trabajo fue: (i) crear un paquete de R que incluya las funciones que permitan analizar los datos provenientes de EMA, incluyendo además metodología re-

¹CRAN (Comprehensive R Archive Network) es el repositorio oficial de paquetes de R, el lugar donde se publican las nuevas versiones del programa, etc. Contiene la lista completa de paquetes oficiales. https://cran.r-project.org/web/packages/available_packages_by_name.html

cientemente publicada que no se encuentra disponible en R; (ii) crear una interfaz gráfica, entre R y el usuario, mediante Shiny con el fin de poder realizar los análisis disponibles en el paquete creado sin necesidad de utilizar el lenguaje de programación.

Elousa, Paula. 2009. “¿EXISTE VIDA MÁS ALLÁ DEL SPSS? DESCUBRE R.” *Revista Psicothema* 21 (4): 652–55. <http://www.psicothema.com/psicothema.asp?id=3686>.

FSF. 2019. “¿Qué Es El Software Libre?” Free Software Foundation. <https://www.fsf.org/es/recursos-es-el-software-libre>.

Capítulo 2

Objetivos

2.1. Objetivo general

Construir un paquete para el programa R, con funciones estadísticas que permitan analizar datos provenientes de EMA, tanto para Windows como Linux, y que cumpla con los estándares de calidad de software. Por otro lado, desarrollar una Shiny Web APP que permita a los usuarios utilizar las funciones del paquete sin tener que utilizar un lenguaje de programación.

2.2. Objetivos específicos

1. Mostrar un flujo de trabajo reproducible para la construcción de paquetes de R.
2. Modificar funciones existentes para el análisis de datos provenientes de EMA de manera que sean más flexibles que las actuales.
3. Incorporar metodología recientemente publicada en el paquete de R.
4. Desarrollar una Shiny Web APP para analizar los datos provenientes de EMA, presentando un flujo de trabajo reproducible.

Capítulo 3

Métodos

3.1. Métodos estadísticos

3.1.1. Modelo AMMI

Modelo AMMI clásico

El modelo AMMI es un modelo multiplicativo en el cual se expresa el fenotipo de un genotipo en un ambiente de la siguiente forma:

$$y_{ij} = \mu + G_i + A_j + \sum_{k=1}^q \lambda_k \alpha_{ik} \gamma_{jk} \quad i = 1, \dots, g; j = 1, \dots, a \quad q = \min(g - 1, a - 1)$$

donde

- y_{ij} es el caracter fenotípico evaluado (rendimiento o cualquier otro caracter de interés) del i -ésimo genotipo en el j -ésimo ambiente,
- μ es la media general,
- G_i es el efecto del i -ésimo genotipo,
- A_j es el efecto del j -ésimo ambiente
- $\sum_{k=1}^q \lambda_k \alpha_{ik} \gamma_{jk}$ es la sumatoria de componentes multiplicativas utilizadas para modelar la IGA. Siendo, λ_k el valor singular para la k -ésima componente principal (PC) α_{ik} y γ_{jk} son los scores de las PC para el i -ésimo genotipo y el j -ésima ambiente para la k -ésima componente, respectivamente;

Los parámetros de IGA en el modelo AMMI se estiman por medio de la DVS de la matriz que contiene los residuos del modelo aditivo luego de ajustar por mínimos cuadrados

el modelo de efectos principales. Generalmente los dos primeros términos multiplicativos son suficientes para explicar los patrones de interacción; la variabilidad remanente se interpreta como ruido.

Los patrones de interacción se pueden visualizar mediante los biplots GE. El concepto del biplot fue presentado por Gabriel (1971), que consiste en una representación de las filas (individuos) y las columnas (variables) de una matriz de datos en un mismo gráfico.

Biplot GE

El biplot GE ayuda a interpretar la variación producida por los efectos de la IGA. Se grafican en un sistema de coordenadas cartesianas de dos dimensiones los scores de los genotipos (α_{ik}) y los ambientes (γ_{jk}), ponderados por la raíz cuadrada del autovalor correspondiente (λ_k).

Dado que los genotipos y los ambientes son definidos como vectores desde el origen (0,0) hasta sus scores, el biplot se interpreta en términos de las direcciones de vectores y sus proyecciones.

El módulo del vector de un ambiente indica la contribución del mismo a la interacción. Los puntos de los genotipos que se encuentran próximos al origen indican que los mismos contribuyen poco a la interacción, es decir, se adaptan de igual manera a todos los ambientes. Puntos cercanos entre sí indican patrones de interacción similares, mientras que puntos alejados entre sí tienen patrones diferentes. Ángulos de $< 90^\circ$ ó $> 270^\circ$ entre los ambientes y genotipos indica que contribuyen positivamente a la interacción (hay una asociación positiva entre ese genotipo y ese ambiente); y mientras más alejados del origen se encuentre los marcadores, más fuerte será esa asociación. Una asociación fuerte y positiva indica que ese ambiente es muy favorable para ese genotipo. De manera similar, cuando un genotipo y un ambiente forman un ángulo entre 90° y 270° se interpreta que ese ambiente es muy desfavorable para ese genotipo.

En la Figura 3.1, se presenta un ejemplo de un biplot GE con 6 ambientes (A, B, C, D, E y F) y 10 genotipos (1, 2, 3, 4, 5, 6, 7, 8, 9 y 10). Se observa los ambientes A y E son los que más contribuyen a la interacción. La cercanía de los genotipos 1 y 2 indica que esos genotipos tienen patrones de interacción similares, y a la vez, muy distintos a los del genotipo 4. Las cercanías entre el genotipo 9 y el ambiente A, entre el genotipo 5 con el ambiente C, entre los genotipos 1 y 2 con el ambiente E, entre los genotipos 6, 7 y 10 con el ambiente F y entre el genotipo 4 con el ambiente D, lo que indica, y la gran distancia con el origen de coordenadas, indica que esos ambientes son muy favorables para esos genotipos. Por otro lado, el ambiente A es considerablemente desfavorable para el genotipo 10. También se observa que los genotipos 3 y 8 se adaptan en igual medida a

todos los ambientes, debido a su proximidad con el origen de coordenadas.

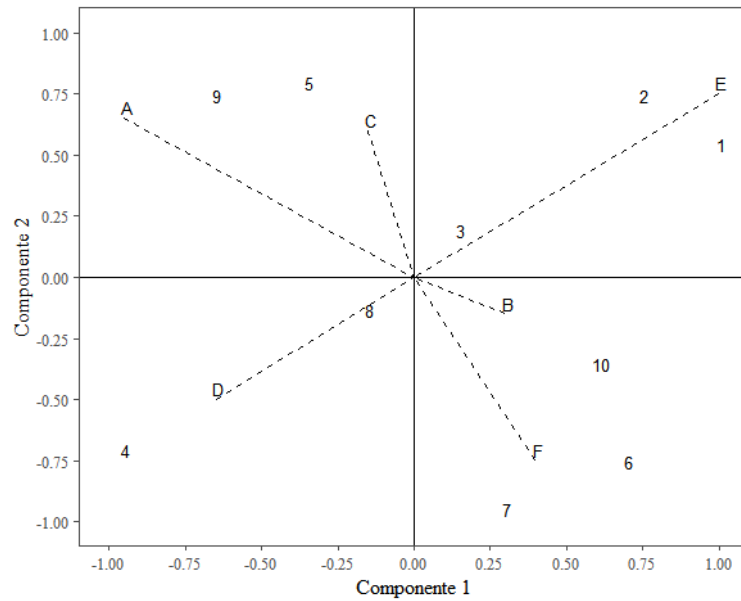


Figura 3.1: Ejemplo de un biplot GE

Modelo AMMI robusto

El modelo AMMI, en su forma estándar, asume que no hay valores atípicos en el conjunto de datos. La presencia de observaciones atípicas es más una regla que una excepción cuando se consideran datos agronómicos debido a errores de medición, algunas plagas / enfermedad que puede influir en algunos genotipos en un ambiente resultando por ejemplo en un rendimiento inferior al esperado, o incluso debido a alguna característica inherente de los genotipos que se evalúan.

Rodrigues et al. (2015) proponen una generalización robusta del modelo AMMI, que se puede obtener en dos etapas: primero ajustar la regresión robusta basada en el estimador M-Huber (Hub) para reemplazar el ANOVA; y luego utilizar un procedimiento DVS / PCA robusto para reemplazar la DVS estándar. En la segunda etapa, consideraron varios métodos dando lugar a total de cinco robustos llamados: R-AMMI, H-AMMI, G-AMMI, L-AMMI, PP-AMMI.

El empleo de la versión robusta del modelo AMMI puede ser extremadamente útil debido a que una mala representación de genotipos y ambientes en los biplots puede dar como resultado un mala decisión con respecto a qué genotipos seleccionar para un conjunto dado de ambientes (es decir, megaambientes; Gauch y Zobel, 1997; Yan et al., 2000). A su vez, la elección de los genotipos incorrectos pueden provocar grandes pérdidas en términos de producción de rendimiento. Los biplots obtenidos de los modelos robustos mantienen las características e interpretación estándar del modelo AMMI clásico (Rodrigues et al.

(2015)).

3.1.2. Modelo SREG

El modelo SREG (Cornelius et al., 1996; Crossa y Cornelius, 1997 y 2002) expresa el fenotipo de un genotipo en un ambiente en función del efecto ambiente aditivo y los efectos genotipo e interacción agrupados y en forma multiplicativa. $y_{ij} = \mu + A_j + \sum_{k=1}^q \lambda_k \alpha_{ik} \gamma_{jk}$ $i = 1, \dots, g; j = 1, \dots, a$ $q = \min(g - 1, a - 1)$ donde

- y_{ij} es la característica fenotípica evaluada (rendimiento u otra variable cuantitativa de interés) del i -ésimo genotipo en el j -ésimo ambiente,
- μ es la media general,
- G_i es el efecto del i -ésimo genotipo,
- A_j es el efecto del j -ésimo ambiente
- $\sum_{k=1}^q \lambda_k \alpha_{ik} \gamma_{jk}$ es la sumatoria de componentes multiplicativas utilizadas para modelar los efectos G e IGA en forma conjunta. Siendo, λ_k el valor singular para la k -ésima PC α_{ik} y γ_{jk} son los scores de las PC para el i -ésimo genotipo y el j -ésimo ambiente para la k -ésima PC, respectivamente;

Para visualizar conjuntamente el efecto de G y IGA Yan et al. (2000) proponen los gráficos biplots GGE (Genotipe plus Genotipe-Environment). A partir de estos gráficos se puede investigar la diferenciación de mega-ambientes entre los ambientes en estudio, seleccionar cultivares superiores en un mega-ambiente dado y seleccionar los mejores ambientes de evaluación para analizar las causas de la IGA. Se define como mega-ambiente a un grupo de ambientes en donde los cultivares de mejor desempeño son los mismos.

Biplot GGE

El biplot GGE, ayuda a interpretar conjuntamente la variación producida por los efectos principales de los G + IGA.

Para la construcción de los biplots GGE, al igual que para los biplots GE, se grafican en un sistema de coordenadas cartesianas de dos dimensiones los scores de los genotipos (α_{ik}) y los ambientes (γ_{jk}), ponderados por la raíz cuadrada del autovalor correspondiente (λ_k).

Para una mejor comprensión de las interpretaciones que se pueden extraer del gráfico biplot GGE, se presenta un ejemplo del mismo para un ensayo de 6 ambientes (A, B, C, D, E y F) con 12 genotipos (1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11 y 12).

Biplot básico

Desempeño de los cultivares en un ambiente dado

Para identificar los mejores genotipos en un ambiente a través del biplot GGE, Yan y Hunt (2002) sugieren trazar una recta que pase por el origen y el punto del ambiente de interés, formando un eje para el ambiente. Luego, ranking de los genotipos en ese ambiente se puede obtener con las proyecciones de los marcadores de los genotipos sobre ese eje. La línea que pasa por el origen y es perpendicular al eje del ambiente separa aquellos genotipos con rendimiento superior e inferior al promedio del ambiente. El genotipo de mayor rendimiento en el ambiente es aquel cuya proyección sobre el eje está más alejada del origen hacia el semi-eje donde se encuentra el marcador del ambiente. Aquel cuya proyección sea la segunda más alejada del origen en ese sentido, es el de segundo mejor rendimiento y así hasta llegar al de menor rendimiento en el ambiente, que es aquel cuya proyección está más alejada del origen en sentido contrario al identificador del ambiente.

Como se observa en la Figura 3.2, el genotipo de mayor rendimiento en el ambiente D es el 4, luego le sigue el 7, luego el 8, y así sucesivamente hasta llegar al genotipo 12, que es el de peor rendimiento en ese ambiente. Además, los genotipos 4, 7, 8, 11, 2 y 5 quedan del lado del marcador del ambiente D, de acuerdo a la división de la perpendicular que pasa por el origen, por que se interpreta que estos genotipos tienen un rendimiento mayor al promedio del ambiente D. Los restantes genotipos tienen un rendimiento inferior al promedio.

Adaptación relativa de un cultivar dado en diferentes ambientes

Para identificar qué ambiente es el más adecuado para un cultivar, Yan y Hunt (2002) proponen graficar una línea que una el marcador del genotipo con el origen, formando un eje para el cultivar, y luego trazar otra línea perpendicular a la primera. Esta última perpendicular es la que separa los sitios favorables y desfavorables para el genotipo. Los sitios cuyos marcadores queden en el mismo lado donde está el genotipo son los mejores para ese genotipo y los restantes son aquellos donde el genotipo rinde por debajo de su promedio.

Como se puede apreciar en la Figura 3.3, la perpendicular al marcador del genotipo 2, determina que los ambientes favorables para ese genotipo son A, E, F y D, en donde el mismo tiene rendimientos mayores a su promedio. Los ambientes desfavorables son B y

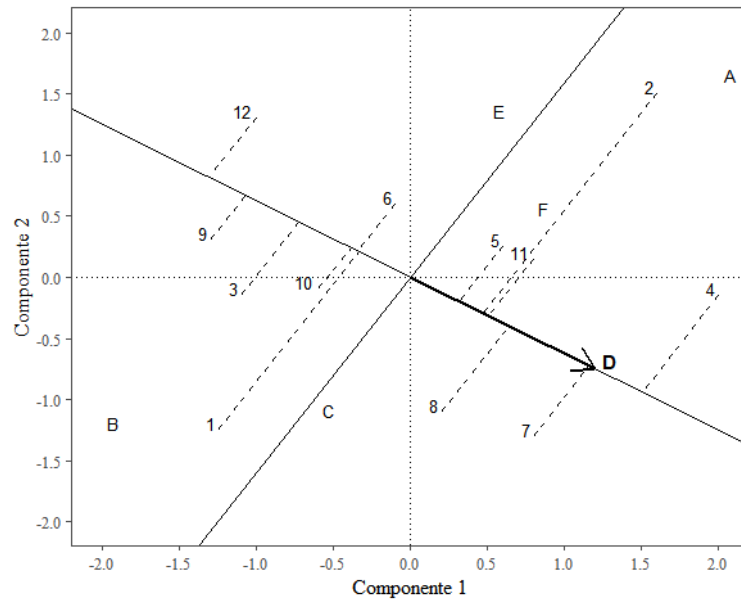


Figura 3.2: Ranking de genotipos en el ambiente D a través del biplot GGE

C. También se observa que el ambiente más favorable es el A, luego le siguen el E y el F. Si bien el ambiente D también es favorable, en ese ambiente el rendimiento del genotipo 2 es apenas superior a su rendimiento medio ya que el marcador del ambiente D, está muy cerca de la perpendicular que pasa por el origen.

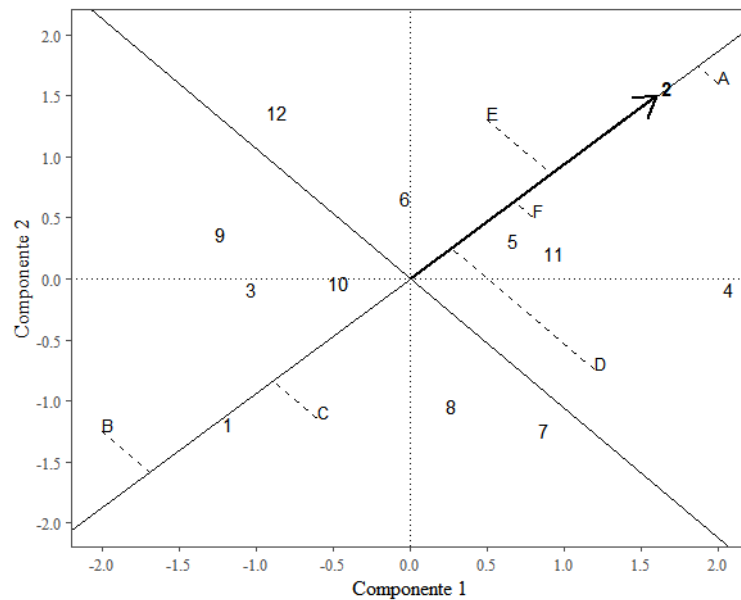


Figura 3.3: Ambientes favorables y desfavorables para el genotipo 2 en el biplot GGE

Relación entre ambientes

Comparación entre dos cultivos

Para comparar dos cultivos, se propone unir mediante una línea recta los genotipos a comparar, luego trazar una línea perpendicular a la anterior que pase por el origen. Esta última línea es la que separa sitios favorables a uno y a otro genotipo.

Al hacer comparaciones de a pares, la longitud de la línea que conecta los cultivos es importante para hacer inferencias. Cuanto más larga sea, más confiable será la comparación. Por el contrario, una línea corta implica que los dos cultivos a comparar son similares en todos los entornos y, por lo tanto, la comparación puede no tener sentido.

Los sitios cuyos marcadores queden en el mismo lado donde está el marcador del genotipo son los mejores para ese genotipo. Si un ambiente queda posicionado sobre la línea perpendicular, los dos genotipos tienen rendimientos similares en ese ambiente. Si dos genotipos están cercanos, sus rendimientos son similares en los ambientes evaluados. Por último, si todos los ambientes quedan a un lado de la línea perpendicular, el genotipo cuyo identificador está de ese lado rinde más que el otro en todos los ambientes.

En la Figura 3.4 se puede ver la comparación de los desempeños de los genotipos 6 y 8. Los ambientes que resultan favorables para el genotipo 6 son el E, el A y el F. Mientras que los favorables para el genotipo 8 son el B, el C y el D.

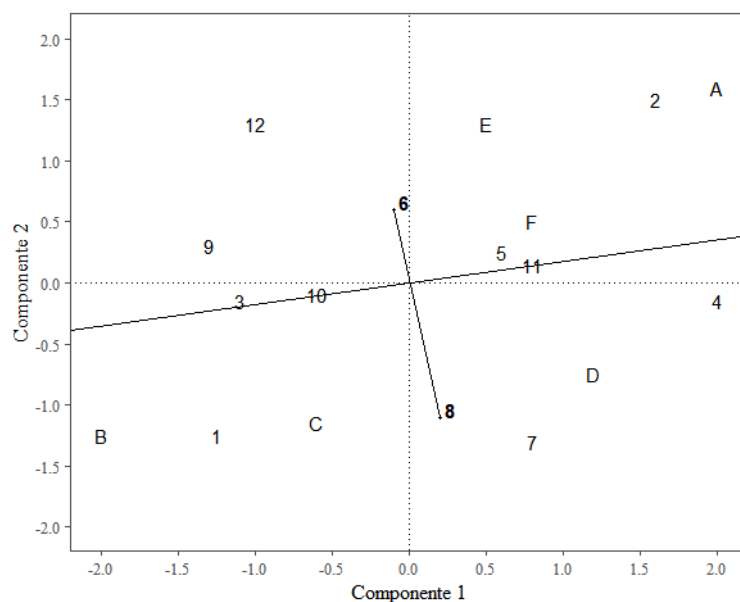


Figura 3.4: Comparación de los genotipos 6 y 8 en el biplot GGE

” Which- Won- Where”

Para poder identificar mega-ambientes y los mejores genotipos en cada uno de ellos se propone graficar un polígono envolvente. Este polígono se forma uniendo los genotipos más extremos en el biplot, de manera que todos los cultivares son contenidos en el mismo. Luego se trazan líneas rectas que pasan por el origen y que son perpendiculares a cada uno de los lados del polígono. De esta forma, el biplot queda dividido en sectores, y los sitios que quedan dentro un mismo sector se consideran pertenecientes a un mismo mega-ambiente. El cultivar que se encuentra en el vertice de cada sector es el de mayor rendimiento en todos los ambientes que comparten el sector con él. A los grupos de ambientes de cada sector se los llama mega-ambientes.

En la Figura 3.5 se presenta el biplot GGE con el polígono envolvente y las perpendiculares a sus lados, que ayudan a la interpretación del mismo.

En primer lugar se observa una mayor variabilidad en los ambientes A y B, es decir, en ellos es donde mejor se diferencian los efectos de los genotipos.

Las perpendiculares a los lados del polígono envolvente determinan tres mega-ambientes:

- uno formado por los ambientes A, E y F, en donde el genotipo de mejor desempeño es el 2 (se encuentra en el vértice del polígono encerrado por las perpendiculares).
- otro está formado solo por el ambiente D y el genotipo ganador en él es el 4, y
- el tercer mega-ambiente lo componen los ambientes B y C, en este caso el genotipo ganador es el 1.

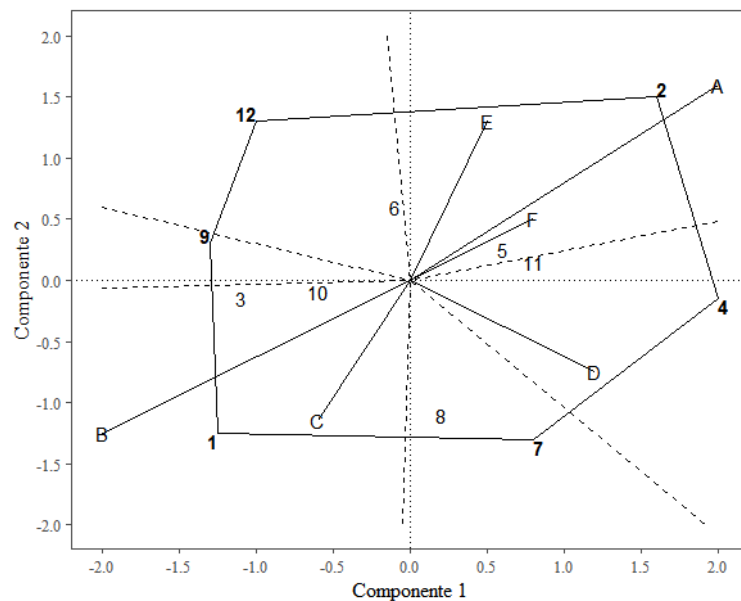


Figura 3.5: Biplot GGE con el polígono envolvente y las perpendiculares a sus lados

Evaluación de los cultivares basada en el desempeño y estabilidad

Si se identifican distintos mega-ambientes, la selección de genotipos debe hacerse para cada mega-ambiente en particular. Los genotipos se seleccionan en base a su desempeño y a su estabilidad a través de los ambientes. Esto se puede visualizar al graficar un eje medio para todos los ambientes pertenecientes a un mismo mega-ambiente. Para ello se traza una línea recta entre la media de scores de la componente 1 y la componente 2 de los ambientes pertenecientes al mega-ambiente y el origen; y una perpendicular a la línea media de scores de ambientes. Estas dos líneas constituyen “el eje de coordenadas de ambiente medio”.

Las proyecciones de los marcadores de los genotipos sobre el eje de la media de scores de ambientes da un ranking de los rendimientos de los genotipos en ese mega-ambiente. A su vez la magnitud de la proyección de los marcadores de los genotipos a la perpendicular al eje de ambientes da una idea de la estabilidad. Cuanto mayor sea esta magnitud, más inestable será el genotipo.

En este ejemplo se calcula el ambiente medio para un mega-ambiente formado por los ambientes A, E y F. El promedio de los scores de la primer componentes es 1,00 y el de la segunda es 1,13, por lo que el punto medio que determina la dirección del eje es (1,00; 1,13), que está graficado con un círculo en la Figura 3.6.

Como se puede observar en el biplot el orden de los genotipos (de mayor a menor rendimiento) es: 2, 4, 12, 11, 5, 6 todos ellos con rendimientos superiores al promedio, seguidos por los de rendimiento menor al promedio: el 10, 9, 7, 3 y por último el 1, el de peor rendimiento medio en ese mega-ambiente.

Debido a que las proyecciones sobre el eje perpendicular al eje medio de ambiente dan una idea de la estabilidad, se observa que el genotipo 12, el 9, el 7 y el 4 son los más inestables. También se observa que el genotipo 2, además de tener el mejor rendimiento medio es de los más estables en el megaambiente.

Comparación de todos los ambientes con respecto a uno ideal

Comparación de todos los genotipos con respecto a uno ideal

El cultivar ideal, representado por el círculo pequeño con una flecha apuntando hacia él, se define como el que tiene el mayor rendimiento en todos los entornos. es decir, tiene el rendimiento medio más alto y es absolutamente estable. los genotipos se clasifican según su distancia del cultivar ideal

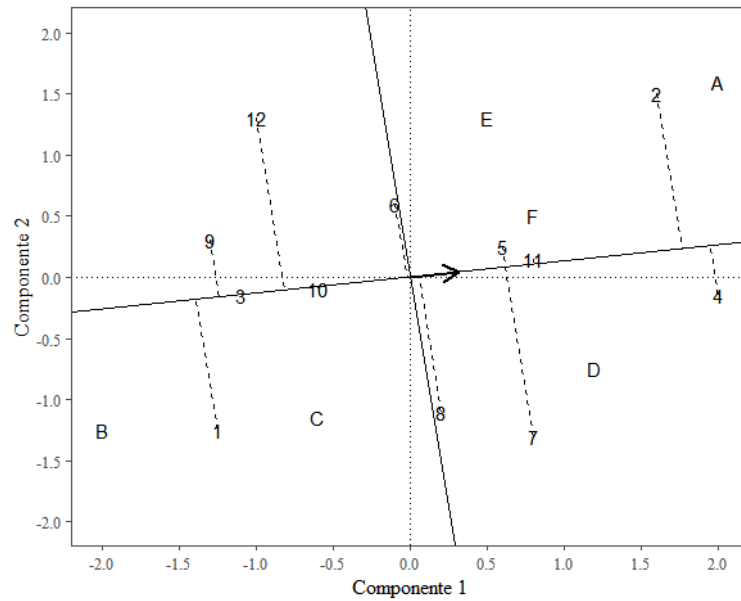


Figura 3.6: Eje de coordenadas de ambiente medio para un mega-ambiente en el biplot GGE

3.1.3. Métodos de imputación

Una limitación importante que presentan los modelos multiplicativos (AMMI y SREG) es que requieren que el conjunto de datos este completo, es decir no admiten valores perdidos. Aunque los EMA están diseñados para que todos los genotipos se evalúen en todos los ambientes, la presencia de valores faltantes es muy común, debido por ejemplo a la incorporación de nuevos genotipos, errores de medición o causas naturales como la destrucción de plantas por animales, inundaciones o durante la cosecha.

Entre las posibles soluciones para tratar un conjunto de datos con observaciones perdidas: el uso de un subconjunto completo de datos, eliminando aquellos genotipos que tienen valores faltantes (Ceccarelli et al., 2007, Yan et al., 2011), completar datos faltantes con la media ambiental, o imputación de los valores faltantes mediante estimaciones utilizando, por ejemplo, un modelo multiplicativo (Kumar et al., 2012).

Se han propuesto numerosas metodologías para superar el problema de valores ausentes en el conjunto de datos, entre las cuales se encuentran:

- EM-AMMI: Gauch y Zobel (1990) desarrollaron este enfoque mediante el cual se imputa utilizando el algoritmo de maximización de la esperanza (EM, del inglés *Expectation-Maximization*) incorporando el modelo AMMI. Consiste en un procedimiento iterativo que funciona de la siguiente forma: Dependiendo del número de términos multiplicativos empleados, el método de imputación puede denominarse EM-AMMI0, EM-AMMI1, etc. (Gauch y Zobel 1990). Los estudios de

Caliński y col. (1992), Piepho (1995), Arciniegas-Alarcón y Dias (2009) y Paderewski y Rodrigues (2014) mostraron que se obtienen los mejores resultados para la imputación con modelos AMMI al incluir como máximo una componente multiplicativa.

- EM-SVD: Perry (2009a) propone un método de imputación que combina el algoritmo EM con DVS. Este método reemplaza los valores faltantes de una matriz $G \times E$ inicialmente por valores arbitrarios para obtener una matriz completa, y luego la DVS se calcula iterativamente en esa matriz. Al final del proceso, cuando las iteraciones alcanzan estabilidad, se obtiene la matriz imputada.Podria ponerlo tambien de manera formal... pero creo que con esto basta
- EM-PCA:
- Gabriel Eigen: Arciniegas-Alarcón et al. (2010) propuso un método de imputación que combina regresión y aproximación de rango inferior usando DVS. El método reemplaza inicialmente las celdas faltantes por valores arbitrarios, y posteriormente el las imputaciones se refinan a través de un esquema iterativo que define una partición de la matriz para cada valor que falta a su vez y utiliza una regresión lineal de columnas (o filas) para obtener el nueva imputación. En esta regresión, la matriz de diseño se aproxima por una matriz de menor rango usando la DVS.
- WGabriel Eigen:

Vale la pena señalar que el modelo de análisis no siempre será el mismo que el modelo de imputación.

3.2. Paquete de R

Una librería o paquete (*package*) es una colección de objetos creados y organizados siguiendo un protocolo fijo que garantiza la ausencia de errores (de sintaxis) en la programación. Éstos son las unidades fundamentales de un código reproducible de R ya que incluyen funciones reutilizables, la documentación que describe cómo usar cada una de ellas y, además datos de ejemplo.

Los pasos necesarios para la creación de un paquete son:

- Creación del esqueleto del paquete.
- Inclusión de los objetos que contendrá el paquete (funciones y/o datos).
- Redacción de la documentación.

-
- Compilación del paquete en Linux y creación de la versión para Windows.
 - Instalación.
 - Prueba y publicación.

Para la creación del paquete se utilizan numerosas funciones incluidas en el paquete *devtools* que permiten realizar diversos aspectos del desarrollo de paquetes. Por lo tanto, antes de comenzar a crear el paquete se deben instalar el mismo como se indica a continuación:

```
# Instalar el paquete devtools desde CRAN
install.packages("devtools")
# Instalar el paquete devtools desde GitHub:
install_github("r-lib/devtools")
```

3.2.1. Esqueleto y estructura del paquete

Para crear la estructura del paquete se utiliza la función `create_package()`. El principal y único argumento requerido por dicha función es el directorio donde el nuevo paquete se alojará. Por lo general, si el directorio se llama “geneticae”, entonces el nombre del paquete también será “geneticae”:

```
# Cargar la libreria devtools
library(devtools)
# Crear el paquete genetiae
create_package("C:/Users/Julia/Desktop/geneticae")
```

El resultado de ejecutar la función `create_package()` es un paquete con los siguientes componentes:

- Un directorio R/.
- DESCRIPTION, un archivo simple cuyo objetivo es almacenar metadatos importantes sobre el paquete, especifica el título, la versión del mismo, identifica al autor y brinda un mail de contacto, una breve descripción del paquete, la lista de los paquetes que el paquete creado necesita para funcionar, la licencia, entre otros.

El contenido básico en un archivo DESCRIPTION es:

```
Package: genetiae
Title: What the Package Does (One Line, Title Case)
Version: 0.0.0.9000
Authors@R:
```

```
person(given = "First",
       family = "Last",
       role = c("aut", "cre"),
       email = "first.last@example.com",
       comment = c(ORCID = "YOUR-ORCID-ID"))
```

Description: What the package does (one paragraph).

License: What license it uses

Encoding: UTF-8

LazyData: true

- Un archivo NAMESPACE

Estas carpetas se irán modificando a medida que el paquete se vaya creando. También puede incluir un archivo de proyecto de RStudio `pkgname.Rproj`, que hace que su paquete sea fácil de usar con RStudio; `.Rbuildignore` enumera los archivos que se necesitan, pero que no deben incluirse al compilar el paquete R desde la fuente; `.gitignore` anticipa el uso de Git. Se crearán, a través de devtools, las siguientes carpetas: `data/` y `man/`.

3.2.2. Creación de funciones y conjuntos de datos

Una vez creada la estructura del paquete se deben incluir las funciones que el mismo contendrá. Cada una de ellas debe ser guardada en un archivo de extensión `.R`, en el subdirectorio `R/`. Para ello, se utiliza la función `use_r()` la cual crea y/o abre un script de la carpeta `R/`.

Una vez creada una función, se realizan pruebas para asegurar que el código realice lo que realmente se desea utilizando la función `load_all()` que simula el proceso de construcción, instalación y conexión del paquete. Permite que las funciones creadas estén disponible rápidamente para uso interactivo, del mismo modo que si se hubiera construido e instalado el paquete y luego cargada en la sesión de R a través de la función `library(geneticae)`.

Muy frecuentemente se utilizan funciones que se encuentran disponibles en otros paquetes, para ello se utiliza la función `use_package()` que agrega el paquete al archivo `DESCRIPTION`.

A menudo es útil incluir datos en un paquete a fin de proporcionar ejemplos de aplicaciones de las funciones incluidas en él. Ellos se almacenan en el directorio `data/`, siendo cada archivo un `.RData` que sólo contiene un objeto. Para esto, se utiliza la función `usethis::use_data()`. Notar que el archivo `DESCRIPTION` creado con la función `create_package()`, mencionada anteriormente, contiene el campo `LazyData: true`, lo cual genera que los conjuntos de datos no ocupen memoria hasta que sean usados.

3.2.3. Documentación

La documentación es uno de los aspectos más importantes del código, sin ella, los usuarios no sabrán cómo usar el paquete. Existen múltiples formas de documentar un paquete, la forma estándar es escribir archivos `.Rd` en la carpeta `man`, los cuales utilizan una sintaxis personalizada, basada en LaTeX. Sin embargo, el paquete *roxygen2*, utilizado en este trabajo, convierte los comentarios con formato especial en archivos `.Rd`. Esta última proporciona una serie de ventajas sobre la estándar:

- El código y la documentación son adyacentes, de modo que cuando el código se modifique, será fácil actualizar la documentación.
- Inspecciona dinámicamente los objetos que está documentando, para que pueda agregar automáticamente los datos que de otra forma se deben escribir a mano.
- Resume las diferencias en la documentación de los métodos S3 y S4, los genéricos y las clases, por lo que necesita aprender menos detalles.

Además de generar archivos `.Rd`, *roxygen2* también creará el archivo `NAMESPACE`. El flujo de trabajo para crear la documentación con el paquete *roxygen2* es el siguiente:

- Agregar comentarios a los archivos `.R`, los cuales comienzan con `#` y preceden a una función. La primera oración se convierte en el título y el segundo párrafo es una descripción de la función. Seguidamente, las funciones son documentadas, la mayoría de las funciones tienen tres etiquetas: `@param`, `@examples` y `@return`.
 - `@param` describe los parámetros de la función, indica de que clase es el parámetro y para que sirve.
 - `@examples` proporciona un código ejecutable que muestra cómo usar la función en la práctica.
 - `@return` describe el resultado de la función.
- Ejecutar `devtools::document()` para convertir los comentarios de *roxygen* en archivos `.Rd`.

Roxygen2 permite utilizar la descripción de los parámetros de otras funciones usando `@inheritParams`. Esta documentará los parámetros que no están documentados en la función actual, pero que si lo están en la función fuente. La fuente puede ser una función en el paquete actual, vía `@inheritParams function`, u otro paquete, vía `@inheritParams package::function`.

A diferencia de las funciones que son documentadas directamente, para los objetos en `data/`, se debe crear un archivo y guardarlo en el directorio `R/`.

Viñetas

A diferencia de la documentación, en la cual se detalla como se utiliza cada una de las funciones del paquete, una viñeta es una descripción el problema que el paquete está diseñado para resolver y muestra al lector cómo resolverlo.

Muchos de los paquetes existentes tienen viñetas la cual puede ser buscada mediante la función `browseVignettes("packagename")`. Cada viñeta proporciona el archivo fuente original, una página HTML o PDF y un archivo de código R. Las viñetas de paquetes que no han sido instalados pueden ser consultados en su página de CRAN, por ejemplo para el paquete `dplyr`: <http://cran.r-project.org/web/packages/dplyr>.

Las Viñetas se pueden construir de diversas formas, en este trabajo se utiliza *devtools* para crear la estructura de la misma y luego se añade el contenido que se desee en formato Rmarkdown. Rmarkdown permite combinar, texto plano, bloques de código y salidas.

Para crear la viñeta, se utiliza `usethis::use_vignette("my-vignette")`. La misma crea un directorio `vignettes/`, agrega las dependencias necesarias a `DESCRIPTION` y redacta la viñeta. Las tres componentes fundamentales de la misma son las siguientes:

- El bloque inicial de metadatos, que contiene la siguiente información:

```
---
title: "Vignette Title"
output: rmarkdown::html_vignette
vignette: >
  %\VignetteIndexEntry{Vignette Title}
  %\VignetteEngine{knitr::rmarkdown}
  \usepackage[utf8]{inputenc}
---
```

- Markdown para formatear texto.
- Knitr para interpretar texto, código y resultados.

3.2.4. Pruebas del flujo de trabajo

Las pruebas resultan fundamentales en el desarrollo de paquetes, asegura que el código haga lo que realmente se desea. Existen pruebas informales como aquellas realizadas con la función `load_all()`. Sin embargo, estas pruebas interactivas pueden convertirse en scripts reproducibles, los cuales resultan superiores debido a que:

-
- Se indica explícitamente cómo debería comportarse el código, provocando que los errores solucionados no vuelvan a ocurrir.
 - El código que es fácil de probar generalmente está mejor diseñado, reduce la duplicación en el código. Como resultado, las funciones serán más fáciles de probar, comprender y trabajar.
 - Si toda la funcionalidad del paquete tiene una prueba asociada, se pueden hacer grandes cambios sin preocuparse por generar errores.

Para ello se utiliza la función `usethis::usetestthat()` (Wickham,2011). Esta crea un directorio `tests/testthat`, agrega `testthat` al campo `Suggests` en el archivo `DESCRIPTION` y además, crea un archivo `tests/testthat.R`.

Las pruebas se organizan jerárquicamente, las expectativas se agrupan en pruebas que se organizan en archivos :

- Una expectativa describe el resultado esperado de un cálculo.
- Una prueba agrupa múltiples expectativas para probar la salida de una función simple, un rango de posibilidades para un solo parámetro de una función más complicada o una funcionalidad estrechamente relacionada de varias funciones.
- Un archivo agrupa múltiples pruebas relacionadas.

Existen tres formas de llevar a cabo las pruebas:

- Ejecutar todas las pruebas en un archivo o directorio `test_file()` o `test_dir()`.
- Ejecutar pruebas automáticamente cada vez que algo cambie con la función `auto-test()`. Estas son útiles cuando las pruebas se ejecutan con frecuencia. Si se modifica un archivo de prueba, probará ese archivo; si se modifica un archivo de código, volverá a cargar ese archivo y volverá a ejecutar todas las pruebas.
- Hacer que `R CMD check` ejecute sus pruebas.

3.2.5. Compilación e instalación

Mediante la función `load_all()` fue utilizado para simular el proceso de construcción, instalación y conexión del paquete, con el fin de ir probando las funciones creadas. Sin embargo, `R CMD check` ejecutado en el shell o la función `check()`, es utilizado para verificar que un paquete R esta en pleno funcionamiento. La misma verificará que no haya errores de sintaxis o no se generen warnings. Está compuesto por más de 50 chequeos individuales

entre los cuales se encuentran: la estructura del paquete, el archivo descripción, namespace, el código de R, los datos, la documentación, entre otros. Se aconseja realizar verificaciones completas de que todo funciona a medida que se van incorporando funciones ya que si se incorporan muchas y luego se verifican será difícil identificar y resolver los problemas. Una vez que las verificaciones completas no encuentran errores, advertencias o notas, se ejecuta la función `install()`, con el objetivo de instalar el paquete en la biblioteca.

3.2.6. Publicación

Un repositorio es el lugar dónde se encuentran alojados los paquetes y desde el cuál se pueden descargarlos. Entre los repositorios más populares de paquetes R se encuentran:

- **CRAN**: es el principal repositorio de paquetes de R, está coordinado por la fundación R. Previa a la publicación en este repositorio el paquete debe pasar por diferentes pruebas para asegurar que cumple con las políticas de CRAN.
- **Bioconductor**: se trata de un repositorio específico para bioinformática. Del mismo modo que CRAN, tiene sus propias políticas de publicaciones y procesos de revisión.
- **GitHub**: a pesar que no es específico para R, github es con toda seguridad el repositorio más popular para la publicación de proyectos *open source* (del inglés, código abierto). Su popularidad procede del espacio ilimitado que proporciona para el alojamiento de proyectos *open source*, la integración con git (un software de control de versiones) y, la facilidad de compartir y colaborar con otras personas. Una de sus desventajas es que no proporciona procesos de control.
- **R-Forge** y **RForge**: son entornos de desarrollo de paquetes y repositorios. Eso significa que incluyen control de fuente, seguimiento de errores y otras características. Puede obtener versiones de desarrollo de paquetes de estos.

El paquete *geneticae* se encuentra en GitHub, para instalar el mismo se deben seguir las siguientes instrucciones:

```
library(devtools)
install_github("jangelini/geneticae")
```

3.3. Shiny APP

Una aplicación web es una aplicación o herramienta informática accesible desde cualquier navegador, bien sea a través de internet (lo habitual) o bien a través de una red

local. Estas aplicaciones son muy populares hoy en día para los usuarios no expertos, debido a la facilidad de su uso, ya que no requiere de una instalación en el ordenador del usuario, simplemente se accede a través de un navegador. Por lo que es posible utilizar una aplicación web desde cualquier dispositivo con conexión a internet, ya sea un ordenador, un smartphone o una tablet, es decir que es independiente del sistema operativo del usuario. Otra gran ventaja es el bajo consumo de recursos, ya que la mayor parte del tiempo estos se consumen en el servidor donde se encuentra alojada la aplicación, que generalmente tiene mucha más potencia de cómputo que cualquier ordenador personal.

Crear aplicaciones web puede resultar difícil para la mayoría de los usuarios de R debido a que se necesita un conocimiento profundo de las tecnologías web como HTML, CSS y JavaScript; y además hacer aplicaciones interactivas complejas requiere un análisis cuidadoso de los flujos de interacción para asegurarse de que cuando una entrada cambie, solo se actualicen las salidas relacionadas. Shiny es un paquete R que te permite crear aplicaciones web interactivas, permitiendo exhibir un trabajo de R a través de un navegador web para que cualquiera pueda usarlo. Este paquete hace que sea mucho más fácil para el programador R crear aplicaciones web al proporcionar un conjunto de funciones de interfaz de usuario (UI para abreviar) que generan el HTML, CSS y JavaScript que necesita para tareas comunes. Esto significa que no se necesita conocer los detalles de HTML / CSS / JS.

Los dos componentes clave de una Shiny APP son:

- *ui* (*user interfaz*): la interfaz de usuario controla el diseño de la aplicación, recibe los inputs y muestra los outputs en el navegador.
- *server*, funciones de R que contienen las instrucciones que se necesitan para construir los resultados de los análisis incluidos en la aplicación.
- *shinyApp*, función que crea objetos de aplicación Shiny a partir de *ui* / servidor.

El esquema interno de la aplicación puede observarse en la Figura 3.7.

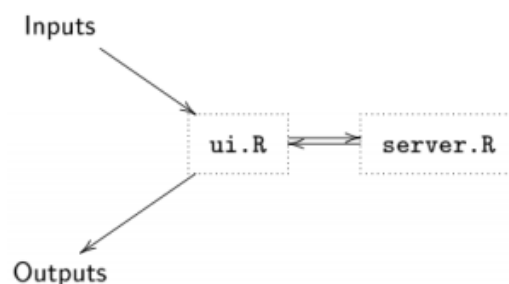


Figura 3.7: Esquema interno de la aplicación.

3.3.1. Flujo de trabajo

En esta sección se mostrará como mejorar dos flujos de trabajo de Shiny importantes: el ciclo de desarrollo básico de crear aplicaciones, realizar cambios y experimentar con los resultados; y la depuración, proceso de identificar y corregir errores de programación.

1. Flujo de trabajo de desarrollo

El objetivo de optimizar el flujo de trabajo de desarrollo es reducir el tiempo entre hacer un cambio y ver el resultado. Cuanto más rápido se pueda iterar, más rápido se podrá experimentar y más rápido se podrá obtener la Shiny APP. Aquí hay dos flujos de trabajo principales para optimizar: crear la aplicación por primera vez y acelerar el ciclo iterativo de ajustar el código y probar los resultados.

Creación de la Shiny APP

Para poder crear una shiny APP se debe tener instalado R, RStudio, y el paquete shiny:

Para crear una Shiny APP, lo más simple es crear un nuevo directorio para la aplicación con un sólo archivo llamado app.R. Este archivo se usará para indicarle a Shiny cómo debería verse la aplicación y cómo debería comportarse.

```
library(shiny)
ui<- ...
server<- ...
shinyApp(ui = ui, server = server)
```

Por lo tanto, en el archivo app.R se realizan las siguientes tareas:

- Carga el paquete shiny: `library(shiny)`
- Define la interfaz de usuario, la página web HTML con la que los usuarios interactúan.
- Especifica el comportamiento de la aplicación definiendo la función `server`.
- Se ejecuta función `shinyApp(ui, server)` para construir e iniciar una aplicación Shiny desde la interfaz de usuario y el servidor.

La sesión de R estará monitoreando la aplicación y ejecutando las reacciones de la aplicación mientras la aplicación Shiny esté activa, por lo que no podrá ejecutar ningún comando.

En todo tipo de programación, es una mala práctica tener código duplicado; puede ser un desperdicio computacional y, lo que es más importante, aumenta la dificultad de

mantener o depurar el código. En la secuencia de comandos R tradicional, se utilizan dos técnicas para lidiar con el código duplicado: capturar el valor usando una variable o capturar el cálculo con una función. Desafortunadamente, ninguno de estos enfoques funciona en una Shiny APP y se necesita un nuevo mecanismo: expresiones reactivas. Una expresión reactiva tiene una diferencia importante con una variable: sólo se ejecuta la primera vez que se llama y luego almacena en caché el resultado de la misma hasta que necesite actualizarse.

La programación reactiva es un estilo de programación que enfatiza valores que cambian con el tiempo, y cálculos y acciones que dependen de esos valores. Esto es importante para las aplicaciones Shiny porque son interactivas: los usuarios cambian los inputs, lo que hace que la lógica se ejecute en el servidor que finalmente resultan en actualización de los outputs/resultados.

Ver los cambios

Al crear o modificar la aplicación, se la ejecuta para poder ver los cambios realizados, por lo que el dominio de flujo de trabajo de desarrollo es especialmente importante. La primera forma de reducir la velocidad de iteración es evitar hacer clic en el botón “Ejecutar aplicación” y, en su lugar, aprender el método abreviado de teclado Cmd/Ctrl+ Shift+ Enter. Esto brinda el siguiente flujo de trabajo de desarrollo:

1. Escribir un código.
2. Iniciar la aplicación con Cmd/Ctrl+ Shift+ Enter.
3. Experimentar interactivamente con la aplicación.
4. Cerrar la aplicación.
5. Ir a 1.

Otra forma de reducir aún más la velocidad de iteración es activar la recarga automática (`options(shiny.autoreload = TRUE)`) y luego ejecutar la aplicación en un trabajo en segundo plano. Con este flujo de trabajo cuando se guarde un archivo, su aplicación se reiniciará: no es necesario cerrarla y reiniciarla. Esto conduce a un flujo de trabajo aún más rápido:

1. Escribir un código y presione Cmd/Ctrl+S para guardar en el archivo.
2. Experimentar interactivamente.
3. Ir a 1.

La principal desventaja de esta técnica es que debido a que la aplicación se ejecuta en un proceso separado, es considerablemente más difícil de depurar.

De manera predeterminada, cuando ejecuta la aplicación, aparecerá en una ventana emergente. Sin embargo, existen otras dos opciones que puede elegir del menú desplegable *Run App*

1. La ejecución en el panel del visor es útil para aplicaciones más pequeñas porque puede verla al mismo tiempo que ejecuta el código de la aplicación.
2. Ejecutar en un navegador externo es útil para aplicaciones más grandes, o si desea verificar que su aplicación se ve exactamente de la manera que espera en el contexto que la mayoría de los usuarios la verán.

2. Depuración

Entre los problemas que pueden surgir al crear una Shiny app se encuentran los siguientes:

- Error inesperado. Este es el caso más fácil, porque obtendrá un rastreo que le permitirá averiguar exactamente de dónde proviene el error. Una vez que haya identificado el problema, deberá probar sistemáticamente su suposición hasta que encuentre una diferencia entre sus expectativas y lo que realmente está sucediendo. El depurador interactivo es una herramienta poderosa para este proceso.
- No obtiene ningún error, pero un valor es incorrecto. Aquí, generalmente es mejor transformar esto en el primer problema utilizando `stop()` para arrojar un error cuando se produce el valor incorrecto.
- Todos los valores son correctos, pero no se actualizan cuando espera. Este es el problema más desafiante porque es exclusivo de Shiny, por lo que no puede aprovechar sus habilidades de depuración de R.

Una vez localizado la fuente del error, la herramienta más poderosa es el depurador interactivo. El depurador detiene la ejecución y le brinda una consola R interactiva donde puede ejecutar cualquier código para descubrir qué salió mal. Hay dos formas de iniciar el depurador:

- Agregar una llamada a la función `browser()` en código fuente. Esta es la forma estándar de R de iniciar el depurador interactivo, y funcionará sin embargo, se está ejecutando brillante.

-
- Agregar un punto de interrupción RStudio haciendo clic a la izquierda del número de línea. Puede eliminar el punto de interrupción haciendo clic en el círculo rojo. La ventaja de los puntos de interrupción es que no son código, por lo que nunca tendrá que preocuparse por registrarlos accidentalmente en su sistema de control de versiones.

3.3.2. Compartiendo una Shiny Web App

Una vez creada la aplicación, resulta conveniente ponerlas a disposición de los usuarios. En este caso la Shiny Web App encuentra disponible en el servidor de CONICET www.cefobi.com. Además el proyecto se encuentra en GitHub https://github.com/jangelini/shinyAPP_geneticae.

Capítulo 4

Resultados

4.1. Paquete de R *geneticae*

El paquete *geneticae* permite analizar datos provenientes de etapas avanzadas de los programas de mejoramiento, donde se evalúan pocos genotipos.

Una vez instalado el paquete, se debe cargar en la sesión de R mediante el comando: `library(geneticae)`

Es posible obtener información detallada sobre las funciones del paquete *geneticae* mediante de los archivos de ayuda indicando `help(package = "geneticae")`. La ayuda para una función, por ejemplo, `imputation()`, en una sesión R se puede obtener usando `?imputation` o `help(imputation)`. Además, a partir de la función `browseVignettes("geneticae")` se obtiene la viñeta del mismo, es decir una descripción el problema que el paquete está diseñado para resolver así como ejemplos de aplicación del mismo.

4.1.1. Conjuntos de datos en *geneticae*

El paquete *geneticae* proporciona dos conjuntos de datos para ilustrar la metodología incluida para analizar los datos obtenidos de MET.

- `yan.winterwheat dataset`: rendimiento de 18 variedades de trigo de invierno cultivadas en nueve ambientes en Ontario en 1993. No hay réplicas disponibles en los datos. Este conjunto de datos se obtuvo del paquete *agridat*.

```
data(yan.winterwheat)
dat_yan <- yan.winterwheat
head(dat_yan)
```

```
##   gen  env yield
## 1 Ann BH93 4.460
## 2 Ari BH93 4.417
## 3 Aug BH93 4.669
## 4 Cas BH93 4.732
## 5 Del BH93 4.390
## 6 Dia BH93 5.178
```

- plrv dataset: rendimiento, peso de planta y parcela de 28 clones de la población del virus del enrollamiento de la papa (PLRV) evaluada en seis ambientes. Las réplicas están disponibles en los datos. Este conjunto de datos se obtuvo del paquete agricolae.

```
data(plrv)
dat_rep <- plrv
head(dat_rep)
```

```
##   Genotype Locality Rep WeightPlant WeightPlot   Yield
## 1   102.18    Ayac   1   0.5100000      5.10 18.88889
## 2   104.22    Ayac   1   0.3450000      2.76 12.77778
## 3   121.31    Ayac   1   0.5425000      4.34 20.09259
## 4   141.28    Ayac   1   0.9888889      8.90 36.62551
## 5   157.26    Ayac   1   0.6250000      5.00 23.14815
## 6   163.9     Ayac   1   0.5120000      2.56 18.96296
```

4.1.2. Funciones en geneticae

Modelo de regresión por sitio

Para ejecutar la función `GGEmodel()`, se debe proporcionar un conjunto de datos con genotipos, ambientes, repeticiones (si hay disponibles), el fenotipo observado y los nombres que dichas variables tienen en el archivo de entrada. Además, se debe indicar el método de centrado, escala y SVD.

Cuando no hay repeticiones disponibles en el conjunto de datos, como es el caso del conjunto de datos `yan.winterwheat`, el modelo GGE se indica de la siguiente manera:

```
GGE1 <- GGEmodel(dat_yan, genotype = "gen", environment = "env", response = "yield"
, centering = "tester")
```

Sin embargo, en el caso de que haya repeticiones disponibles, como el conjunto de datos `plrv`, se indica de la siguiente manera:

```
GGE1_rep <- GGEmodel(dat_rep, genotype = "Genotype", environment = "Locality",  
  response = "Yield", rep = "Rep", centering = "tester")
```

La salida de la función `GGEmodel()` es una lista con los siguientes elementos:

- `coordgenotype`: trazado de coordenadas para genotipos de todos los componentes.
- `coordenviroment`: trazado de coordenadas para entornos de todos los componentes.
- `valores propios`: vector de valores propios de cada componente.
- `vartotal`: varianza general.
- `varexpl`: porcentaje de varianza explicado por cada componente.
- `labelgen`: nombres de genotipo.
- `labelenv`: nombres de entorno.
- `ejes`: etiquetas de eje.
- `Datos`: datos de entrada escalados y centrados.
- `centrado`: nombre del método de centrado.
- `escala`: nombre del método de escala.
- `SVP`: nombre del método SVP.

Biplot GGE

Para ejecutar la función `GGEPlot()`, se requiere un objeto de la clase `GGEmodel()`. La salida es un biplot construido a través de los componentes principales generados por `GGEmodel()`. Los diferentes biplots que se pueden obtener usando la función `GGEPlot()` se muestran usando el conjunto de datos `yan.winterwheat`. Si hay repeticiones disponibles en el conjunto de datos, como es el caso del conjunto `plrv`, se debe indicar el nombre de la columna que contiene las réplicas en el archivo de entrada.

- Biplot básico

```
GGEPlot(GGE1, type = "Biplot")
```

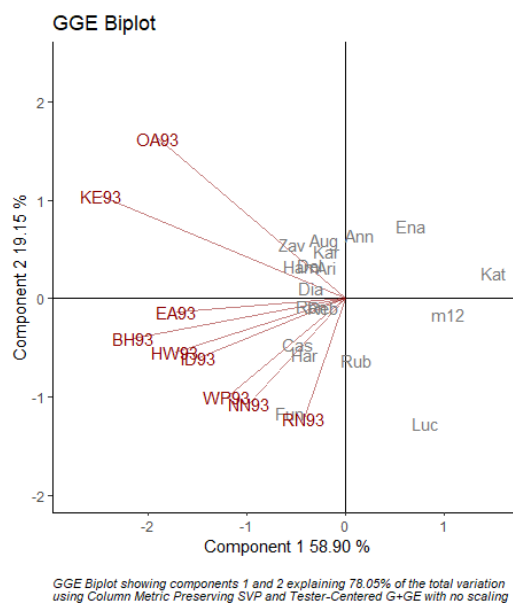


Figura 4.1: Biplot básico obtenido de la función `GGEPlot()`

- Ranking de los cultivares en función de su rendimiento en el ambiente OA93.

```
GGEPlot(GGE1, type = "Selected Environment", selectedE = "OA93")
```

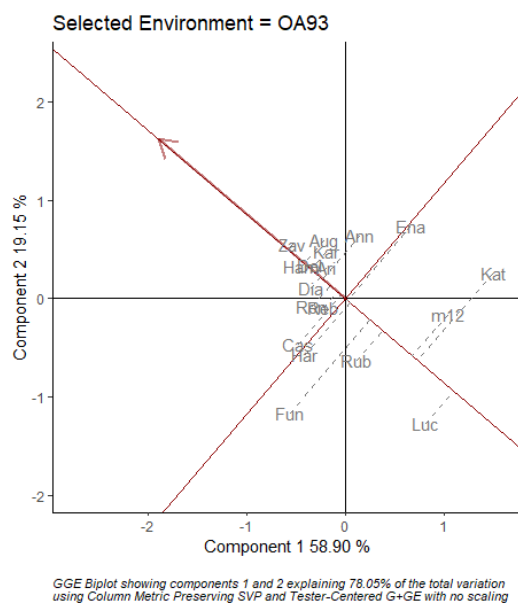
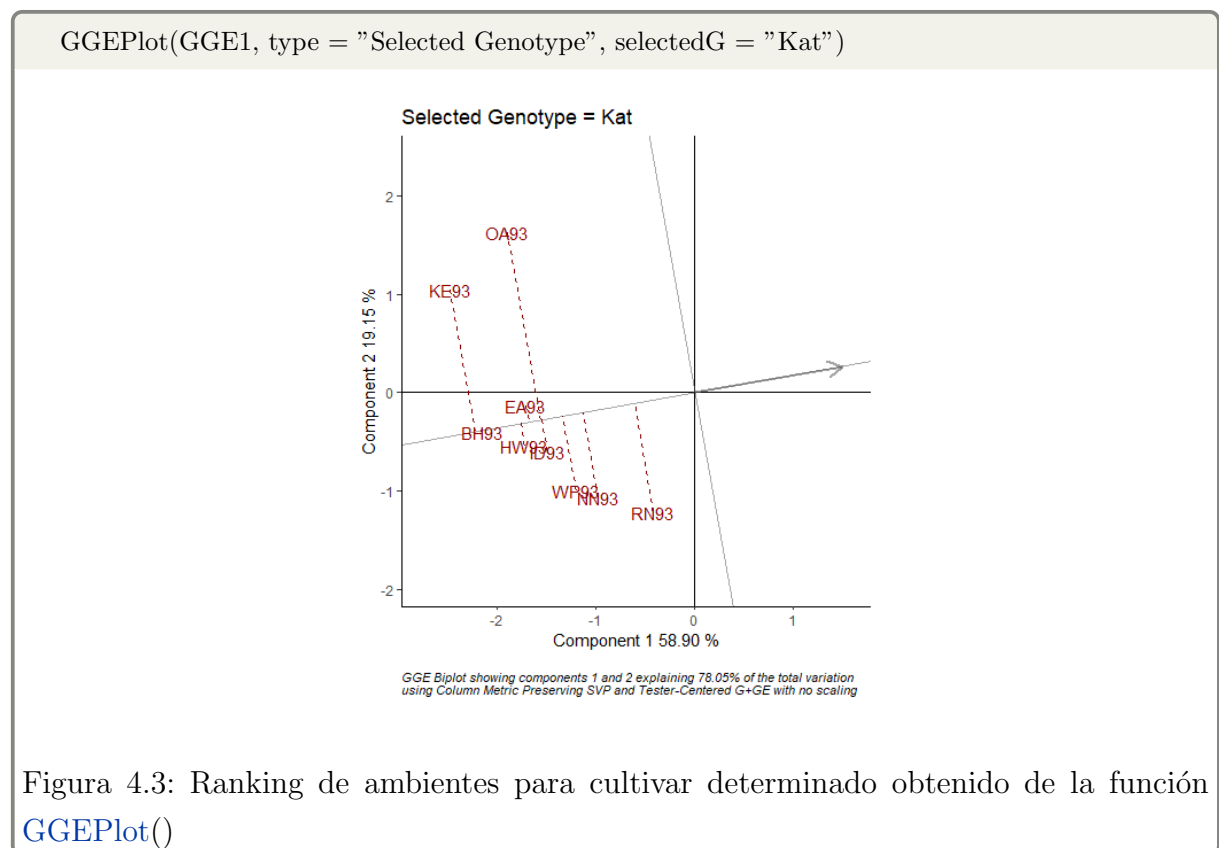


Figura 4.2: Ranking de cultivares para un ambiente determinado obtenido de la función `GGEPlot()`

- Ranking de los ambientes en función del rendimiento relativo del cultivar Kat.



- Relación entre ambientes.

GGEPlot(GGE1, type = "Relationship Among Environments")

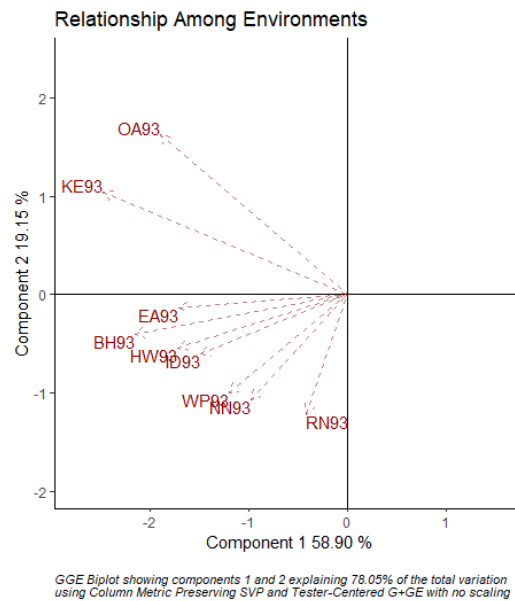


Figura 4.4: Relación entre ambientes obtenido de la función `GGEPlot()`

- Comparación entre los genotipos Kat y Cas.

GGEPlot(GGE1, type = "Comparison of Genotype", selectedG1 = "Kat", selectedG2 = "Cas")

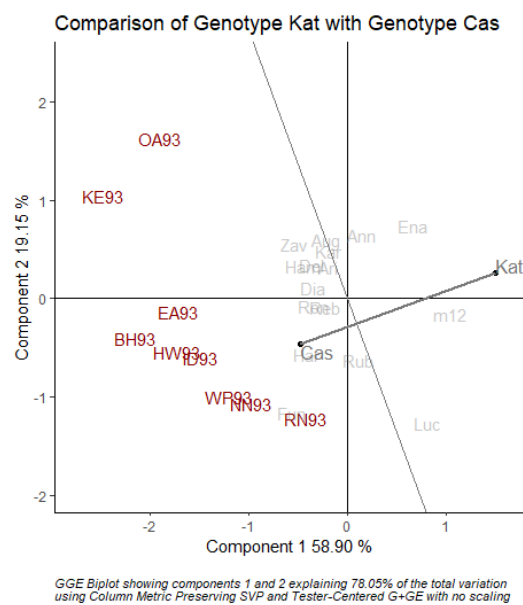
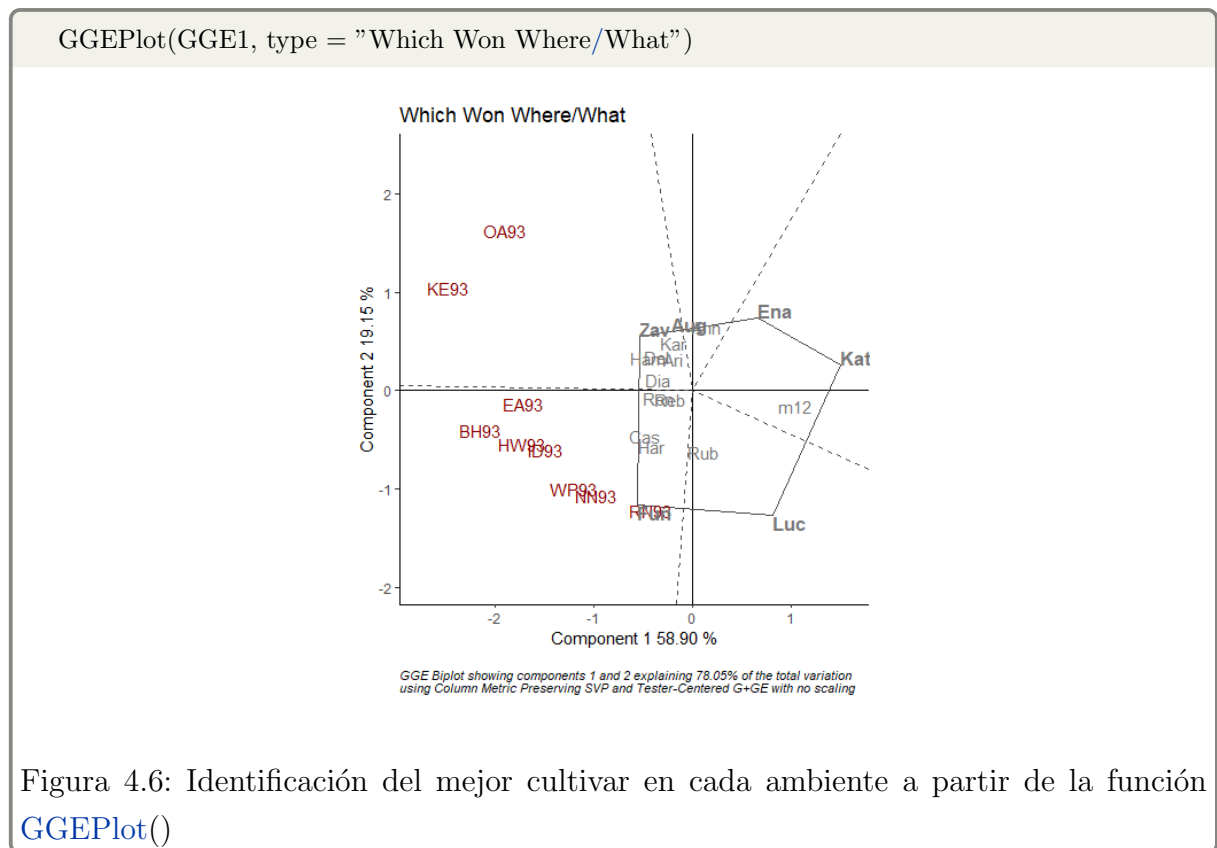


Figura 4.5: Comparación entre dos genotipos obtenido de la función `GGEPlot()`

- Identificación del mejor cultivar en cada ambiente.



- Evaluación de los ambientes basados tanto en la capacidad de discriminación como en la representatividad.

GGEPlot(GGE1, type = "Discrimination vs. representativeness")

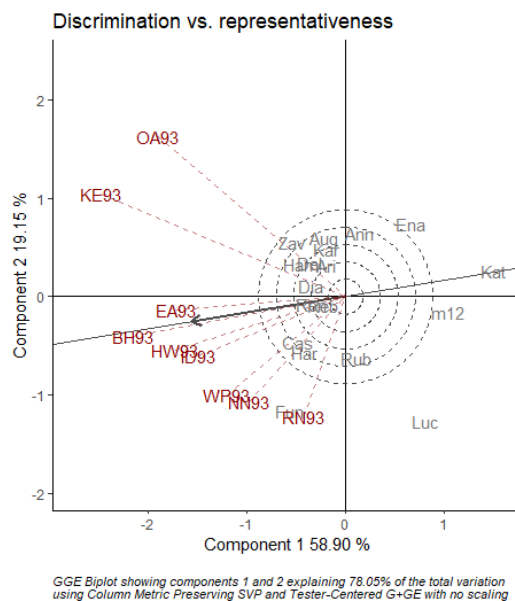


Figura 4.7: Evaluación de los ambientes basados tanto en la capacidad de discriminación y representatividad a partir de la función `GGEPlot()`

- Clasificación de ambientes con respecto al ambiente ideal.

GGEPlot(GGE1, type = "Ranking Environments")

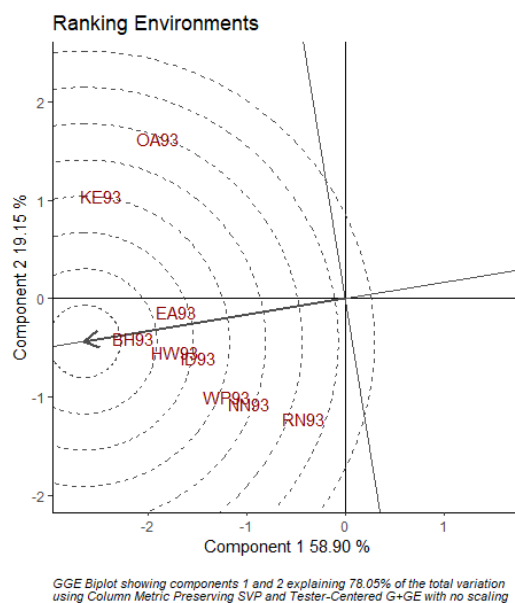


Figura 4.8: Clasificación de ambientes con respecto al ambiente ideal a partir de la función `GGEPlot()`

- Clasificación de genotipos con respecto al genotipo ideal.

GGEPlot(GGE1, type = "Ranking Genotypes")

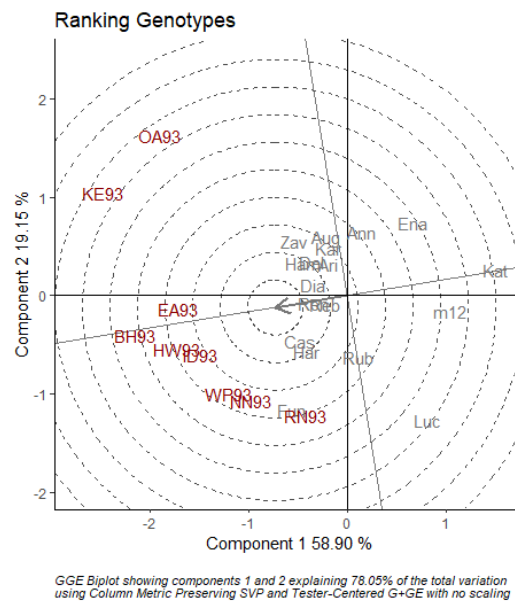
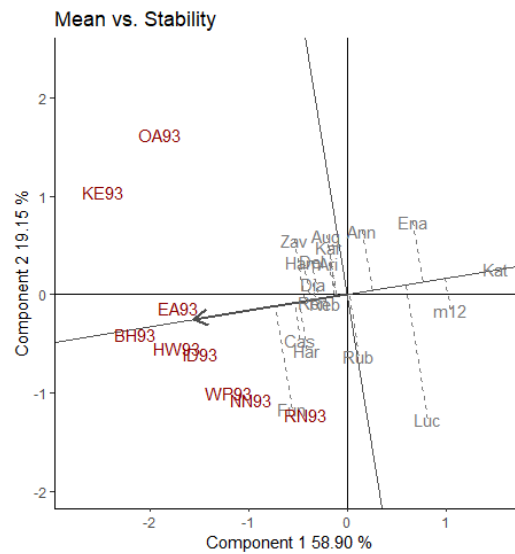


Figura 4.9: Clasificación de genotipos con respecto al genotipo ideal a partir de la función `GGEPlot()`

- Evaluación de los cultivares con base en el rendimiento promedio y la estabilidad.

```
GGEPlot(GGE1, type = "Mean vs. Stability")
```



GGE Biplot showing components 1 and 2 explaining 78.05% of the total variation using Column Metric Preserving SVP and Tester-Centered G+GE with no scaling

Figura 4.10: Evaluación de los cultivares con base en el rendimiento promedio y la estabilidad a partir de la función `GGEPlot()`

Classic AMMI model

Para ejecutar la función `rAMMI()`, como en la función `GGEmodel()`, se debe proporcionar un conjunto de datos con genotipo, entorno, repeticiones (si las hay) y la variable de respuesta. Se debe indicar el nombre de las columnas que contienen cada una de estas variables en el conjunto de datos de entradas. La salida de la función es un biplot.

A continuación se muestra el biplot GE obtenido del modelo AMMI clásico obtenido con el conjunto de datos `yan.winterwheat`.

```
rAMMI(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "AMMI")
```

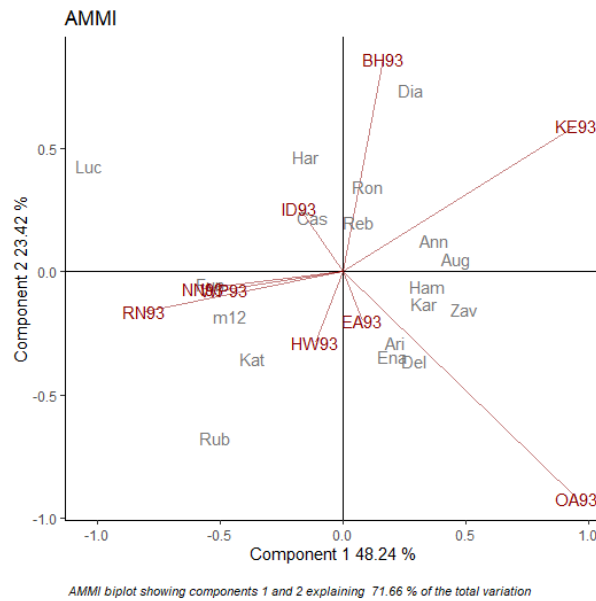


Figura 4.11: Biplot GE obtenido del modelo clasico AMMI

Robust AMMI model

Como se dijo anteriormente, el modelo AMMI clasico, en su forma estándar, no funciona bien en presencia de observaciones atípicas. Dado que los outliers son muy comun en los datos agronómicos, Rodrigues et al. (2015) proponen cinco modelos AMMI robustos, que permiten superar el problema de la contaminación de datos con observaciones atípicas. Los biplots de los cinco modelos AMMI robustos propuestos por Rodrigues et al. (2015), se pueden obtener utilizando la función `rAMMI()` A continuación se muestran los biplots obtenidos con dichos modelos robustos usando el conjunto de datos yan.winterwheat.

- modelo rAMMI"

```
rAMMI(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "rAMMI")
```

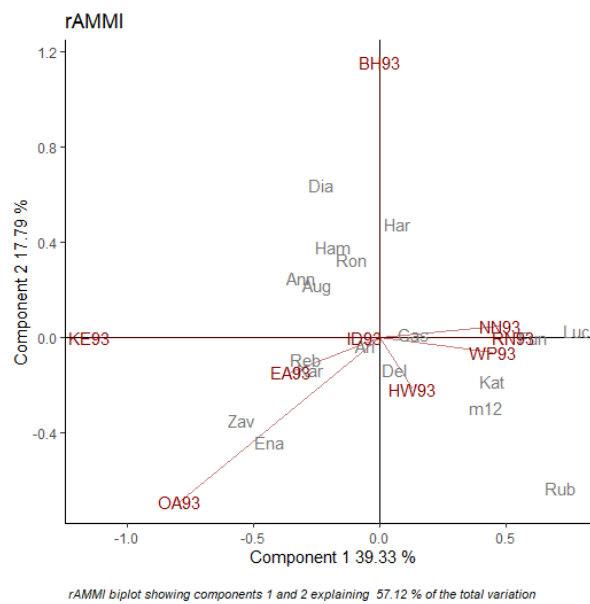


Figura 4.12: Biplot GE obtenido del modelo robusto rAMMI

■ modelo "hAMMI"

```
rAMMI(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "hAMMI")
```

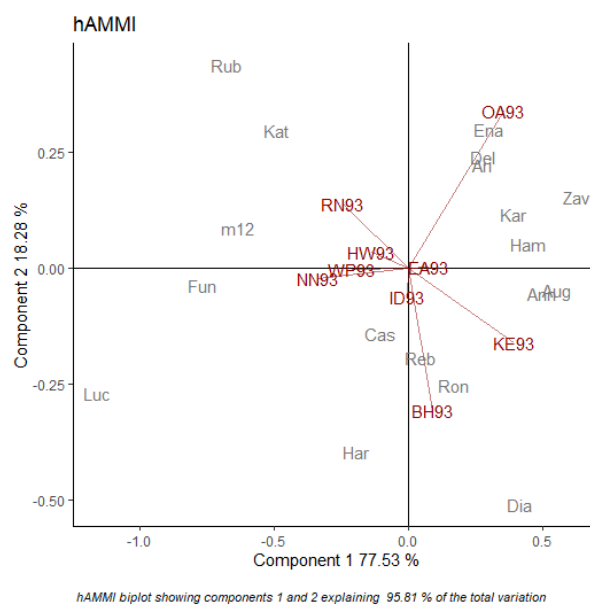


Figura 4.13: Biplot GE obtenido del modelo robusto hAMMI

- modelo "gAMMI"

```
rAMMI(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "gAMMI")
```

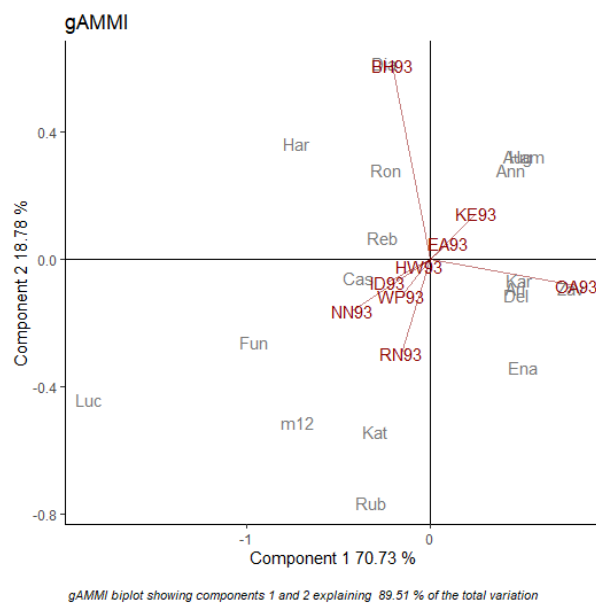


Figura 4.14: Biplot GE obtenido del modelo robusto gAMMI

- modelo "lAMMI"

```
rAMMI(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "IAMMI")
```

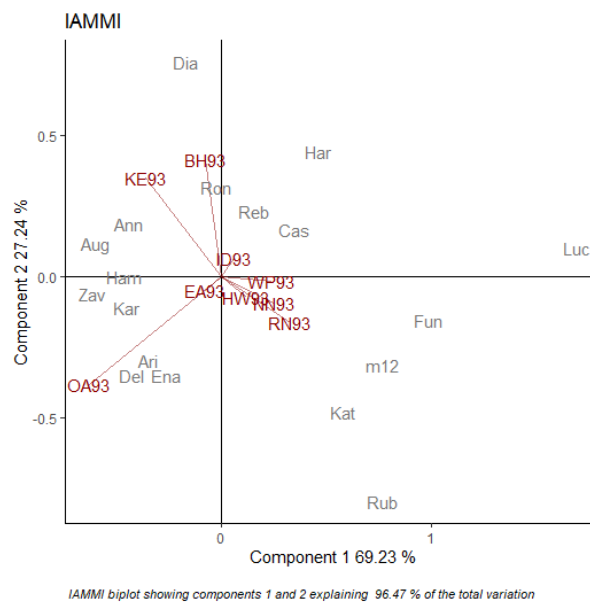


Figura 4.15: Biplot GE obtenido del modelo robusto IAMMI

■ modelo "ppAMMI"

```
rAMMI(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "ppAMMI")
```

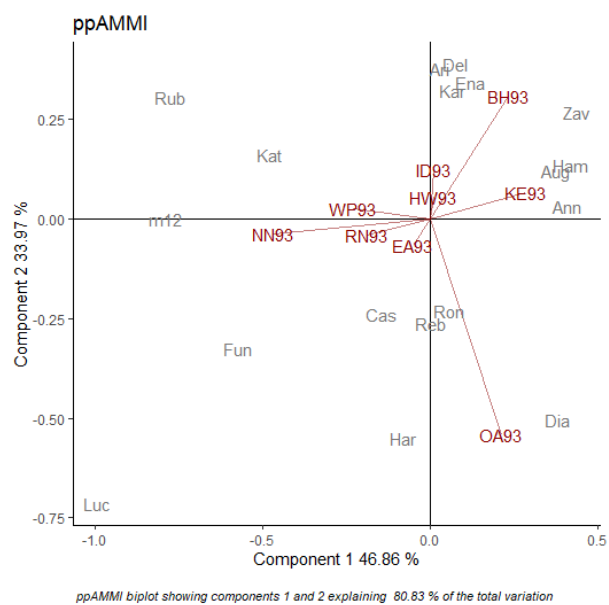


Figura 4.16: Biplot GE obtenido del modelo robusto ppAMMI

Métodos de imputación Una limitación importante de los modelos presentados anteriormente es que requieren una que el conjunto de datos este completo. Por lo tanto, en el paquete se incluyen una serie de metodologías propuestas, algunas de las cuales no se encuentran disponible en R, para superar el problema de falta de equilibrio.

El conjunto de datos `yan.winterwheat` se utilizó como ejemplo. Como el conjunto de datos no contaba con observaciones perdidas, algunas fueron eliminadas con el objetivo de mostrar las metodologías de imputación incluidas.

```
# generates missing data
dat_yan[1, 3] <- NA
dat_yan[3, 3] <- NA
dat_yan[2, 3] <- NA
```

- GabrielEigein proposed by Arciniegas-Alarcón S., et al. (2010).

```
imputation(dat_yan, PC.nb = 2, genotype = "gen", environment = "env", response = "yield", type = "EM-AMMI")

##           BH93  EA93  HW93  ID93  KE93  NN93  OA93  RN93  WP93
## Ann 4.150120 4.150 2.849 3.084 5.940 4.450 4.351 4.039 2.672
## Ari 4.035814 4.771 2.912 3.506 5.699 5.152 4.956 4.386 2.938
## Aug 4.305244 4.578 3.098 3.460 6.070 5.025 4.730 3.900 2.621
## Cas 4.732000 4.745 3.375 3.904 6.224 5.340 4.226 4.893 3.451
## Del 4.390000 4.603 3.511 3.848 5.773 5.421 5.147 4.098 2.832
## Dia 5.178000 4.475 2.990 3.774 6.583 5.045 3.985 4.271 2.776
## Ena 3.375000 4.175 2.741 3.157 5.342 4.267 4.162 4.063 2.032
## Fun 4.852000 4.664 4.425 3.952 5.536 5.832 4.168 5.060 3.574
## Ham 5.038000 4.741 3.508 3.437 5.960 4.859 4.977 4.514 2.859
## Har 5.195000 4.662 3.596 3.759 5.937 5.345 3.895 4.450 3.300
## Kar 4.293000 4.530 2.760 3.422 6.142 5.250 4.856 4.137 3.149
## Kat 3.151000 3.040 2.388 2.350 4.229 4.257 3.384 4.071 2.103
## Luc 4.104000 3.878 2.302 3.718 4.555 5.149 2.596 4.956 2.886
## Reb 4.375000 4.701 3.655 3.592 6.189 5.141 3.933 4.208 2.925
## Ron 4.940000 4.698 2.950 3.898 6.063 5.326 4.302 4.299 3.031
## Rub 3.786000 4.969 3.379 3.353 4.774 5.304 4.322 4.858 3.382
## Zav 4.238000 4.654 3.607 3.914 6.641 4.830 5.014 4.363 3.111
## m12 3.340000 3.854 2.419 2.783 4.629 5.090 3.281 3.918 2.561
```

- EM-AMMI proposed by Gauch and Zobel (1990).


```
imputation(dat_yan, PC.nb = 1, genotype = "gen", environment = "env", response = "yield", type = "EM-AMMI")
```

```
##           BH93  EA93  HW93  ID93  KE93  NN93  OA93  RN93  WP93
## Ann 4.136249 4.150 2.849 3.084 5.940 4.450 4.351 4.039 2.672
## Ari 4.474249 4.771 2.912 3.506 5.699 5.152 4.956 4.386 2.938
## Aug 4.386299 4.578 3.098 3.460 6.070 5.025 4.730 3.900 2.621
## Cas 4.732000 4.745 3.375 3.904 6.224 5.340 4.226 4.893 3.451
## Del 4.390000 4.603 3.511 3.848 5.773 5.421 5.147 4.098 2.832
## Dia 5.178000 4.475 2.990 3.774 6.583 5.045 3.985 4.271 2.776
## Ena 3.375000 4.175 2.741 3.157 5.342 4.267 4.162 4.063 2.032
## Fun 4.852000 4.664 4.425 3.952 5.536 5.832 4.168 5.060 3.574
## Ham 5.038000 4.741 3.508 3.437 5.960 4.859 4.977 4.514 2.859
## Har 5.195000 4.662 3.596 3.759 5.937 5.345 3.895 4.450 3.300
## Kar 4.293000 4.530 2.760 3.422 6.142 5.250 4.856 4.137 3.149
## Kat 3.151000 3.040 2.388 2.350 4.229 4.257 3.384 4.071 2.103
## Luc 4.104000 3.878 2.302 3.718 4.555 5.149 2.596 4.956 2.886
## Reb 4.375000 4.701 3.655 3.592 6.189 5.141 3.933 4.208 2.925
## Ron 4.940000 4.698 2.950 3.898 6.063 5.326 4.302 4.299 3.031
## Rub 3.786000 4.969 3.379 3.353 4.774 5.304 4.322 4.858 3.382
## Zav 4.238000 4.654 3.607 3.914 6.641 4.830 5.014 4.363 3.111
## m12 3.340000 3.854 2.419 2.783 4.629 5.090 3.281 3.918 2.561
```

- EM-SVD proposed by Perry (2009)

```
imputation(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "
EM-SVD")
```

```
##           [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9]
## [1,] 4.332467 4.150 2.849 3.084 5.940 4.450 4.351 4.039 2.672
## [2,] 4.332467 4.771 2.912 3.506 5.699 5.152 4.956 4.386 2.938
## [3,] 4.332467 4.578 3.098 3.460 6.070 5.025 4.730 3.900 2.621
## [4,] 4.732000 4.745 3.375 3.904 6.224 5.340 4.226 4.893 3.451
## [5,] 4.390000 4.603 3.511 3.848 5.773 5.421 5.147 4.098 2.832
## [6,] 5.178000 4.475 2.990 3.774 6.583 5.045 3.985 4.271 2.776
## [7,] 3.375000 4.175 2.741 3.157 5.342 4.267 4.162 4.063 2.032
## [8,] 4.852000 4.664 4.425 3.952 5.536 5.832 4.168 5.060 3.574
## [9,] 5.038000 4.741 3.508 3.437 5.960 4.859 4.977 4.514 2.859
## [10,] 5.195000 4.662 3.596 3.759 5.937 5.345 3.895 4.450 3.300
## [11,] 4.293000 4.530 2.760 3.422 6.142 5.250 4.856 4.137 3.149
## [12,] 3.151000 3.040 2.388 2.350 4.229 4.257 3.384 4.071 2.103
## [13,] 4.104000 3.878 2.302 3.718 4.555 5.149 2.596 4.956 2.886
## [14,] 4.375000 4.701 3.655 3.592 6.189 5.141 3.933 4.208 2.925
## [15,] 4.940000 4.698 2.950 3.898 6.063 5.326 4.302 4.299 3.031
## [16,] 3.786000 4.969 3.379 3.353 4.774 5.304 4.322 4.858 3.382
## [17,] 4.238000 4.654 3.607 3.914 6.641 4.830 5.014 4.363 3.111
## [18,] 3.340000 3.854 2.419 2.783 4.629 5.090 3.281 3.918 2.561
```

- WGabriel proposed by Alarcon.

```
imputation(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "
WGabriel")
```

```
##          BH93  EA93  HW93  ID93  KE93  NN93  OA93  RN93  WP93
## Ann 4.004664 4.150 2.849 3.084 5.940 4.450 4.351 4.039 2.672
## Ari 4.455727 4.771 2.912 3.506 5.699 5.152 4.956 4.386 2.938
## Aug 4.328095 4.578 3.098 3.460 6.070 5.025 4.730 3.900 2.621
## Cas 4.732000 4.745 3.375 3.904 6.224 5.340 4.226 4.893 3.451
## Del 4.390000 4.603 3.511 3.848 5.773 5.421 5.147 4.098 2.832
## Dia 5.178000 4.475 2.990 3.774 6.583 5.045 3.985 4.271 2.776
## Ena 3.375000 4.175 2.741 3.157 5.342 4.267 4.162 4.063 2.032
## Fun 4.852000 4.664 4.425 3.952 5.536 5.832 4.168 5.060 3.574
## Ham 5.038000 4.741 3.508 3.437 5.960 4.859 4.977 4.514 2.859
## Har 5.195000 4.662 3.596 3.759 5.937 5.345 3.895 4.450 3.300
## Kar 4.293000 4.530 2.760 3.422 6.142 5.250 4.856 4.137 3.149
## Kat 3.151000 3.040 2.388 2.350 4.229 4.257 3.384 4.071 2.103
## Luc 4.104000 3.878 2.302 3.718 4.555 5.149 2.596 4.956 2.886
## Reb 4.375000 4.701 3.655 3.592 6.189 5.141 3.933 4.208 2.925
## Ron 4.940000 4.698 2.950 3.898 6.063 5.326 4.302 4.299 3.031
## Rub 3.786000 4.969 3.379 3.353 4.774 5.304 4.322 4.858 3.382
## Zav 4.238000 4.654 3.607 3.914 6.641 4.830 5.014 4.363 3.111
## m12 3.340000 3.854 2.419 2.783 4.629 5.090 3.281 3.918 2.561
```

- EM-PCA proposed by

```
imputation(dat_yan, genotype = "gen", environment = "env", response = "yield", type = "EM-PCA")
```

```
##           BH93  EA93  HW93  ID93  KE93  NN93  OA93  RN93  WP93
## Ann 3.980317 4.150 2.849 3.084 5.940 4.450 4.351 4.039 2.672
## Ari 4.463093 4.771 2.912 3.506 5.699 5.152 4.956 4.386 2.938
## Aug 4.327731 4.578 3.098 3.460 6.070 5.025 4.730 3.900 2.621
## Cas 4.732000 4.745 3.375 3.904 6.224 5.340 4.226 4.893 3.451
## Del 4.390000 4.603 3.511 3.848 5.773 5.421 5.147 4.098 2.832
## Dia 5.178000 4.475 2.990 3.774 6.583 5.045 3.985 4.271 2.776
## Ena 3.375000 4.175 2.741 3.157 5.342 4.267 4.162 4.063 2.032
## Fun 4.852000 4.664 4.425 3.952 5.536 5.832 4.168 5.060 3.574
## Ham 5.038000 4.741 3.508 3.437 5.960 4.859 4.977 4.514 2.859
## Har 5.195000 4.662 3.596 3.759 5.937 5.345 3.895 4.450 3.300
## Kar 4.293000 4.530 2.760 3.422 6.142 5.250 4.856 4.137 3.149
## Kat 3.151000 3.040 2.388 2.350 4.229 4.257 3.384 4.071 2.103
## Luc 4.104000 3.878 2.302 3.718 4.555 5.149 2.596 4.956 2.886
## Reb 4.375000 4.701 3.655 3.592 6.189 5.141 3.933 4.208 2.925
## Ron 4.940000 4.698 2.950 3.898 6.063 5.326 4.302 4.299 3.031
## Rub 3.786000 4.969 3.379 3.353 4.774 5.304 4.322 4.858 3.382
## Zav 4.238000 4.654 3.607 3.914 6.641 4.830 5.014 4.363 3.111
## m12 3.340000 3.854 2.419 2.783 4.629 5.090 3.281 3.918 2.561
```

4.2. Geneticae Shiny Web App

La aplicación Geneticae se organiza en las siguientes pestañas:

- Los datos
- Análisis descriptivo
- ANOVA
- Biplot GGE
- Biplot GE
- Ayuda

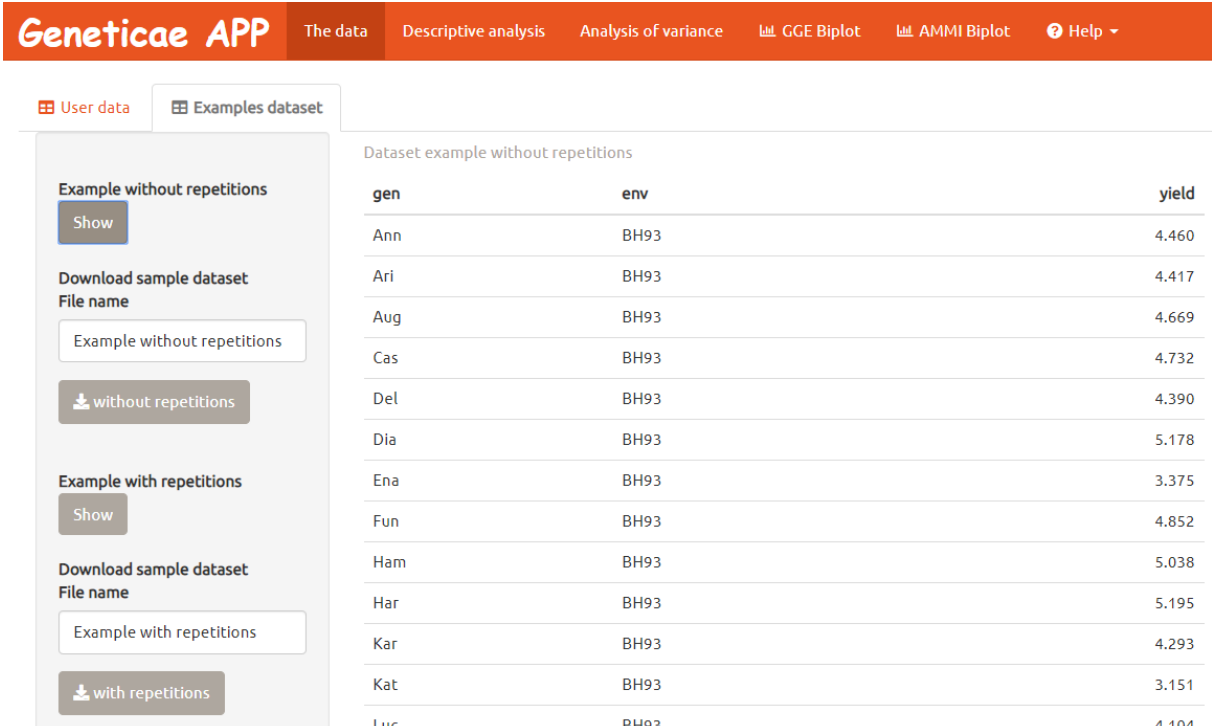
En muchos casos, algunos atributos estilísticos de salida pueden personalizarse para que el usuario obtenga la salida a su gusto. A su vez, los gráficos obtenidos pueden ser descargados.

4.2.1. Los datos

Al iniciar la aplicación Geneticae, se muestra una pantalla en la cual se carga el conjunto de datos a analizar. La aplicación admite datos en formato .csv, delimitados por coma o punto y coma; y permite el siguiente formato de datos:

- Cada fila contiene una observación, en la cual deben estar presentes los siguientes datos: nombre del cultivar, ambiente, repetición si está disponible y valor fenotípico medido. Pueden estar presentes otras variables que no serán utilizadas por la aplicación.
- La primera fila de encabezado contiene los nombres de cada variable. Los encabezados pueden dar cualquier nombre que elija, y deben indicarse al cargar el archivo de datos.
- El número de repeticiones puede diferir con los genotipos y los ambientes.

Se utilizan dos conjuntos de datos, incluidos en el paquete Geneticae, para ilustrar la aplicación. Estos conjuntos de datos, uno de los cuales tiene repeticiones (plr dataset) y el otro no (yan.winterwheat dataset), los cuales se pueden ver y descargar en la pestaña *The data* → *Example datasets* (Figura 4.17, 4.18).



Geneticae APP | The data | Descriptive analysis | Analysis of variance | GGE Biplot | AMMI Biplot | Help

User data | **Examples dataset**

Example without repetitions
Show
Download sample dataset
File name
Example without repetitions
without repetitions

Example with repetitions
Show
Download sample dataset
File name
Example with repetitions
with repetitions

Dataset example without repetitions

gen	env	yield
Ann	BH93	4.460
Ari	BH93	4.417
Aug	BH93	4.669
Cas	BH93	4.732
Del	BH93	4.390
Dia	BH93	5.178
Ena	BH93	3.375
Fun	BH93	4.852
Ham	BH93	5.038
Har	BH93	5.195
Kar	BH93	4.293
Kat	BH93	3.151
Luc	BH93	4.104

Figura 4.17: yan.winterwheat dataset disponible en Shiny Web App

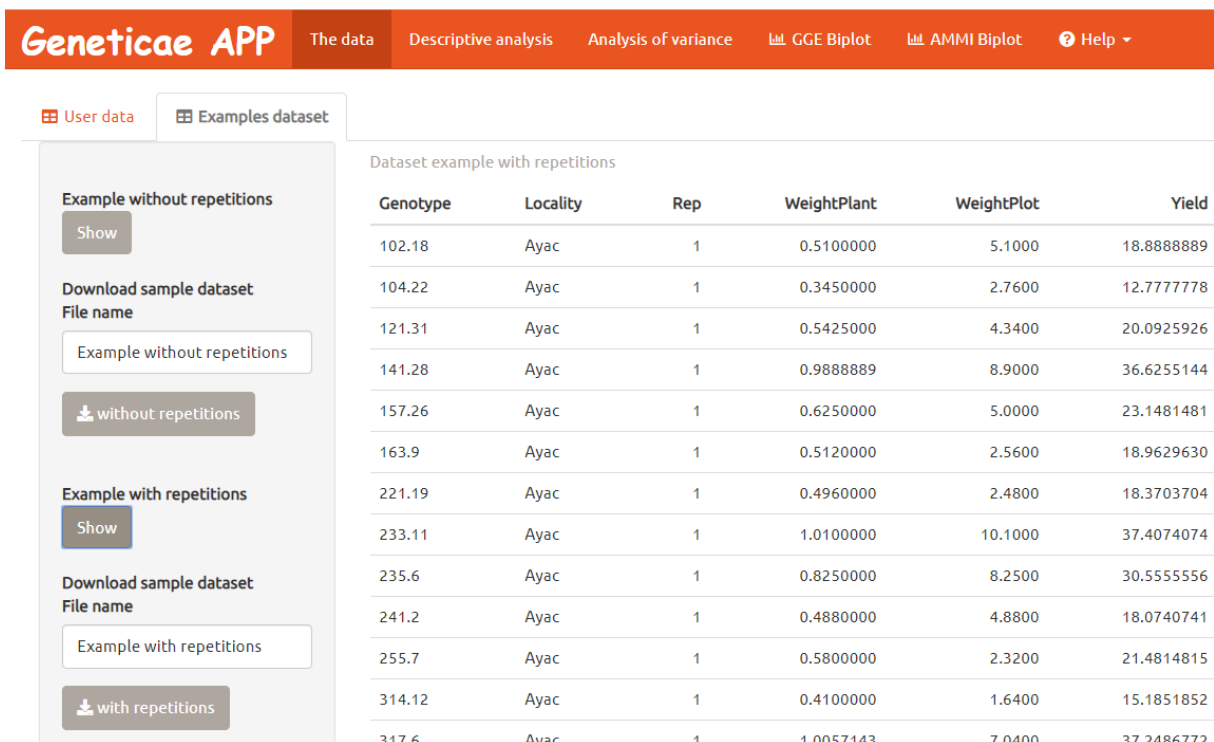


Figura 4.18: plrv dataset disponible en Shiny Web App

4.2.2. Análisis descriptivo

El menú *Descriptive analysis* le permite describir un conjunto de datos utilizando diagrama de caja (o *boxplot*), gráfico y matriz de correlación y gráfico de interacción.

4.2.2.1. *Boxplot*

El *boxplot* proporciona una medida central, la mediana y una idea de la dispersión a través del rango y el rango intercuartil. La posición de la mediana dentro de la caja y la similitud en la longitud de los bigotes nos dan una idea de la simetría de la distribución.

Un *boxplot* interactivo que compara el carácter cuantitativo de interés a través de genotipos, así como a través de los ambientes se pueden obtener (Figura 4.19, 4.20). A partir de los mismos se pueden obtener medidas resumen en forma interactiva usando el *Toggle Spike Lines* como se muestra en la figura 4.19. Estos gráficos se pueden descargar en formato interactivo (.HTML) a partir del botón *Download* (Figura 4.19 y 4.20), así como también en formato .png como se muestra en la Figura 4.20.

Boxplot

Correlation plot

Correlation matrix

Interaction Plot

Variable

☒ Environment

☐ Genotype

Fill color

orange

X-axis label:

Environments

Y-axis label:

Yield

Run

File name

Boxplot

Download

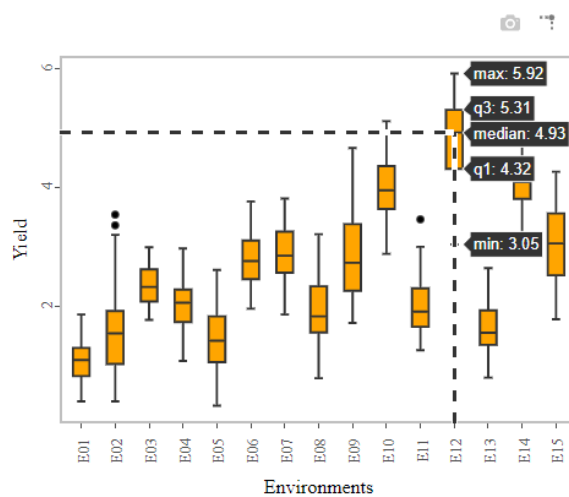


Figura 4.19: Boxplot de ambientes a través de los genotipos para el conjunto de datos Plrv



Figura 4.20: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

4.2.2.2. Gráfico de correlación

El correlograma o gráfico de correlación muestra la correlación tanto entre los genotipos como entre los ambientes (Figura 4.21 y 4.22). Se pueden mostrar las correlaciones de Pearson y Spearman. Las correlaciones positivas se muestran en azul y las negativas en rojo. La intensidad del color y el tamaño del círculo son proporcionales a los coeficientes de correlación.

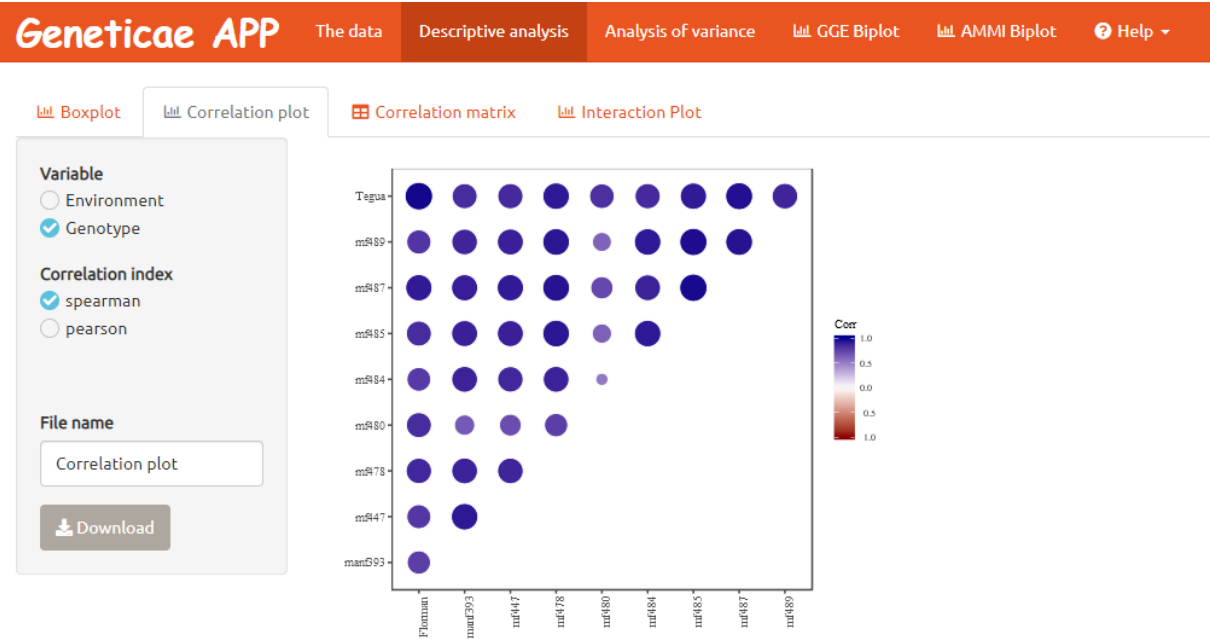


Figura 4.21: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

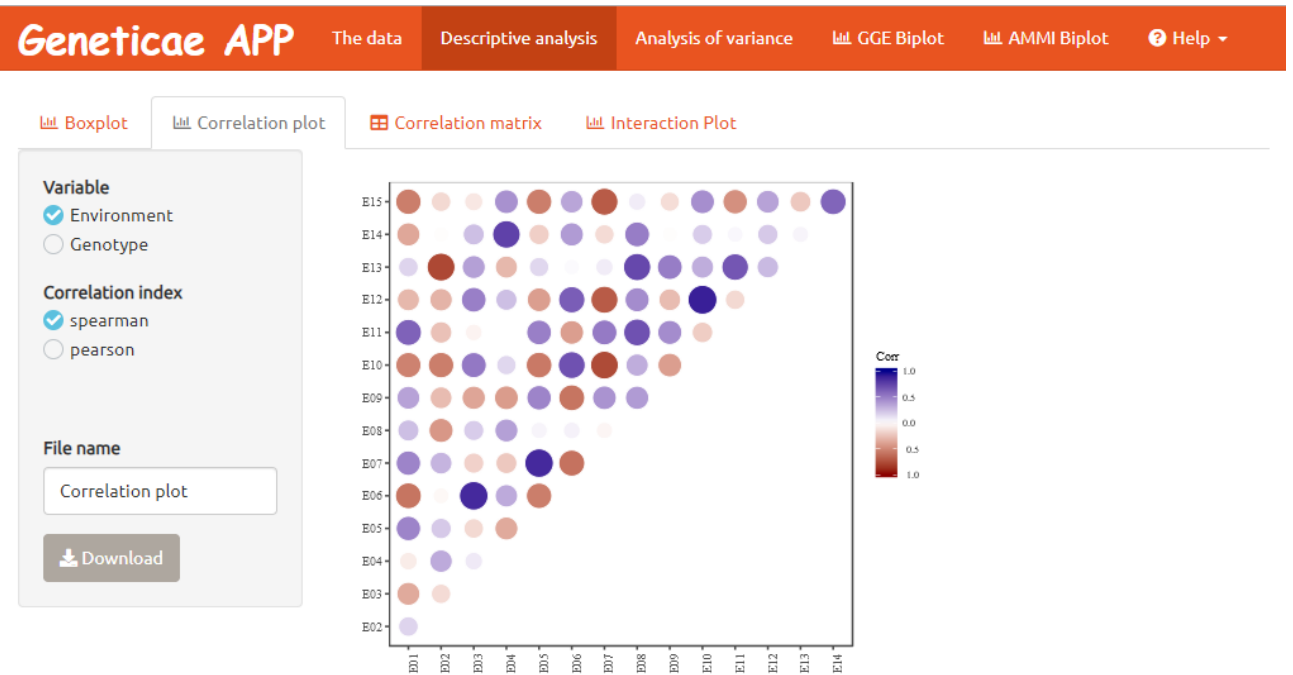


Figura 4.22: Boxplot de ambientes a través de los genotipos para el conjunto de datos Plrv

4.2.2.3. Matriz de correlación

Una matriz de correlación se utiliza como una forma de resumir datos. Muestra los coeficientes de correlación de pares de variables. Las correlaciones de Spearman o Pearson se pueden calcular tanto para ambientes como para genotipos (Figura 4.23).

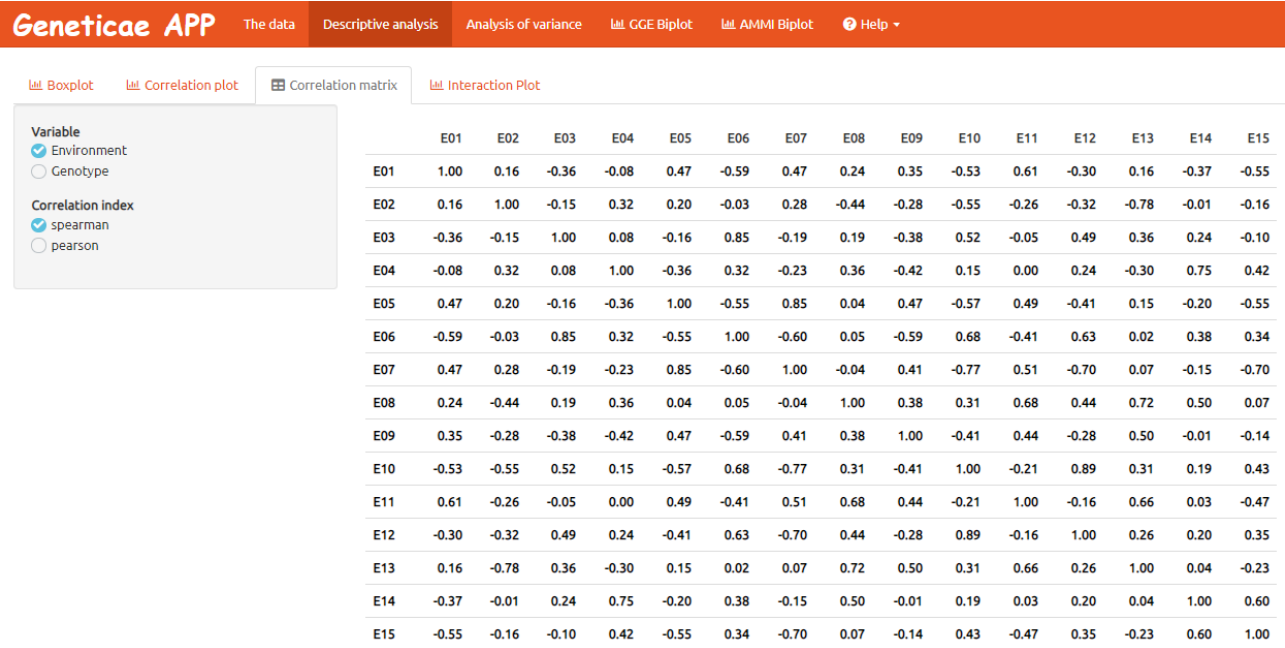


Figura 4.23: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

4.2.2.4. Gráfico de interacción

Un diagrama de interacción es una representación visual de la interacción entre los efectos de dos factores, o entre un factor y una variable numérica.

Se puede obtener el gráfico interactivo que muestra el cambio en el efecto genotípico a través de los entornos y también el que muestra el cambio en el efecto ambiental a través de los genotipos (Figura 4.25,4.26). Es posible descargarlo en formato interactivo (.HTML) a partir del boton *Download* (Figura 4.25), así como también en formato .png como se muestra en la Figura 4.26.



Figura 4.24: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

4.2.3. Análisis de la variancia

Cuando se pretende llevar a cabo el análisis de la variancia si el conjunto de datos tiene repeticiones entonces saldrá un mensaje en el cual se aclara que la interacción puede ser testada debido a la presencia de repeticiones "The interaction effect can be tested since there are repetitions in the data set", si no hay repeticiones disponibles entonces el mensaje será que la interacción no puede testarse.



Figura 4.25: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

El ANOVA depende del cumplimiento de los supuestos de que los errores tengan distribución normal con media cero y variancia constante. Por ello, tres pestañas: *Check normality*, *Check homoscedasticity* y *Outliers* se encuentran disponibles para la verificación de los supuestos mencionados.

Para verificar el supuesto de normalidad, se puede realizar un histograma, un gráfico de probabilidad normal y la prueba de shapiro-wilks sobre los residuos del ANOVA.

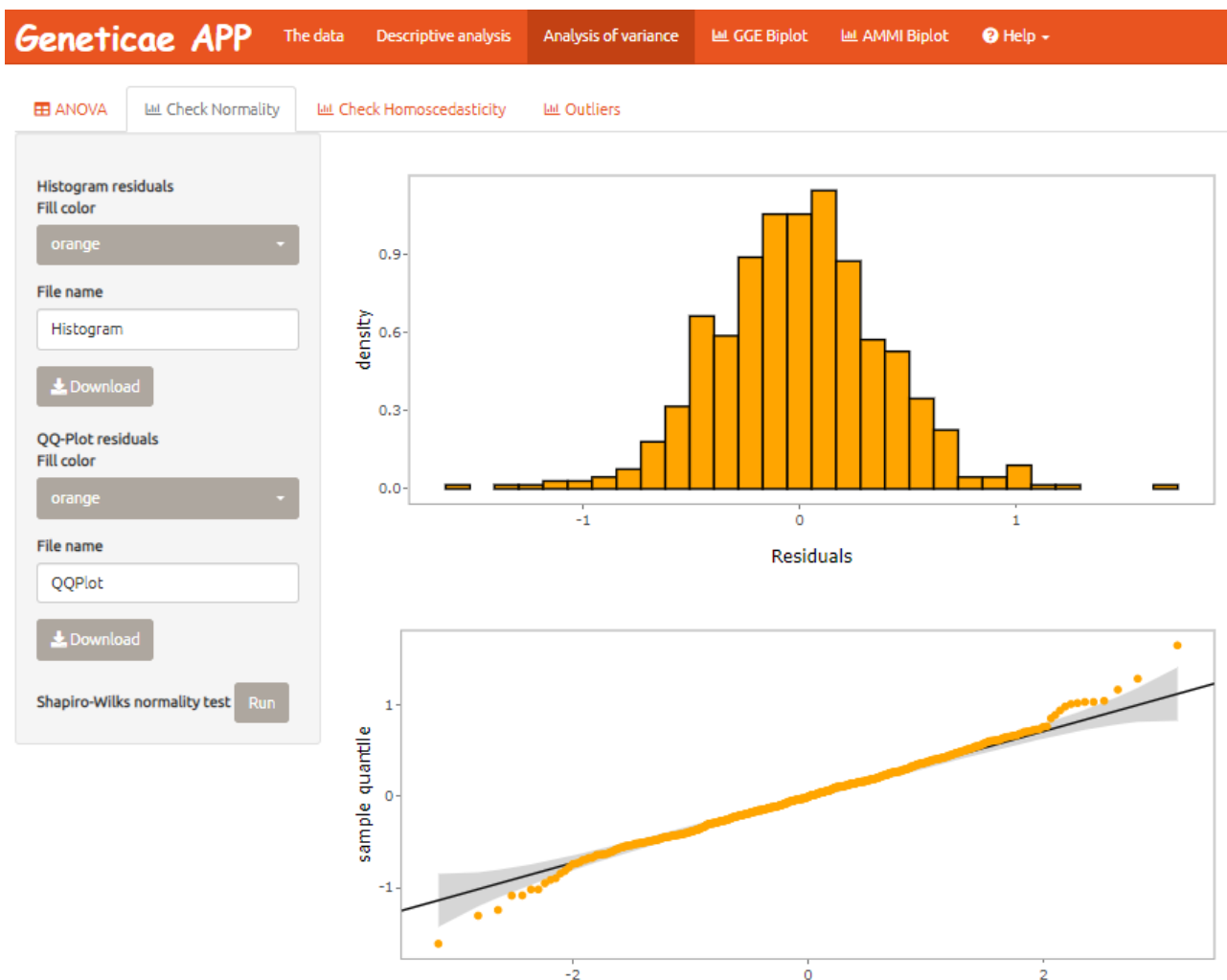


Figura 4.26: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

El grafico de residuos vs. valores predichos y las pruebas de levene permiten verificar el supuesto de variancia constante u homocedasticidad.

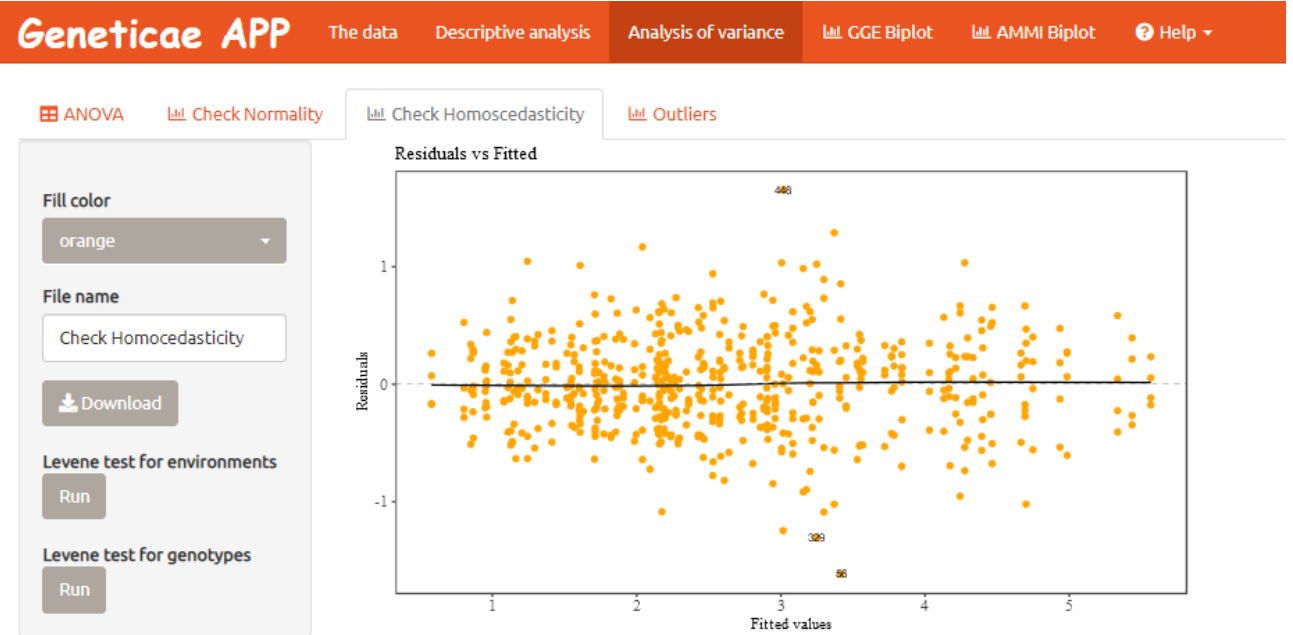


Figura 4.27: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

Por último, la presencia de observaciones atípicas u outliers provoca que el ANOVA no de buenos resultados, un grafico para detectar outliers es posible realizarlo.

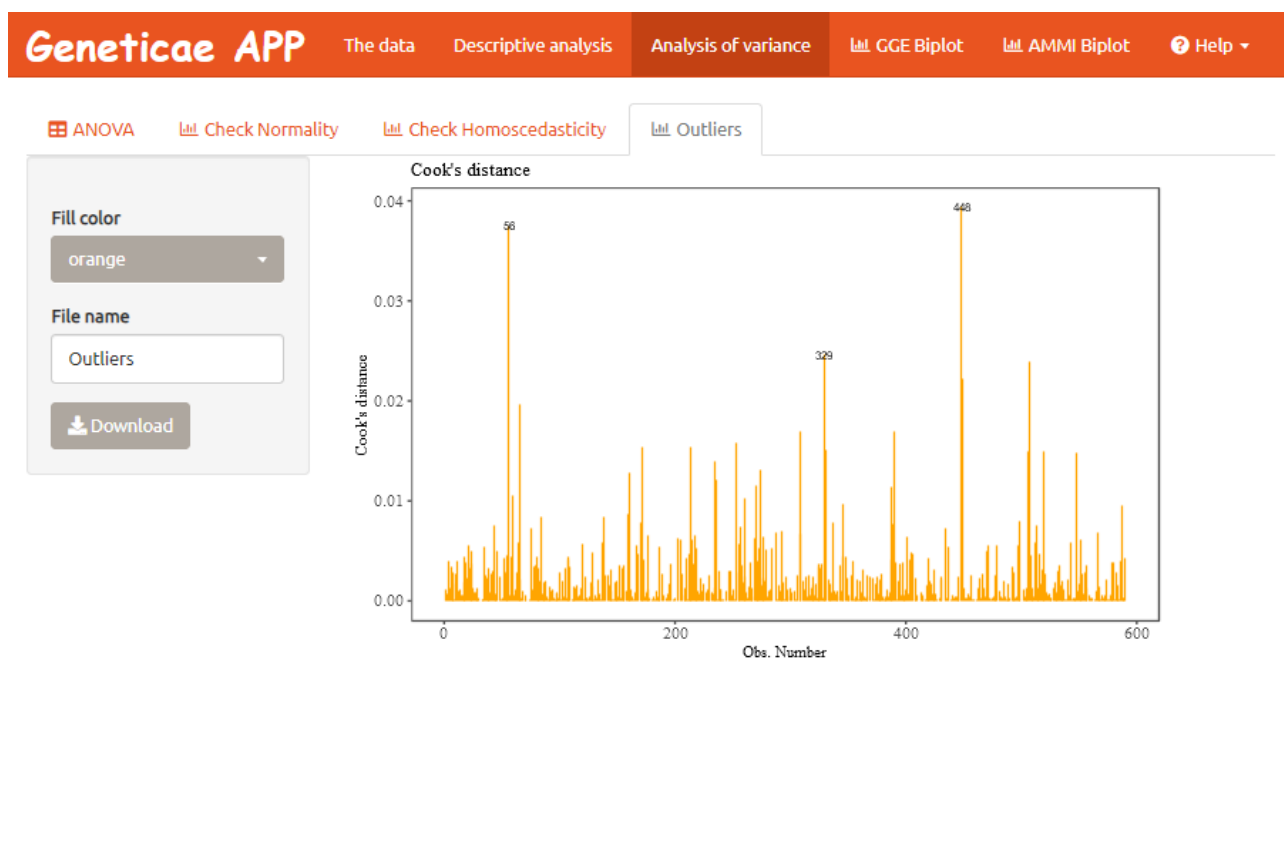


Figura 4.28: Boxplot de genotipos a través de los ambientes para el conjunto de datos Plrv

4.2.4. Biplot GGE

El biplot GGE aborda visualmente muchos problemas relacionados con la evaluación de los genotipo y ambientes de prueba. En el caso de repeticiones disponibles en el conjunto de datos, se obtiene el valor fenotípico promedio para cada combinación de genotipo y ambiente. Los valores faltantes no están permitidos.

4.2.5. Biplot GE

4.2.6. Ayuda

Capítulo 5

Conclusiones

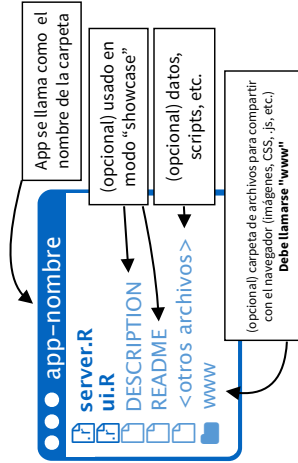
Apéndice A

Hoja de referencia Shiny

2. server.R Instrucciones que constituyen los componentes R de tu app. Para escribir server.R:

- Provee server.R con el mínimo de código necesario, **shinyServer(function(input, output) {})**.
- Define los componentes en R para tu app entre las llaves {} después de **function(input, output)**.
- Guarda cada componente R destinados para tu interfaz (UI) como **output\$<nombre componente>**.
- Crea cada componente de salida con una función **render***.
- Dale a cada función **render*** el código R que el servidor necesita para construir el componente. El servidor notará valores reactivos que aparecen en el código y reconstruirá el componente cada vez que estos valores cambian.
- Has referencia a valores en "widgets" con **input\$<nombre del widget>**.

1. Estructura Cada app es una carpeta que contiene un archivo server.R y comúnmente un archivo ui.R (opcionalmente contiene archivos extra)



server.R

```
# carga paquetes, scripts, datos
A shinyServer(function(input, output) {B
  # crea variables específicos para usuario
  output$texto <- renderText({
    input$titulo
  })
  C output$gráfica <- renderPlot({
    x <- mtcars[, input$x]E
    y <- mtcars[, input$y]
    plot(x, y, pch = 16)
  })
})
```

funciones render*

función	espera	crea
renderDataTable	objetos como tablas	tabla DataTables.js
renderImage	lista atributos imágenes	imagen HTML
renderPlot	gráfica	gráfica
renderPrint	salida impresa	texto
renderTable	objetos como tablas	tabla simple
renderText	cadena de caracteres	texto
renderUI	objeto "tag" o HTML	elemento UI (HTML)

valores de entrada (input) son reactivos.

Deben estar rodeados por uno de:

- render*** - crea un componente shiny UI (interfaz)
- reactive** - crea una expresión reactiva
- observe** - crea un observador reactivo
- isolate** - crea una copia no-reactiva de un objeto reactivo

3. Ejecución Coloca código en el lugar donde correrá la menor cantidad de veces

Corre una vez - código puesto fuera de **shinyServer** solo corre una vez cuando inicias tu app. Usalo para instrucciones generales. Crea una sola copia en memoria.

Corre una vez por usuario - código puesto dentro de **shinyServer** corre una vez por cada usuario que visita tu app (o refresca su navegador). Usalo para instrucciones que necesitas dar por cada usuario del app. Crea una copia por cada usuario.

Corre a menudo - código puesto dentro de una función **render***, **reactive**, o **observe** correrá muchas veces. Usalo solo para código que el servidor necesita para reconstruir un componente UI después de que un widget cambia.

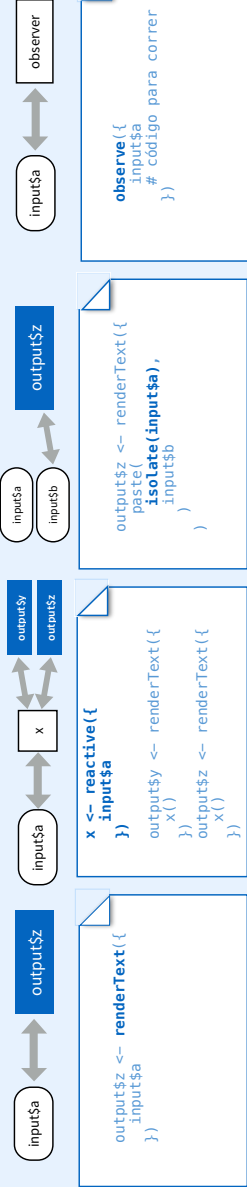
observe - usa **observe** para crear código que corre cuando una entrada cambia, pero que no crea un objeto de salida.

isolate - usa **isolate** para usar una entrada sin dependencia. Shiny no reconstruirá la salida cuando una entrada aislada cambia

reactive - usa **reactive** para crear objetos que se usaran en múltiples salidas.

render* - Una salida se actualiza automáticamente cuando una entrada en su función **render*** cambia.

4. Reactividad Cuando una entrada (input) cambia, el servidor reconstruye cada salida (output) que depende de ella (también si la dependencia es indirecta). Puedes controlar este comportamiento a través de la cadena de dependencias.



5. ui.R

```
titlePanel("datos mtcars"),
B sidebarLayout(
  sidebarPanel(
```

```
selectInput("x", "Escoge una var x:",
  choices = names(mtcars),
  selected = "disp"),
```

```
),
mainPanel(
  h3(textOutput("texto")),
  plotOutput("plot")
)
```

Componentes R - Los objetos de salida que has definido en **server.R**. Para colocar un componente:

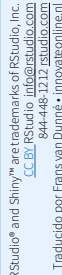
2. Pasa la función ***Output** a una cadena de caracteres correspondiente al nombre del objeto en **server.R**:

funciones *Output

dataTableOutput	tableOutput
htmlOutput	textOutput
imageOutput	uiOutput
plotOutput	verbatimOutput

runApp - corre archivos locales

runGitHub - corre archivos alojados en [www.GitHub.com](https://www.github.com)
runGist - corre archivos guardados como gist (gist.github.com)
runURL - corre archivos guardado en algún URL

Traducido por Frans van Dunné • innovateonline.nl

Presenta tu app como una página web accesible en línea

ShinyApps.io
Aloja tus apps en el servidor
de RStudio. Opciones gratis y
pagas.

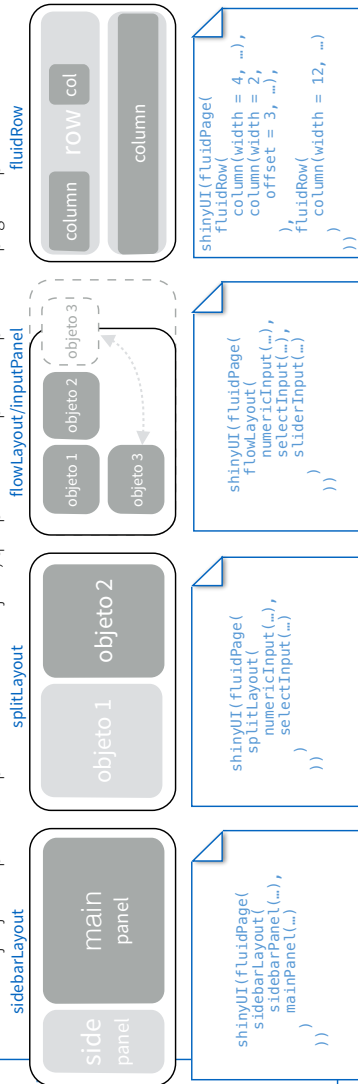
shiny.rstudio.com/deploy

Shiny Server Pro
Construye un servidor
comercial con autenticación,
gestión de recursos y más.
shiny.rstudio.com/deploy

A Incluye el mínimo de código necesario para `ui.R`, `shinyUI(fluidPage())`

* nota: usa `navigationBar` en vez de `fluidPage` si quieres que tu app tenga múltiples páginas conectados con un navbar

B Construye el plano para tu UI. `sideBarLayout` da una composición estándar cuando se usa con `sideBarPanel` y `mainPanel`. `splitLayout`, `flowLayout`, e `inputLayout` dividen la página en regiones equidistantes. `fluidRow` y `column` trabajan juntos para crear planos basados en reillas, que puedes usar para componer una página o panel.



Widgets - El primer argumento de cada función de *widget* es el **<nombre>** del widget. Puedes acceder a su valor actual en **server.R** con **input\$<nombre>**

widget	función	argumentos comunes
Botón de acción	actionButton	inputId, label
casilla	checkboxBoxInput	inputId, label, value
grupo de casillas	checkboxGroupBox	inputId, label, choices, selected
selección de fechas	dateInput	inputId, label, value, min, max, format
selección rango fechas	dateRangeInput	inputId, label, start, end, min, max, format
subarchivo	fileInput	inputId, label, multiple
campo numérico	numericInput	inputId, label, value, min, max, step
botón de selección	radioButtons	inputId, label, choices, selected
casilla de selección	selectInput	inputId, label, choices, selected, multiple
deslizador	sliderInput	inputId, label, min, max, value, step
botón de envío	submitButton	text
campo de texto	textInput	inputId, label, value



Elementos HTML - Añade elementos html con funciones shiny similares a etiquetas HTML comunes.

[illegible]

Apéndice B

Guías para usuario de Geneticae APP

Apéndice C

Código R de Geneticae APP

Bibliografía

R.W. Allard. *Principios de la mejora genética de las plantas*. Ediciones Omega, 1967.

H. G. Gauch y R. W. Zobel. Identifying mega-environments and targeting genotypes. *Crop Science*, 37:311—326, 1997.

W. Yan, L. A. Hunt, Q. Sheng, y Z. Szlavnic. Cultivar evaluation and mega-environment investigation based on the GGE biplot. *Crop Science*, 40:597—605, 2000.