



Gestión de macrodatos para el programa de mejora genética de la caña de azúcar

Big data management for sugarcane breeding program

 Reynaldo Rodríguez-Gross^{1*},  Yaquelin Puchades-Isaguirre¹,
 Wilfre Aiche-Maceo¹,  Héctor García-Pérez²

¹Estación Territorial de Investigaciones de la Caña de Azúcar Oriente Sur. Cuba

²Instituto de Investigaciones de la Caña de Azúcar (INICA). Cuba

RESUMEN: El objetivo de este trabajo fue diseñar y establecer un modelo de gestión de macrodatos para facilitar la toma de decisiones en el Programa de Mejoramiento de la caña de azúcar en Cuba e incrementar su eficiencia. Para esto se utilizaron las fuentes de información disponibles del proceso de selección del período 2000 al 2017 y las procedentes de la respuesta agroproductiva de los cultivares en áreas de producción. Se diseñó y aplicó un modelo que incluye los componentes: infraestructura, colección, validación, almacenamiento, procesamiento, análisis y visualización. Se realizó un estudio de caso del cruce C86-12 x CP70-1133 para caracterizar sus antecedentes de selección. Como resultado, el enfoque *big data* permitió obtener una compilación de datos primarios y resultados de selección, estimar el valor genético de progenitores y cruces, clasificar los cruzamientos y facilitar la toma de decisiones en el programa de mejoramiento de la caña de azúcar en Cuba para la obtención de nuevos cultivares comerciales. Su aplicación en el caso de estudio garantizó acceder a toda la información disponible al respecto y recomendar su mejor manejo.

Palabras clave: hibridación, TIC, métodos estadísticos.

ABSTRACT: The aim of this work was to design and establish a model of macro data management to facilitate decision-making in the Sugarcane Breeding Program in Cuba and increase its efficiency. For this, available sources of information from the selection process from the period 2000 to 2017 and those coming from the agroproductive response of cultivars in production areas were used. A model was designed and applied that includes components such as: infrastructure, collection, validation, storage, processing, analysis and visualization. A case study of the C86-12 x CP70-1133 cross was conducted to characterize its selection background. As a result, the big data approach made it possible to obtain a compilation of primary data and selection results, estimate the genetic value of parents and crosses, classify the crosses and facilitate decision-making in the sugarcane breeding program in Cuba to obtain new commercial cultivars. Its application in the case study guaranteed access to all available information on the subject and recommending its best management.

Key words: hybridization, ICT, statistical methods.

INTRODUCCIÓN

El empleo de técnicas y métodos de macrodatos a la agricultura constituyen una gran oportunidad para el uso de tecnologías en función de la inversión y de la percepción del valor adicional de estas dentro del sector agroalimentario (1,2). El paradigma de macrodatos, datos masivos o *big data* es bastante reciente y desempeña un papel importante en la mejora de la eficiencia de toda la cadena de suministro y en la mitigación de los problemas de seguridad alimentaria (3).

Las aplicaciones de macrodatos en la agricultura no se refieren estrictamente a la producción primaria de grandes conjuntos de datos, sino que basan las tareas de gestión en la interconexión de formas automatizadas de recopilación y almacenamiento con funciones de reconocimiento de patrones (4). Si bien el enfoque *big data* tiene en cuenta la variabilidad en el campo, necesita herramientas de reconfiguración en tiempo real para agilizar la toma de decisiones, y generalmente incluyen asistencia inteligente en la implementación, mantenimiento y uso de la tecnología (5).

*Autor para correspondencia: reynaldo.rodriguez@inicasc.azcuba.cu

Recibido: 20/01/2021

Aceptado: 25/07/2021

Este artículo se encuentra bajo los términos de la licencia Creative Commons Attribution-NonCommercial (CC BY-NC 4.0).
<https://creativecommons.org/licenses/by-nc/4.0/>



Los programas de mejoramiento genético de la caña de azúcar siguen cuatro pasos claves: (i) la generación de una gran población de progenie a partir de cruces específicos, (ii) la evaluación de esa descendencia en diferentes etapas, (iii) la selección de clones con características superiores, y (iv) recombinación de los clones seleccionados para cerrar el ciclo, o para iniciar uno nuevo. Este es un proceso relativamente largo (nunca inferior a los 10 años), que consume recursos y genera gran cantidad de información que debe ser resumida en diferentes formas (6,7).

El Instituto de Investigaciones de la Caña de Azúcar de Cuba (INICA) desarrolla un programa de mejoramiento genético para dar respuesta a la obtención de cultivares (8). Esto implica que anualmente maneja un extenso volumen de datos provenientes de una población grande de clones y cultivares en diferentes etapas de selección.

Para capturar, almacenar y procesar la información obtenida en el Programa de Mejora de la caña de azúcar en Cuba se dispone de un programa informático (SASEL) que permite mejorar la eficiencia de todo el proceso (9). Sumado a lo anterior existe información proveniente de la respuesta de los cultivares en condiciones de producción, tanto de rendimiento como de reacción frente a las plagas, que no se toma en consideración, de manera dinámica, para corregir el programa de cruzamiento de la caña de azúcar.

El objetivo de este trabajo fue diseñar y establecer un modelo de gestión con un enfoque de macrodatos asistido por la interface SASEL para facilitar la toma de decisiones en el programa de mejoramiento de la caña de azúcar en Cuba e incrementar la eficiencia de este.

MATERIALES Y MÉTODOS

Para la realización del modelo de gestión de los macrodatos del programa de mejora de la caña de azúcar se tuvo en cuenta el desarrollo de sus componentes como son:

Infraestructura y seguridad digital: El Programa de Mejora con una red de ambientes de prueba (once), que son representativos de las principales condiciones en las que se cultiva la caña de azúcar en el país (Figura 1). La misma es atendida y mantenida por siete estaciones experimentales provinciales y seis grupos de extensión y servicios agrícolas, todos dotados del personal (investigadores, técnicos y especialistas) y del equipamiento necesario (invernaderos, laboratorios de

biología molecular, diagnóstico, semilla, fisiología y análisis azucarero). Posee además un centro exclusivamente para generar variabilidad a través de la hibridación. Esta infraestructura asegura la colección de datos tanto de los nuevos estudios de cultivares como los plantados en áreas de producción de caña de azúcar.

Validación de la información: Todo el proceso de mejora genética (incluye generación de variabilidad, selección, pruebas de resistencia a plagas, extensión en áreas comerciales) se rige por Normas y Procedimientos (10), que permiten estandarizar los procedimientos y estructurar las bases de datos de la información generada, auditable por la dirección del programa, al mismo tiempo que es validada.

Colección, almacenamiento y organización de los datos: Los datos tienen un registro físico en el expediente o protocolo de cada ensayo. Para facilitar su digitalización y manejo se dispone de una aplicación informática denominada SASEL (9). La misma a nivel de provincia estratifica la información por año de selección o serie, combinación, etapa y cepa y para el análisis nacional adiciona a la provincia.

Herramientas para el análisis, modelación y visualización: A las bases de datos de selección ya sea de una serie en particular o compilación en tiempo y espacio (varias series, etapas y localidades) se le determinan parámetros genéticos-estadísticos como la media, desviación estándar, varianza y diferencial de selección.

Adicionalmente a los parámetros genéticos-estadísticos, se utilizaron los datos individuales de los cruces de una serie o compilados de varias series para determinar un índice de selección simultánea en el estimado del valor genético (EVG) de cruces y progenitores (11). Con los valores del EVG se confeccionó una metodología de evaluación de cruces para clasificar los mismos desde muy comprobados a muy descartado (12). Estas herramientas matemáticas de análisis multivariados fueron utilizadas en su conjunto para el procesamiento de la cadena de macrodatos generados en el proceso de mejora genética.

Para la visualización de los datos y sus análisis se contó con los informes predeterminados de la interface SASEL, además de las consultas personalizadas que se pueden realizar. Por otra parte, se reestructuró el modelo de relaciones de bases de datos de la plataforma informática SASEL a un nuevo enfoque de big data que incluyó las



Figura 1. Red experimental para el trabajo de mejoramiento genético de la caña de azúcar en Cuba

bases de datos relacionales, semi relacionales y no relacionales:

- Compilación de datos primarios de selección en campo de la primera y segunda etapa clonal por series y territorios.
- Compilación de resultados de selección en campo de la primera y segunda etapa clonal por series y territorios.
- Estimado del valor genético de progenitores y cruces y clasificación

Estudio de caso. Resultados de selección en la primera y segunda propagación clonal.

Como ejemplo del enfoque Big-Data se procesó una cadena de datos resultantes del proceso de selección en las etapas clonales I y II del esquema de selección vigente, correspondientes a las progenies obtenidas en el período 2000-2007, del cruce biparental entre los progenitores femenino y masculino respectivamente C86-12 y CP70-1133, evaluado en los sitios de prueba o ambientes de selección correspondientes a las provincias Villa Clara, Holguín y Santiago de Cuba.

El universo de datos abarcó 160 806 clones (no entiendo bien lo de los 1738 cruces, a menos que se haya tenido en cuenta dentro de esos cruces la participación de los progenitores C86-12 y CP70-1133) procesados con el auxilio de la aplicación SASEL que permitió obtener los resúmenes de la selección por series y territorios.

RESULTADOS Y DISCUSIÓN

El esquema de gestión de macrodatos para el programa de mejora genética de la caña de azúcar en Cuba tiene en cuenta la información que se genera en el proceso de selección de nuevos cultivares, así como las evaluaciones realizadas en áreas de producción comercial (Figura 2). Este diseño se basó en la integración de bases de datos relacionales, semirelacionales y no relaciones.

La información primaria está almacenada en diferentes bases de datos como son: programa de cruzamientos, evaluaciones, selección en campo y pruebas de resistencia. Todas estas determinaciones se obtienen por series o años, sitio de mejoramiento y etapas del esquema de selección.

Las evaluaciones de la respuesta agroproductiva de los cultivares en áreas de producción y porcentajes de áreas ocupadas, pruebas de validación comercial de nuevos cultivares y encuestas fitosanitarias también forman parte del esquema *big data*. Estas bases de datos poseen entre ellas y el programa de mejora una relación no siempre relacional o estructurada; es decir, que pueden proveerse en diferentes formatos y de diferentes fuentes.

La gestión de macrodatos permite el almacenamiento, organización y validación de los datos, para luego analizarlos, modelarlos y visualizarlos. El núcleo de análisis y procesamiento (Figura 2), está sostenido por la aplicación de análisis estadísticos, análisis de series en tiempo y espacio, modelos matemáticos y estimado del valor genético (11), así como metodologías de clasificación de cruces y progenitores (12). Con los resultados obtenidos en estos análisis se facilita la toma de decisiones en el programa de mejoramiento de la caña de azúcar en Cuba, lo que trae consigo un incremento en la eficiencia de este y en las recomendaciones de la composición de cultivares en áreas de producción comercial.

Por otra parte, se obtuvo una nueva opción de la interface SASEL para el manejo y visualización de los datos en entorno *big data* (Figura 3). Es importante destacar los cuatro elementos fundamentales que la integran. El recuadro a1 lo componen las bases de datos con el registro de progenitores y cruces a procesar. Esta información puede provenir del listado de cruce o hibridación, fichero de la serie en uso o datos introducidos manualmente.

El recuadro a2 y a3 se distingue la fuente o base datos de donde se obtiene la información de los progenitores y

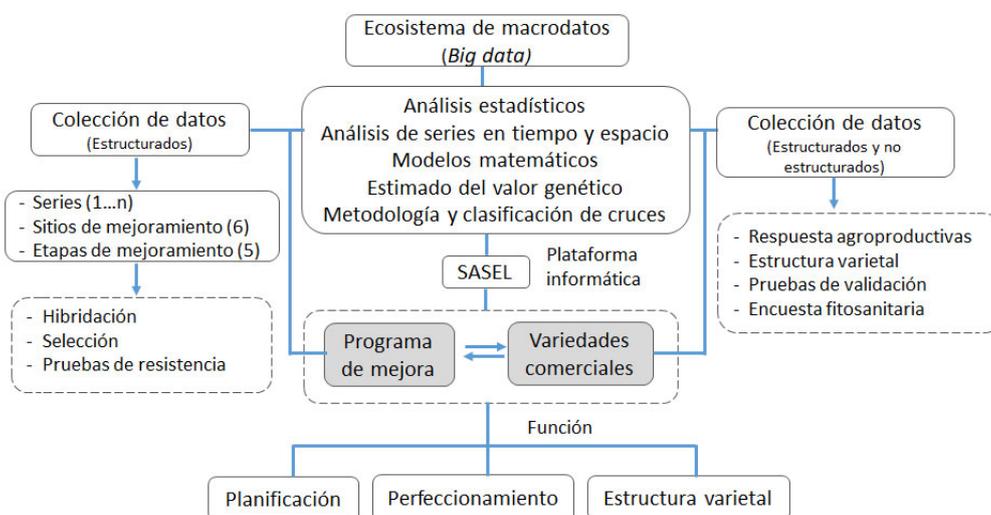


Figura 2. Esquema de gestión de macrodatos para el programa de mejora genética de la caña de azúcar en Cuba

(a1- registro de progenitores, a2- fuente de información, a3- Datos colectados o procesados, a4- Salidas del informe)

Figura 3. Interface para el manejo y visualización de macrodatos

cruces seleccionados, y si esta información es primaria o procesada previamente con la aplicación análisis estadísticos y modelos matemáticos. Las bases de datos que se encuentran integradas en esta opción son: (i) etapa de posturas; (ii) etapa de lotes clonal 1 y 2; (iii) estudios replicados; (iv) valor genético de los cruces; (v) variedades comerciales y (vii) pruebas de resistencia que incluye las encuestas fitosanitarias.

El recuadro a4 revela el tipo de salida y el alcance de la información consultada. Es posible realizar análisis individuales por progenitores, serie, territorio o etapa y también se pueden obtener resúmenes integrales. Es en esos últimos donde mejor se aprecia el enfoque *big data* pues la información se visualiza resumida lo que facilita su comprensión y con ello la toma de decisiones.

Esta interconexión de bases de datos ha permitido obtener una compilación de resultados sobre la cantidad de cruzamientos evaluados por el programa de mejoramiento genético, su frecuencia y sus resultados. También se obtuvo un modelo matemático para estimar el valor genético de progenitores y cruces (11), y el mismo constituye la base de la metodología para la clasificación de los cruzamientos (12). Este último resultado tiene en consideración la interacción genotipo ambiente, lo que sin dudas constituye una valiosa herramienta para incrementar la eficiencia del programa de mejora.

La incorporación de nuevas bases de datos provenientes de evaluaciones en áreas de producción comercial brinda la posibilidad de nuevas salidas en función de mejorar las recomendaciones de la composición de cultivares y su manejo.

Caso de estudio. Antecedentes de resultados de selección y valor genético en etapas clonales 1 y 2

La consulta realizada sobre los antecedentes de selección en las bases de datos de los lotes clonales 1 y 2 del cruce C86-12 x CP70-1133 mostró los resultados de su

participación en el período de estudio de las series del 2000 al 2017 (Tabla 1). Esta salida es producto de la recopilación, procesamiento, análisis e interconexión de seis cadenas de datos:

a1) Datos colectados en las evaluaciones de selección de las etapas clonales 1 y 2. En el período de estudio, solamente se encontraron siete series donde participó el cruce C86-12 x CP70-1133, lo que representa la gestión de datos en dos etapas, tres provincias o territorios (Villa Clara, Holguín y Santiago de Cuba) y tres años o series (2008, 2009 y 2011).

a2) Muestra la reacción en pruebas estatales de resistencia de las plagas roya parda, carbón, escaldadura foliar y virus de la caña de azúcar de los cultivares C86-12 y CP70-1133 utilizados como progenitores femenino y masculino, respectivamente. Las categorías de reacción a estas cuatro enfermedades se resumen en: susceptible (S), Intermedia (I) y Resistente (R).

a3) Visualiza los resultados de la clasificación del cruce (caso de estudio) almacenado en una base de datos, previa compilación de resultados de selección, aplicación de un modelo matemático o índice simultáneo de selección y algoritmo de clasificación, donde se obtiene el valor genético o clasificación del cruce en ocho categorías (13): 0-sin evaluación; 1-muy descartada; 2-moderadamente descartada; 3-descartada; 4-exploratorio; 5-moderadamente comprobado; 6-comprobado y 7-muy comprobado. En este caso el cruce C86-12 x CP70-1133 resultó muy comprobado y comprobado en las provincias Villa Clara y Santiago de Cuba, respectivamente y sin evaluación en la provincia Holguín.

a4) Refleja los resúmenes de selección de las siete bases de datos colectadas en cada serie, etapas y territorio en que participó. En la misma se pudo determinar que en el

y *big data* para facilitar la toma de decisiones y contribuir al mejoramiento de la agroindustria azucarera. Lo más importante es que el enfoque descrito en este documento se puede extender fácilmente a otras áreas de investigación del cultivo e industrias agrícolas para recomendar mejores prácticas agrícolas.

CONCLUSIONES

- Se obtuvo un modelo de gestión con un enfoque de macrodatos, asistido por la interface SASSEL, para facilitar la toma de decisiones en el Programa de Mejoramiento de la caña de azúcar en Cuba e incrementar la eficiencia de este.
- A través de este modelo de gestión de información se pudo obtener una compilación de datos primarios, resultados de selección, estimar el valor genético de progenitores y cruces y clasificar los cruzamientos utilizados en el Programa de Mejoramiento.

RECOMENDACIONES

Extender el uso del enfoque de macro datos a partir de la compilación de otras fuentes de información que contribuyan a incrementar la eficiencia del programa obtención de nuevos cultivares comerciales.

BIBLIOGRAFÍA

1. Sun J, Zhou Z, Bu Y, Zhuo J, Chen Y, Li D. Research and development for potted flowers automated grading system based on internet of things. *Journal of Shenyang Agricultural University* [Internet]. 2013;44(5):687-91. Available from: <https://www.cabdirect.org/cabdirect/abstract/20133406244>
2. Yang C. Big Data and its potential applications on agricultural production. *Crop, Environment & Bioinformatics* [Internet]. 2014;11(1):51-6. Available from: <https://www.cabdirect.org/cabdirect/abstract/20143190953>
3. Chen M, Mao S, Liu Y. Big data: A survey. *Mobile networks and applications* [Internet]. 2014;19(2):171-209. Available from: <https://dl.acm.org/doi/10.1007/s11036-013-0489-0>
4. Mazzocchi F, Lapointe FJ. Sobre el 'big data': ¿Cómo podríamos dar sentido a los macrodatos? *Métode: Revista de difusión de la Investigación* [Internet]. 2020;1(104):34-41. Available from: <https://dialnet.unirioja.es/servlet/articulo?codigo=7391781>
5. Wolfert S, Ge L, Verdouw C, Bogaardt M-J. Big data in smart farming-a review. *Agricultural systems* [Internet]. 2017;153:69-80. Available from: <https://www.sciencedirect.com/science/article/pii/S0308521X16303754>
6. Park S, Jackson P, Berding N, Inman-Bamber G. Conventional breeding practices within the Australian sugarcane breeding program. In: *Proceedings of the Australian Society of Sugar Cane Technologists* [Internet]. 2007. p. 113-21. Available from: https://www.researchgate.net/profile/Nils-Berding/publication/305390103_Conventional_breeding_practices_within_the_Australian_sugarcane_breeding_program/links/578c603608ae5c86c9a14e99/Conventional-breeding-practices-within-the-Australian-sugarcane-breeding-program.pdf
7. Yadav S, Jackson P, Wei X, Ross EM, Aitken K, Deomano E, et al. Accelerating genetic gain in sugarcane breeding using genomic selection. *Agronomy* [Internet]. 2020;10(4):585. Available from: <https://www.mdpi.com/2073-4395/10/4/585>
8. González R. Variedades de caña de azúcar cultivadas en Cuba. *Cronología, legislación, metodologías y conceptos relacionados*. Instituto Cubano de Investigaciones de Derivados de La Caña de Azúcar; 2019.
9. Rodríguez R, Puchades Y, Abiche W, Rill S, García H. SASSEL: software for data management generated in the Cuban sugarcane-breeding program. In: *Proceedings of the International Society of Sugar Cane Technologists*. 2016. p. 63-6.
10. Jorge H, González R, Casas M, Jorge I. Normas y procedimientos del programa de mejoramiento genético de la caña de azúcar en Cuba. *PUBLINICA*, La Habana. 2011;
11. Rodríguez-Gross R, Puchades-Isaguirre Y, Aiche-Maceo W, Cornide-Hernández MT. Modelo matemático para estimar el valor genético de progenitores y cruces en caña de azúcar. *Cultivos Tropicales* [Internet]. 2018;39(2):81-8. Available from: http://scielo.sld.cu/scielo.php?pid=S0258-59362018000200011&script=sci_art-text&tlang=en
12. Rodríguez-Gross R, Puchades-Isaguirre Y, Aiche-Maceo W. Metodología de validación y manejo de cruces en la mejora genética en caña de azúcar. *Cultivos Tropicales* [Internet]. 2020;41(1). Available from: http://scielo.sld.cu/scielo.php?script=sci_arttext&pid=S0258-59362020000100002
13. Kumar H, Menakadevi T. A review on big data analytics in the field of agriculture. *International Journal of Latest Transactions in Engineering and Science* [Internet]. 2017;1(4):1-10. Available from: <http://www.ijltes.com/wp-content/uploads/2017/02/1.pdf>
14. Everingham Y, Sexton J, Skocaj D, Inman-Bamber G. Accurate prediction of sugarcane yield using a random forest algorithm. *Agronomy for sustainable development* [Internet]. 2016;36(2):27. Available from: <https://link.springer.com/content/pdf/10.1007/s13593-016-0364-z.pdf>
15. Everingham Y, Sexton J, Robson A. A statistical approach for identifying important climatic influences on sugarcane yields. In: *Proceedings of the 37th Conference of the Australian Society of Sugar Cane Technologists*, 28-30 April 2015, Bundaberg, Queensland, Australia [Internet]. Australian Society of Sugar Cane Technologists; 2015. Available from: https://www.researchgate.net/publication/287208302_A_Statistical_Approach_for_identifying_Important_Climatic_Influences_on_Sugarcane_Yields
16. Biqing L, Yongfa L, Miao T, Shiyong Z. Design and Implementation of Sugarcane Growth Monitoring System based on RFID and ZigBee. *International Journal of Online Engineering* [Internet]. 2018;14(3). Available from: https://www.researchgate.net/publication/324114641_Design_and_Implementation_of_Sugarcane_Growth_Monitoring_System_based_on_RFID_and_ZigBee