

Propuesta de un algoritmo híbrido para la corrección de imágenes RAPD

Ricardo Contreras A., M. Angélica Pinninghoff J. y Alejandra Alvarado V.

Departamento de Ciencias de Computación,
Universidad de Concepción, Chile
{rcontrer,mpinning}@udec.cl

Resumen Este artículo describe una experiencia que combina Algoritmos Genéticos y Simulated Annealing como un mecanismo de corrección de carriles en imágenes ADN obtenidas a través de la técnica RAPD (Random Amplified Polymorphism DNA). Las imágenes RAPD son afectadas por diversos factores; entre ellos el ruido y la distorsión que impactan la calidad de las imágenes, y consecuentemente, la precisión en la interpretación de los datos. Este trabajo propone un método híbrido que emplea algoritmos genéticos, para tratar con la característica altamente combinatorial de este problema, y simulated annealing, para tratar con los óptimos locales. Los resultados obtenidos con esta aproximación sobre este problema particular muestran una mejora, respecto del uso exclusivo de algoritmos genéticos, tanto en el fitness, o calidad de los individuos, como en los tiempos de ejecución.

Palabras clave: Imágenes RAPD, cocción simulada (*simulated annealing*), algoritmos genéticos.

1. Introducción

RAPD (Randomly Amplified Polymorphism DNA) [11] es un tipo de marcador molecular que se utiliza para la verificación de identidad genética. En los últimos años, se ha utilizado RAPD para estudiar relaciones filogenéticas [1][10], mapeamiento genético [4], marcadores asociados a rasgos [9], y mapeamiento de enlaces genéticos [2]. Esta técnica ha sido usada como apoyo en muchos programas asociados a la agricultura, al sector forestal y a la reproducción de animales [5].

En la figura 1, se muestra la fotografía de una imagen RAPD. En este caso se consideraron 12 muestras, usando los carriles 1 y 14 para indicar pesos moleculares estándar. Cada uno de los carriles mostrados en la imagen está compuesto por un conjunto de trazos de diferente brillo, que corresponden a las bandas de cada carril. Se estudiaron cuatro diferentes genotipos de *Eucalyptus globulus*, incluyendo tres copias idénticas de cada genotipo (conocidos como rametos). Si los rametos son idénticos, entonces es de esperar que se obtengan patrones de bandas similares, al ser analizados en las mismas condiciones. Sin embargo, esto

no ocurre debido a, por ejemplo, errores en la rotulación de las muestras u otro tipo de deformaciones producidas durante el proceso.

La técnica RAPD consiste en la amplificación randómica de secuencias del ADN genómico usando *primers* (cebadores), los cuales comunmente tienen una longitud de 10 pb (pares base). Este proceso es llevado a cabo por la reacción en cadena de la polimerasa, conocido como PCR (polymerase chain reaction) y genera un patrón típico para una sola muestra y diferentes *primers*. Un *primer* es una cadena de ácido nucleico que sirve como punto de partida para la síntesis del ADN. Los productos PCR son separados en un gel de agarosa, lo que permite que, bajo la acción de un campo eléctrico, los fragmentos más pequeños de los productos PCR se desplacen más rápido, mientras que los más grandes lo hagan más lentamente. El gel es coloreado con una tintura (típicamente bromuro de etidio) y fotografiado para el análisis posterior de los datos. Una manera de analizar la imagen obtenida es la simple comparación visual de las diferentes bandas que se obtiene, para cada muestra. Sin embargo, puede ser un proceso tedioso cuando varias muestras con diferentes combinaciones de *primers* deben ser analizadas. Al mismo tiempo, ya que en este caso la presencia o ausencia de bandas es cuantificada, esa cuantificación puede ser subjetiva y no hay un nivel de umbral confiable, ya que las intensidades de las bandas son afectadas por varios factores (tintura, calidad del gel, reacción PCR, calidad del ADN, entre otros).

Durante el proceso de generación de una imagen RAPD, muchos factores físico-químicos afectan la electroforesis produciendo diferentes tipos de ruido, rotaciones, deformaciones y otras distorsiones anormales en la imagen. El efecto de esto es, desafortunadamente, propagado a las diferentes etapas en el análisis posterior, incluyendo visualización, extracción del background, detección de bandas y clustering, que puede conducir a conclusiones biológicas erróneas. Por esta razón, técnicas eficientes de procesamiento de imágenes pueden tener un impacto positivo sobre esas conclusiones biológicas.

Errores típicos consideran rotación de carriles, que es el problema que intentamos resolver. Este es el primer paso; una vez que los carriles son corregidos, es decir, la imagen completa muestra pendientes mínimas para cada carril, es necesario trabajar en la corrección de las bandas, un problema difícil debido a la naturaleza de las distorsiones, diferente a la distorsión de carriles. Este segundo paso, la corrección de las bandas, puede ser atacado de manera análoga; sin embargo, va más allá del alcance de este artículo.

La base para este trabajo es la experiencia descrita en [7], donde se usan algoritmos genéticos para tratar con una población de soluciones potenciales, donde las soluciones se entienden como los mejores *templates* (esquemas). Un *template* es un conjunto de líneas que representan carriles. Así, el mejor *template* es uno en el cual las líneas presentan pendientes que se aproximan de manera más cercana a las pendientes que aparecen en la imagen original. Este trabajo presentó buenas soluciones, aunque los tiempos de ejecución eran superiores a

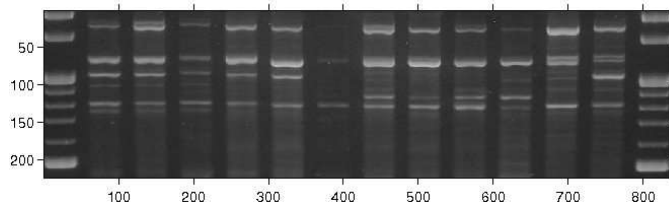


Figura 1. Una muestra de imagen RAPD con dos carriles de referencia, y 12 carriles que representan a cuatro rametos.

los esperados. El otro problema que este enfoque presenta, es la presencia de mínimos locales.

El objetivo de este trabajo es corregir las distorsiones presentes en los carriles, usando algoritmos genéticos híbridos con *Simulated Annealing* (Recocido simulado). Esto permite la comparación de dos estrategias: la primera, que considera solamente algoritmos genéticos, y la segunda, que usa algoritmos genéticos y simulated annealing como mecanismos colaborativos.

Este artículo está estructurado de la siguiente forma; la primera sección consta de la presente introducción; la segunda sección describe el problema específico a ser enfrentado; la tercera sección está dedicada a consideraciones sobre algoritmos genéticos y simulated annealing, mientras la cuarta sección muestra los resultados obtenidos con nuestra propuesta. La sección final presenta las conclusiones del trabajo.

2. El enfoque propuesto

El problema enfrentado en este trabajo puede ser formalmente establecido de la siguiente forma.

Considérese una imagen (una matriz) $A = \{a_{ij}\}, i = 1, \dots, n$ y $j = 1, \dots, m$, donde $a_{ij} \in Z^+$, y A es una imagen RAPD. Usualmente, a_{ij} está en el rango $[0, 255]$ en una imagen en escala de grises, y se utiliza a_{ij} para denotar a un elemento $A(x, y)$, donde x e y son las coordenadas del pixel.

Para tratar las distorsiones de carriles, se usa un conjunto de *templates*. Estos *templates* son líneas creadas en forma random que presentan diferentes grados de distorsión, y que están en correspondencia uno-a-uno con los carriles en la imagen RAPD original. Un buen *template*, al igual que lo considerado en [7], es aquel que refleja de manera más precisa las distorsiones que tiene la imagen RAPD bajo consideración.

El *template* es representado por una matriz L (los carriles) donde $L = \{l_{ij}\}, i = 1, \dots, n$ y $j = 1 \dots, m; l_{ij} = 0$ o $l_{ij} = 1$ (una imagen binaria), con 1 indicando que l_{ij} pertenece a una línea y 0 en caso contrario. Un procedimiento descrito en [8] se usa para detectar en forma aproximada la posición inicial de los carriles. De esta forma la generación de la matriz L se limita a aquellas regiones que corresponden a carriles en la matriz A . Debido a la rotación de los carriles, es

necesario considerar diferentes configuraciones alternativas. Si se está tratando con una imagen de 12 carriles, y si por cada línea se consideran 14 rotaciones posibles, entonces se está hablando de 12^{14} configuraciones diferentes a evaluar. Esto provoca una explosión combinatoria, que es lo que justifica el uso de una técnica como los algoritmos genéticos.

Los algoritmos genéticos permiten manipular un gran número de *templates*, escogiendo aquellos que son similares a la imagen original. Así, es necesario encontrar una función objetivo que refleje esta similitud de una manera precisa. Esta función es usada como una medida de la calidad del *template* seleccionado.

Cuando se aplica el procedimiento de corrección de carriles, los *templates* contienen líneas rectas. Diferentes *templates* mostrarán diferentes pendientes para cada línea, como se muestra en la figura 2. Un *template* contiene líneas verticales que no se intersectan, pero que no son necesariamente paralelas (en el marco restringido a la imagen).

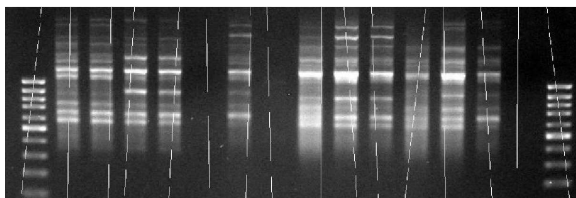


Figura 2. Un *template* de muestra para la corrección de carriles.

Los resultados obtenidos en trabajos previos eran prometedores, pero no lo suficientemente buenos como para ser aceptados como una alternativa. Con esta consideración en mente, se decidió hibridar la estrategia solución incorporando la técnica de simulated annealing. Simulated annealing es un procedimiento para optimizar una función en base a perturbaciones randómicas de una potencial solución y a una decisión probabilística sobre la posibilidad de mantener la solución mutada [3].

3. Algoritmos genéticos y *simulated annealing*

Los algoritmos genéticos (AG) son una clase particular de algoritmos evolutivos, usados para encontrar soluciones óptimas, o buenas, examinando sólo una fracción reducida del espacio de soluciones posibles. Los AGs están inspirados en la teoría de la evolución de Darwin.

En lo fundamental, los AGs se diseñan para simular procesos necesarios para la evolución en sistemas naturales, específicamente aquellos que siguen los principios de supervivencia del más adecuado. De esta forma, representan una explotación inteligente de una búsqueda randómica dentro de un espacio de búsqueda definido para resolver un problema.

La estructura de un algoritmo genético consiste de un procedimiento iterativo simple sobre una población de individuos genéticamente diferentes. Los fenotipos son evaluados de acuerdo a una función de fitness predefinida; los genotipos de los mejores individuos son copiados varias veces y luego modificados por medio de los operadores genéticos; los nuevos genotipos obtenidos en este proceso son introducidos en la población, en reemplazo de los antiguos. Este procedimiento continúa hasta que una solución *suficientemente buena* es encontrada [3].

En este trabajo, los *templates* son los cromosomas, las líneas en los templates son los genes y una línea con una pendiente particular representa el valor (alelo) que tiene un gen.

Un buen fitness indica que un *template* particular (matriz L) tiene una mayor similitud con la imagen RAPD original (matriz A).

Para evaluar un *template*, las imágenes que corresponden a las matrices A y L se superponen, obteniéndose una suma de intensidades que considera los píxeles de una vecindad dentro de un rango para cada línea. En este trabajo, este rango se determina considerando el ancho de la porción más brillante del carril. El objetivo de esto es obtener mayor precisión en la función fitness. Si una línea en el *template* coincide con un carril, se obtiene un mayor valor para la suma. En contraste con lo anterior, si no existe coincidencia, el valor es menor que en el primer caso, porque se está agregando intensidades de píxeles que corresponden al fondo de la imagen (valores cercanos a cero).

Un aspecto interesante en el presente enfoque resulta del hecho que el valor obtenido para la evaluación de cada línea es almacenado como parte del *gen*. De esta forma, la suma de las intensidades se hace solamente al momento de crear una nueva línea. Como consecuencia de esto, el tiempo de ejecución se reduce de manera considerable, en comparación con experiencias previas.

Operadores genéticos: Para este trabajo se consideraron diferentes operadores genéticos. Estos operadores se describen brevemente a continuación:

- **Selección.** La selección se lleva a cabo utilizando el mecanismo de *ruleta* [3]. Esto significa que los individuos con mejores valores de fitness tendrán mayor probabilidad de ser escogidos como padres en el proceso de reproducción.
- **Cruzamiento.** El cruzamiento se usa para intercambiar material genético, permitiendo que parte de la información genética de un individuo se combine con parte de la información genética de un individuo diferente. Por ejemplo, si se tiene dos *templates*, cada uno de ellos con $r + s$ líneas, después del cruzamiento los hijos generados resultan en: hijo 1 tendrá las primeras r líneas provenientes del *template* 1, mientras que las siguientes s líneas provendrán del *template* 2. Para el hijo 2, el proceso es levemente diferente, en cuanto a que el orden de los *templates* considerados se modifica.
- **Mutación.** Al usar este operador genético, se introduce una pequeña variación en la población de manera que se crea nuevo material genético. En este trabajo, la mutación se realiza reemplazando en forma aleatoria, con baja probabilidad, una línea particular en un *template*.

Simulated Annealing (SA) es una técnica que se inspira en la obtención de metales de diversas formas, a partir del metal en estado líquido, a través de un proceso de enfriamiento gradual, permitiendo que el material transite desde un estado inestable, desordenado, de alta energía, a un estado estable, ordenado y de baja energía.

En SA, el material es una solución candidata. Esta solución sufre una mutación y, si su energía es menor que la energía de la solución en el estado previo, la solución mutada reemplaza a la original. Si, por otra parte, la energía es mayor, entonces se produce el reemplazo de la anterior dependiendo de una probabilidad que es proporcional a la diferencia de energías. Inicialmente, cuando la temperatura del sistema es alta, las soluciones mutadas con energía relativamente altas (bajo fitness) tienen alguna posibilidad de ser mantenidas en el espacio de soluciones. La temperatura del sistema es disminuida después de n evaluaciones, reduciendo de forma efectiva la probabilidad de mantener soluciones mutadas con valores de energía altos [3].

Para este trabajo se hace uso de una variante del algoritmo SA que permite que se produzca *reannealing* (aumento de la temperatura), ya que entrega más libertad en la búsqueda por soluciones, tanto en tiempo como en espacio. El problema de encontrar la inclinación de carriles a través de SA consiste en, dada una determinada temperatura T , buscar durante a lo más K iteraciones una cantidad S de soluciones. Cada configuración o posible solución es generada a través de una determinada función, a partir de una solución encontrada durante la iteración anterior. La solución de partida es aquella proporcionada por la salida del AG, un individuo representado por un conjunto de rectas donde cada una se ajusta en cierta medida a la inclinación de cada carril de la imagen RAPD bajo análisis. Cada vez que se crea una nueva configuración ésta es evaluada; de ser mejor que la antigua, esta última se desecha, mientras que si ocurre lo contrario, se acepta bajo cierta probabilidad, que es proporcional a la temperatura.

Un parámetro, asociado a la temperatura, da cuenta de la magnitud del movimiento a realizar, a mayor temperatura, los cambios son más drásticos; a menor temperatura, los cambios son más sutiles.

Cuando el proceso finaliza, se espera seleccionar una solución donde el conjunto de rectas asociado a la imagen representa la inclinación de los carriles de mejor manera que la entregada por la configuración de partida. A partir de una configuración x , durante la aplicación de SA se producen dos posibles movimientos, para generar soluciones candidatas que representan un aumento o una disminución de la energía.

Una solución es un *template* que representa una imagen RAPD, la solución x , que debe transformarse en una solución x' a través de una modificación del *template*. Las dos posibilidades de modificación del template recién mencionadas son: en primer lugar, el cambio en la pendiente de una o más líneas del *template*, es decir, un movimiento de rotación; en segundo lugar un movimiento de traslación, es decir, una línea es movida hacia la izquierda o derecha en el

template, sin alterar su pendiente. En otras palabras, si llamamos x_{inf} e x_{sup} a los valores de los puntos inferiores y superiores en una línea, respectivamente, un movimiento de rotación se consigue cambiando los valores de esos puntos a $x_{inf} - \delta$ y $x_{sup} + \delta$ (o modificando a $x_{inf} + \delta$ y $x_{sup} - \delta$). De manera análoga, el movimiento de traslación se consigue desplazando los puntos originales a $x_{inf} + \delta$ y $x_{sup} + \delta$ (cambiando el signo si el movimiento es hacia el otro lado de la imagen). La Figura 3 ilustra el movimiento de rotación y el movimiento de traslación. Los valores permitidos tanto para la rotación como para el desplazamiento son disminuidos consecuentemente con la disminución de la temperatura, para evitar cambios severos en la calidad de las soluciones. Esto pensando en que temperaturas más bajas se acercan más a óptimos globales.

La elección de cuál movimiento realizar se toma aleatoriamente, aunque es posible definir un parámetro que permita privilegiar una opción por sobre la otra. En la rotación y el desplazamiento de las líneas se considera un valor máximo para δ de 10 píxeles.

El criterio de parada se definió en función de un cierto número de iteraciones. Esta decisión se debió a que la elección de una temperatura (mínima) como argumento, se podía distorsionar debido a eventuales aumentos de temperatura, realizados durante el proceso (reannealing, como se mencionó más arriba).

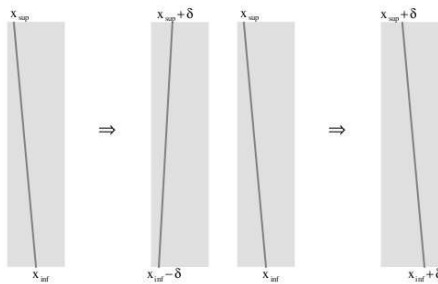


Figura 3. Esquema para rotación y desplazamiento de una línea.

4. Pruebas y resultados

Para verificar la hipótesis original, se construye un prototipo basado en AGs, hibridado con SA, obteniéndose mejoras en fitness y en tiempo de ejecución con respecto a trabajos previos, como el presentado en [6].

El problema de la detección de la inclinación de carriles se aborda haciendo uso de un algoritmo híbrido que combina ambos métodos mencionados. Se plantean dos estrategias de hibridación, una simple (HS) y una asistida (HA). La primera plantea partir buscando una solución con AG, la cual al no presentar

mejoras durante un cierto número de iteraciones (generaciones en AG), rescata al mejor individuo alcanzado hasta el momento y lo utiliza como entrada a SA, quien finalmente explora mejores soluciones. La segunda estrategia también comienza aplicando AGs y, al igual que la estrategia anterior, pasado un número de generaciones sin mejora, activa el componente de SA. SA recibe como entrada un individuo y le retorna al AG una cantidad de soluciones (no solamente la mejor solución, sino que además una parte de todas aquellas soluciones que aceptó en su ejecución) que pasarán a reemplazar parte de la población y provocarán la variabilidad de la misma. Este ciclo se repite hasta que AG completa una cantidad límite de generaciones.

Los parámetros utilizados para el algoritmo genético se muestran en la Tabla 1.

Tabla 1. Parámetros usados en el algoritmo genético.

Parámetros fijos	
Número generaciones	1000
Elitismo	10 %
Cruzamiento	70 %
Parámetros variables	
Mutación	2 %, 4.5 %, 6 %
Tamaño población	20, 50, 100 individuos

Y los parámetros utilizados para simulated annealing se muestran en la Tabla 2. En esta tabla, GSM (Generaciones sin mejora) corresponde a la condición de activación del proceso de simulated annealing, a partir del número de generaciones del algoritmos genético en que no se evidencian mejoras.

Tabla 2. Parámetros usados para simulated annealing.

Parámetros fijos	
α (factor de enfriamiento)	0.95
η (factor de reannealing)	1.5
Número de iteraciones	40
Parámetros variables	
Temperatura inicial	50, 100, 150, 200
GSM	100, 500

La siguiente secuencia de imágenes permite mostrar cómo evoluciona el proceso de corrección, a partir de la imagen a la que se le han superpuesto las líneas correspondientes al *template* creado usando AGs.

Dados los resultados de HA y HS, se tiene que ambos enfoques son capaces de provocar mejoras tanto en fitness como en tiempo por sobre los resultados

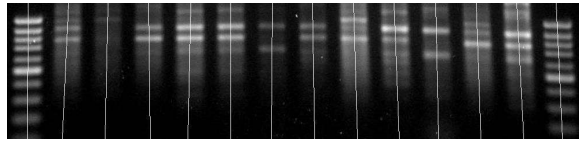


Figura 4. Template producido por la acción del algoritmo genético.

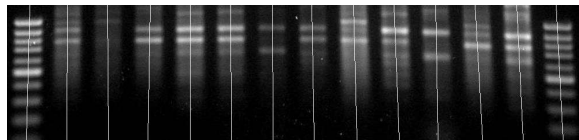


Figura 5. Template modificado por la acción de la hibridación asistida.

obtenidos con AG, y la magnitud de esta mejora varía dependiendo de la combinación de valores que tomen los parámetros en cuestión, razón por la cual a veces los mejores resultados se asocian a HA y otra veces a HS.

El conjunto de pruebas consideró una familia de 13 imágenes, seleccionadas como elementos representativos de una muestra disponible de 150 imágenes. Esto quiere decir que, a juicio de un experto, las imágenes seleccionadas contienen todos los atributos que resultan de interés en todas las imágenes de la muestra.

Aplicando HA, las mejoras se obtuvieron para un 85% de la muestra de imágenes, mientras que con HS se obtuvieron mejoras para un 62%. Un aspecto importante a señalar tiene que ver con la forma en que se fue alcanzando la mejora final de la función objetivo en cuestión: Con HA, la totalidad del valor de la función objetivo fue creciendo a lo largo de casi la totalidad de sus iteraciones. No así con HS, por el cual la totalidad de la función objetivo se enfoca en la partida, y en iteraciones posteriores crece pero más lentamente.

5. Conclusiones

El tratamiento de corrección de imágenes RAPD a través del uso de algoritmos evolutivos, ya se había enfrentado anteriormente con una hibridación entre

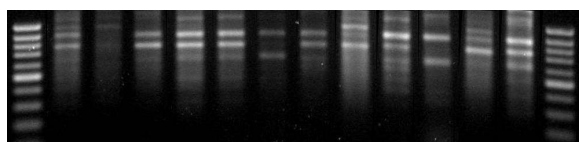


Figura 6. Imagen corregida en base al template generado.

algoritmos genéticos y la técnica de tabu search. El trabajo que aquí se ha desarrollado pretende ampliar el análisis en el uso de mecanismos evolutivos para resolver este problema. Tanto el tratamiento con algoritmos genéticos operando aisladamente, como la hibridación de algoritmos genéticos, ya sea con simulated annealing o con tabu search, se encuentran todavía en fase experimental.

La propuesta aquí explicada permite automatizar el proceso de corrección, con una disminución en tiempo que va desde 8 horas a los 25 minutos (tiempos aproximados), para el tratamiento de una sola imagen. Más allá del tiempo, la calidad de los resultados es el aspecto relevante a considerar, y lo que justifica el esfuerzo en la búsqueda de alternativas efectivas y replicables.

Simulated annealing es capaz de producir mejoras sobre los resultados obtenidos con AG, particularmente con el caso de HA, que cumple con la tarea de sacar constantemente a AG de eventuales óptimos locales. Para la corrección de imágenes RAPD se utilizan las configuraciones asociadas a HA puesto que este presentó mejoras para una mayor cantidad de imágenes que con HA.

Cabe señalar que uno de los factores más complejos de definir fue la temperatura, puesto que fue necesario construir un mecanismo para calcular la temperatura asociada a cada imagen según sus características. Esto último destaca la importancia que posee lograr definir aspectos que caractericen a una imagen, pues en base a eso se pueden generalizar formas de buscar mejoras para distintos casos que se presenten.

Referencias

1. Cao, W., Scoles, G., Hucl, P., Chibbar, R.: *Phylogenetic Relationships of Five Morphological Group of Hexaploid wheat Based on RAPD Analysis*. *Genome* 43, 724-727 (2000)
2. Casasoli, M., Mattioni, C., Cherubini, M., Villani, F.: *A Genetic Linkage Map of European Chestnut (Castanea Sativa Mill.) Based on RAPD, ISSR and Isozyme Markers*. *Theoretical Applied Genetics* 102, 1190-1199 (2001)
3. Floreano, D., Mattiussi, C.: *Bio-Inspired Artificial Intelligence. Theories, Methods, and Technologies*. The MIT Press (2008)
4. Groos, C., Gay, G., Perrenant, M., Gervais, L., Bernard, M., Dedryver, F., Charmet, G.: *Study of the Relationships Between Pre-harvest Sprouting and Grain Color by Quantitative Trait Loci Analysis in the White X Red Grain Bread-wheat Cross*. *Theoretical Applied Genetics* 104, 39-47 (2002)
5. Herrera, R., Cares, V., Wilkinson, M., Caligaris, D.: *Characterization of Genetic Variations Between Vitis vinifera Cultivars from Central Chile Using RAPD and Inter Simple Sequence Repeat Markers*. *Euphytica* 124, 139-145 (2002)
6. Muñoz, A.: *Algoritmos Genéticos en la Corrección de Imágenes RAPD*. Informe Técnico, Departamento de Informática, Facultad de Ingeniería, Universidad de Concepción, Chile (2007)
7. Pinninghoff, M. A., Contreras, R., Rueda, L.: *An Evolutionary Approach for Correcting Random Amplified Polymorphism DNA Images*. *Bioinspired Applications in Artificial and Natural Computation. Lecture Notes in Computer Science* 5602, pp. 467-477, Springer (2009)

8. Rueda, L., Uyarte, O., Valenzuela, S. and Rodriguez, J.: Processing Random Amplified Polymorphism DNA Images Using the Radon Transform and Mathematical Morphology. M. Kamel and A. Campilho (Eds.): ICIAR 2007, LNCS 4633, pp. 1071-1081 (2007)
9. Saal, B., Struss, D.: RGA-and RAPD-derived SCAR Markers for a Brassica B-Genome Introgression Conferring Resistance to Blackleg Oil Seed in Oil Seed Rape. *Theoretical Applied Genetics* 111, 281-290 (2005)
10. Sudapak, M., Akkaya, M., Kence, A.: Analysis of Genetic Relationships Among Perennial and Annual Cicer Species Growing in Turkey Using RAPD Markers. *Theoretical Applied Genetics* 105, 1220-1228 (2002)
11. Tripathi, S., Mathish, N., Gurumurthi, K.: Use of Genetic Markers in the Management of Micropropagated Eucalyptus Germplasm. *New Forests* 31, 361-372 (2006)