

Modelo de Saliencia Utilizando el Descriptor de Covarianza

Cristobal Undurruga Rius*, Domingo Mery Quiroz†
Departamento Ciencias de la Computación
Pontificia Universidad Católica de Chile
Av. Vicuña Mackenna 4860, Santiago, Chile
Email: caundurr@uc.cl*, dmery@ing.puc.cl†

Luis Enrique Sucar Succar
Departamento Ciencias Computacionales
Instituto Nacional de Astrofísica, Óptica y Electrónica
Luis Enrique Erro 1, Tonantzintla, Puebla, México
Email: esucar@inaoep.mx

Abstract—Uno de los grandes desafíos de la visión por computador es mejorar los sistemas automáticos para la detección y seguimiento de objetos o regiones en un conjunto de imágenes. Un enfoque que ha cobrado importancia recientemente se basa en la extracción de descriptores de covarianza, ya que logran permanecer invariantes en las regiones de estas imágenes a pesar de los cambios de posición, traslación, rotación y escala. Utilizando el mismo descriptor de covarianza proponemos, en este trabajo, un novedoso sistema de saliencia capaz de encontrar en una imagen aquella(s) región(es) más relevantes que pueden ser utilizadas tanto en la detección como en el seguimiento de objetos. Nuestro método se basa en la cantidad de información de un punto en una imagen y nos permite adaptar las regiones para maximizar la diferencia de información con su entorno. Los resultados muestran que con esta herramienta se puede obtener por una parte un mapa de saliencia, y por otra parte se pueden detectar buenos candidatos a ser seguidos en una secuencia de imágenes.

Index Terms—Saliency, Edge Detector, Tracking, Object Detection, Covariance Descriptor.

I. INTRODUCCIÓN

El ser humano, al ver una imagen o escena, se concentra en ciertas zonas más que otras. La ley general de la percepción dice que las formas simples y sin excesos de información son más fáciles de percibir. Pero en la realidad esto no siempre es cierto, y puede llegar incluso a ser que estas formas sean las más fáciles de confundir. Una característica de la percepción humana es que procesamos las imágenes selectivamente, concentrándonos en zonas que contienen información más relevante.

El ser humano para detectar un objeto trata de utilizar la característica que da mayor información sobre el objeto y en una mayor área. Por ejemplo para detectar a una persona en una multitud, tenderíamos a usar como característica el color de la ropa, ya que es la característica que más área abarca. Pero si ésta no da la información necesaria para diferenciar pasaríamos a la siguiente que mas información diera y así sucesivamente. Para detectar objetos en una imagen existen diferentes enfoques para definir un objeto a través de: puntos de interés (descriptores con información relevante) [1] y [2]; bolsa de palabras (zonas que definen al objeto) [3]; características de una región (variación de las

características de la imagen) [4]; aspecto local (zonas de saliencia)[5]; entre otros. Uno de los métodos que mejores resultados ha obtenido se basa en caracterizar regiones a través del *descriptor de covarianza* propuesto por Porikli et al. [4]. Este descriptor representa una región o ventana por la matriz de covarianza de la matriz formada a partir de las características de la imagen. Nuestro enfoque fue inspirado por los buenos resultados obtenidos en diversas aplicaciones que utilizan este descriptor [6], [7], [8], [9].

Uno de los problemas de los sistemas que utilizan este descriptor u otro es la correcta elección de una ventana que otorgue información significativa para la detección. Dependiendo de la aplicación que se le va a dar al sistema es como se determina esta región. Por ejemplo para hacer seguimiento de una persona se podría elegir un detector de caras como el de Viola-Jones [10] o un detector de personas [11], [12]. De esto surge el problema que dentro de la región seleccionada se encuentra el objeto descrito y también el fondo de la imagen. Esto provoca que mucha información de la región no sea del objeto, lo que provoca un bajo rendimiento del descriptor. Por ende si el objeto tiene poca información significativa y el fondo mucha información, los seguidores se confundirían al considerar que es más importante detectar el fondo.

A través de la cuantificación de la información en la imagen, buscamos resolver este problema. Esto se lograría adaptando la región objetivo para disminuir el ruido provocado por el fondo, maximizando el contraste de información entre la ventana y su entorno. Actualmente para definir si una zona o un punto contiene información relevante se utilizan sistemas de saliencia. La saliencia es un método para predecir donde miraría por primera vez un humano. Por lo tanto es un método para encontrar que es lo que contiene mayor información. Itti et al. [13] nos presenta el modelo posible para el reconocimiento de saliencias a través de un modelo biológico llamado “Teoría de integración de características”. Este método busca saliencias en 3 diferentes capas: una de colores, una de intensidades y una de orientación. Luego combina linealmente las zonas encontradas por las 3 capas para obtener el mapa de saliencia de la imagen. Por otra lado,

Walter & Koch [14] proponen utilizar los llamados proto-objetos, que son elementos volátiles de información visual que pueden ser relacionados a un objeto real. Finalmente Rosin [15] propone efectuar saliencia utilizando la información proporcionada por los bordes.

En nuestro método de saliencia, a diferencia del método de Itti et al. [13], la búsqueda de zonas de saliencia se hace integrando las capas de colores e intensidad, evaluando su covarianza en la vecindad de un punto. Con esta información cuantificamos la cantidad de variación en un punto permitiéndonos formar un sistema de puntos de saliencia. Además introducimos un nuevo sistema de puntos de interés. Esto se hace a través de la saliencia que corresponde a una magnitud de información la cual es local (depende del punto y de su vecindad), a diferencia de otros métodos de saliencia los cuales dan un valor global (dependen de la imagen completa). Si en un punto, su descriptor tiene una gran magnitud entonces es definido como de interés y su vecindad se denomina zona de interés. Podemos crear zonas de diferentes tamaños dependiendo del tamaño de la vecindad del punto.

Por último este trabajo presenta un novedoso sistema para mejorar los sistemas de reconocimiento y de seguimiento a través de la cuantificación de la información en un píxel. Es simple y rápido, utilizando las propiedades del descriptor de covarianza para determinar la varianza de distintas características en un píxel encontramos zonas que contienen mayor información. Este artículo se organiza de la siguiente forma: en la sección 2 se describe el estado del arte actual del problema abordado; en la sección 3 se abordan las bases matemáticas, la hipótesis y la implementación del problema; a continuación, en la sección 4, se presenta la metodología y los resultados; finalmente, en la sección 5 se presentan las conclusiones.

II. MÉTODO PROPUESTO

A. Descriptor de Covarianza

El descriptor de covarianza propuesto por Porikli et al. en [4], se define formalmente como:

$$F(x, y, i) = \phi_i(I, x, y) \quad (1)$$

Donde I es una imagen (la cual puede estar en RGB, tonos de grises, infrarrojo, etc.), F es una matriz de $W \times H \times d$, donde W es el ancho de la imagen, H el alto de la imagen y d es el número de sub-características utilizadas y ϕ_i es la función que relaciona la imagen con la i -ésima característica, es decir la función que obtiene la i -ésima características a partir de la imagen I . Es importante destacar que las características se obtienen a nivel del píxel.

El objetivo es representar el objeto a partir de la matriz de covarianza de la matriz F , construida a partir de estas

características. La covarianza es la medición estadística de la variación o relación entre dos variables aleatorias, esta puede ser negativa, cero o positiva, dependiendo de la relación entre ellas. En nuestro caso las variables aleatorias representarían las sub-características. En la matriz de covarianza las diagonales representan la varianza de cada característica, mientras que el resto representa la correlación entre las características.

La matriz de covarianza tiene las siguientes ventajas como descriptor: 1) unifica información tanto espacial como estadística del objeto; 2) provee una elegante solución para fusionar distintas características y modalidades; 3) tiene una dimensionalidad muy baja; 4) es capaz de comparar regiones, sin estar restringido a un tamaño de ventana constante o fija, ya que no importa el tamaño de la región, el descriptor es la matriz de covarianza, que es de tamaño constante $d \times d$; 5) la matriz de covarianza puede ser fácilmente calculable, para cualquier región o sub-región.

A pesar de todos los beneficios que trae la representación del descriptor a partir de la matriz de covarianza, el cálculo para cualquier sub-ventana o región dado una imagen, utilizando los métodos convencionales, la hace computacionalmente prohibitiva. Porikli et al. en [16] proponen un método computacionalmente superior, para calcular la matriz de covarianza de cualquier sub-ventana o región (rectangular) de una imagen a partir de la formulación de la imagen integral. El concepto de la imagen integral fue inicialmente introducida por Viola and Jones et al. en [17], para el cómputo rápido de características de Haar.

Sea P una matriz de $W \times H \times d$, el tensor de la imagen integral

$$P(x', y', i) = \sum_{x < x', y < y'} F(x, y, i) \quad i = 1 \dots d \quad (2)$$

Sea Q una matriz de $W \times H \times d \times d$, el tensor de segundo orden de la imagen integral

$$Q(x', y', i, j) = \sum_{x < x', y < y'} F(x, y, i) F(x, y, j) \quad (3)$$

$$i, j = 1 \dots d$$

Ahora, sea

$$P_{x,y} = [P(x, y, 1) \quad \dots \quad P(x, y, d)]^T \quad (4)$$

III. EXPERIMENTOS Y RESULTADOS

La finalidad de los experimentos descritos a continuación es mostrar las características de nuestro método y luego mostrar una exitosa aplicación para el mejoramiento de sistemas de seguimiento.

A. Algoritmo de Saliencia

El fin de este algoritmo es determinar la cantidad de información de un punto. Para ello utilizamos el descriptor de covarianza para determinar si un punto es saliente o no. Para cada punto de la imagen, le asignamos como región su vecindad de n vecinos. Con esta región calculamos el descriptor de Covarianza y determinamos su magnitud utilizando el determinante de ésta. Al utilizar solamente el determinante surgen altos valores para ciertas matrices, lo que resolvimos utilizando la función logarítmica. Esto se debe a que ya previamente el logaritmo se ha utilizado como función para determinar la cantidad de información de un elemento. Por lo tanto al punto le asignamos el valor obtenido de (9).

De esto obtenemos un mapa de saliencia utilizando la variación de la características que forman la matriz de covarianza. Este mapa de saliencia en comparación con otros algoritmos es mucho más entendible visualmente, ya que señala la saliencia de cada punto y no la saliencia de una zona (Figura 1).

$$Q_{x,y} = \begin{pmatrix} Q(x,y,1,1) & \dots & Q(x,y,1,d) \\ \vdots & \ddots & \vdots \\ Q(x,y,d,1) & \dots & Q(x,y,d,d) \end{pmatrix} \quad (5)$$

Hay que notar que la matriz $Q_{x,y}$ es simétrica y que para calcular P y Q se necesitan $d + (d^2 + d)/2$ pasos. La complejidad de calcular la imagen integral es de $O(d^2WH)$. Utilizando el método de la imagen integral vemos que la covarianza de cualquier región de la imagen se calcula como:

$$R_Q = Q_{x,y} + Q_{x',y'} - Q_{x',y} - Q_{x,y'} \quad (6)$$

$$R_P = P_{x,y} + P_{x',y'} - P_{x',y} - P_{x,y'} \quad (7)$$

$$C_{R(x,y;x',y')} = \frac{1}{n-1} [R_Q - \frac{1}{n} R_P R_P^T] \quad (8)$$

Donde $n = (x' - x)(y' - y)$. De esta forma, después de construir el tensor de primer orden P y el tensor de segundo orden Q , la covarianza de cualquier región se puede computar en $O(d^2)$.

B. Modelo de Saliencia

Para determinar la cantidad de información que contiene un píxel, primero creamos la matriz F con (1), a continuación obtenemos los tensores de primer y segundo orden a partir de (2) y (3). La idea es obtener la cantidad de información para un píxel, por eso definimos la región del descriptor como la vecindad al punto, el cual se obtiene a partir de (8). Por último, definimos la magnitud de la matriz C_R obtenida como el valor absoluto del determinante de ésta. En teoría de la información es común utilizar el logaritmo para determinar la dispersión de la información. Por lo tanto definimos la cantidad de información I para un punto (x, y) como el logaritmo del determinante de la matriz de covarianza en la vecindad V del punto:

$$S(x, y) = \log(|\det(C_{R(V)})| + c) \quad (9)$$

como deseamos obtener valores positivos definimos la variable c como 1 para obtener solo valores positivos de información.

Algorithm 1 Modelo de Saliencia

- 1: Cálculo de los tensores P y Q de la imagen I
 - 2: Definición de V como el tamaño de la vecindad
 - 3: Definición de $v = \text{floor}(V/2)$
 - 4: $\forall x, y \in I$
 $\text{SaliencyMap}(x,y) = \log(|\det(C_{R(x-v,y-v;x+v,y+v)})| + 1)$
-

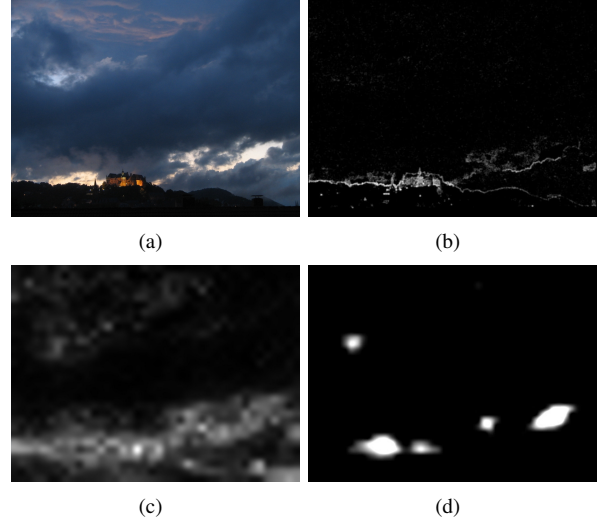


Fig. 1. Ejemplo del mapa de saliency obtenido a través de la matriz de Covarianza: (a) imagen original; (b) imagen del mapa de saliency de la matriz de covarianza; (c) imagen del mapa de saliency para escenas rapidas [13]; (d) imagen del mapa de saliency para proto-objects [14].

Una duda que surgió fue si lo que obteníamos estaba demasiado influenciado por las características de intensidad de la imagen provocando un simple reconocedor de bordes. Por lo tanto hicimos experimentos para determinar la influencia del color en el mapa de saliencia y la influencia de la intensidad (Figura 2). De los resultados determinamos que la intensidad resalta tanto bordes como texturas, mientras que los colores

resalta mayormente los bordes donde existen un cambio de color. Por lo tanto la suma de las características nos daba un detector de bordes resaltando las zonas con cambios de color y opacando las zonas con textura. Por otro lado los detectores de bordes son más puntuales mientras que el descriptor utiliza una región provocando que los zonas de alta variación de colores sean más saliente para nuestro método que para los que utilizan los bordes [15].

B. Algoritmo de Puntos de Interés

Una de las ventajas de nuestro algoritmo frente a otros algoritmos de saliencia es que los puntos contienen la cantidad de información que existe a su alrededor. Si un punto tiene un alto valor, él y su vecindad son de interés para nosotros ya que implica que visualmente es una zona de altos cambios. Por lo tanto los puntos de interés son los puntos con mayor valor. El tamaño de las zonas de interés puede ser regulado a través del tamaño n de la región del descriptor de Covarianza. En realidad un punto es de interés si su zona es de interés, la cual es equivalente a la región del descriptor.

Si utilizamos pequeñas regiones para el descriptor de covarianza, obtenemos zonas de interés pequeñas las cuales asemejan más a puntos de interés que a zonas. En cambio si utilizamos regiones grandes, obtenemos posibles regiones iniciales para algoritmos de seguimiento. Decimos posibles ya que como nos basamos en un sistema de saliencia, es muy posible que lo más saliente sea el fondo de la imagen (ver Figura 3).

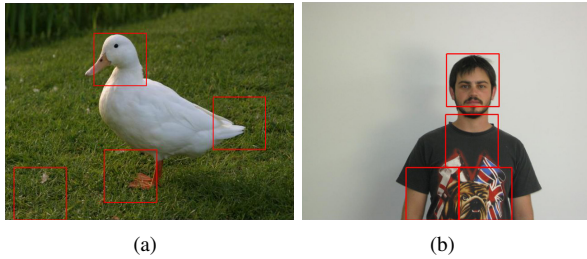


Fig. 3. Ejemplo de las 4 zonas de mayor interés obtenidos con una región cuadrada de tamaño 115 pixel por lado.

C. Detector de Zonas de Saliencia

Ya obtenido el mapa de saliencia, buscamos determinar la ventana donde se concentra la mayor cantidad de información. Para esto creamos un algoritmo que reduce el tamaño de una ventana maximizando la información dentro de la ventana. Para hacerlo de una manera rápida, efectuamos el mismo método que para calcular la matriz de covarianza: primero, creamos la matriz integral del mapa de saliencia I_S ; y segundo, calculamos la información de una región:

$$I_{S(R)} = I_S(x, y) + I_S(x', y') - I_S(x', y) - I_S(x, y') \quad (10)$$

Recordemos que una línea es un rectángulo con un costado de un píxel. Luego establecemos la ventana como toda la imagen y comenzamos a reducirla. Nos ponemos un punto de parada: definimos que porcentaje de la información de la imagen queremos que quede dentro de la ventana. Luego calculamos para cada costado cuanta información proporciona, el que entregue menos información es reducido, y así sucesivamente hasta obtener una región que contenga el porcentaje de información definido de la imagen total (ver Figura 4).

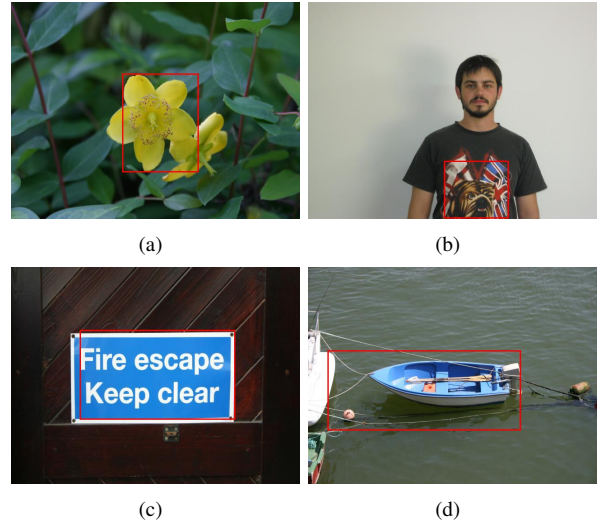


Fig. 4. Ventanas resultantes dejando un 66 por ciento de la información total de la imagen dentro de ellas.

Pero en la realidad los fondos son complejos y tienen mucha información, lo que lleva a que el objeto a rastrear sea el que tiene menos información. Por eso también se puede efectuar el proceso de minimizar la información dentro de la ventana, de esta forma dejamos la mayor parte de la información que corresponde al fondo, la cual produce ruido y errores en los algoritmos de seguimiento, afuera. Eliminando los costados que contienen la mayor cantidad de información, logramos reducir la ventana para que contenga menos información (ver Figura 5).



Fig. 5. Ventanas resultantes dejando un 66 por ciento de la información total dentro de ellas: (a) imagen resultante de un detector de caras; (b) imagen resultante de un detector de personas.

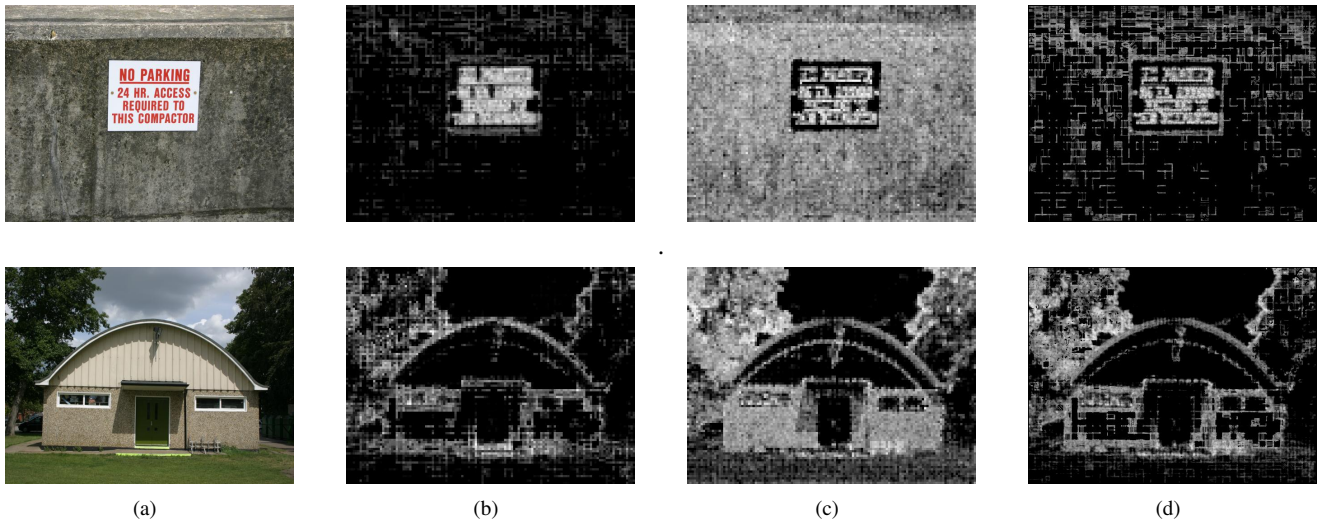


Fig. 2. Ejemplo del mapa de saliencia obtenido a través de la matriz de Covarianza: (a) imagen original; (b) imagen del mapa de saliencia con características de colores; (c) imagen del mapa de saliencia con característica de intensidad; (d) imagen del mapa de saliencia con ambas características.

D. Mejoramiento de Algoritmos de Seguimiento

Una aplicación a la cual sometimos nuestro algoritmo de saliencia fue en algoritmos de seguimiento. Realizamos experimentos con vídeos obtenidos en un supermercado, buscando corroborar si existe una mejoría aplicándolo al inicio de un sistema de seguimiento. El problema surge cuando una mala elección de la ventana de inicio provoca que un algoritmo de reconocimiento o de seguimiento confunda la región con el fondo. Para ello determinamos un método para discernir entre el fondo y el objeto a seguir: determinamos la cantidad de información que contiene un punto en particular utilizando la matriz de covarianza.

Lo que realmente provoca que los sistemas de seguimiento no se confundan es mejorar las ventanas de detección. Utilizando la cantidad de información de cada pixel, queremos maximizar la diferencia de información entre la zona dentro de la ventana y la exterior. Finalmente nuestra hipótesis es que maximizando la diferencia de información mejoraremos el sistema de selección del objetivo para los algoritmos de reconocimiento o de seguimiento, logrando aumentar la eficiencia de estos.

Como algoritmo de seguimiento utilizamos el on-line naive bayes nearest neighbor para el descriptor de covarianza [18]. Testeamos en diferentes ambientes: donde el fondo contiene poca información; donde el fondo contiene mucha información y la región contiene poco fondo; y por último, donde el fondo contiene mucha información y la región contiene mucho fondo. Efectuamos mayoritariamente experimentos donde el fondo contiene mucha información (el fondo es más saliente que el objeto a seguir) ya que en el caso contrario, el sistema de seguimiento no presenta problemas. En el caso en que la región contiene poco fondo no se perciben mayores cambios, no existen mejoras ni empeoramiento. Pero en el último caso,

donde la región contiene harto fondo, se perciben notables mejoras. Por ejemplo logramos evitar que la región se quede pegada en la region inicial ya que no se confunde con el fondo (Ver Figura 6).

En escenarios como el supermercado, donde todo esta hecho para resaltar a la vista, los algoritmos de saliencia seleccionan el fondo como zonas de interés y no las personas, las cuales tienden a tener colores más uniformes. En estos casos usamos el algoritmo de saliencia para restringir la zona inicial de los algoritmos de seguimiento con el fin de maximizar las diferencia entre la region a seguir y el fondo.

IV. CONCLUSIONES Y TRABAJO FUTURO

Para el mejoramiento de los sistemas de seguimiento hemos elaborado un novedoso sistema de detección de saliencia que utiliza el descriptor de covarianza. Este sistema de saliencia, nos permite: determinar si un objeto resultara fácil o difícil de seguir en un vídeo y extraer la información suficiente para poder mejorar los sistemas de detección y seguimiento.

A través de este sistema hemos podido entender porque algunos elementos son más fáciles de seguir que otros. Esto se debe a que los fondos contienen menos información que el objeto o contienen más información que él. También hemos mejorado las regiones iniciales para los detectores y disminuyendo el ruido provocado por los fondos. Aunque no siempre se ven mejoras en el sistema de seguimiento, hemos podido mejorar los casos donde el seguimiento se perdía y manteniendo los casos donde daba buenos resultados.

Por otro lado, nuestro método propuesto entrega un eficiente set de regiones altamente distintivas las cuales pueden ser usadas en algoritmos de detección y/o seguimiento de objetos.



Fig. 6. Resultados obtenido en con una region con mucho fondo en una imagen con mucha información en el fondo. Primera fila corresponde al resultado de referencia obtenido utilizando el método NBNN [18]. La segunda fila es el resultado donde la región inicial fue modificada con nuestro algoritmo.

Nuestras futuras investigaciones se ven enfocadas en hacer mejoramientos a este algoritmo de zonas de interés para ser utilizado como método de reconocimiento de clientes en supermercados. Otra meta para el futuro de esta investigación es evaluar la retro-alimentación del algoritmo de seguimiento para que este se adapte a posibles cambios del objeto o del fondo.

AGRADECIMIENTOS

Agradecemos a Fernando Betteley de Cencosud por facilitar las instalaciones de Supermercados Santa Isabel para la adquisición de vídeos. Agradecemos al Instituto Nacional de Astrofísica, Óptica y Electrónica por recibirme en el marco del Programa de Estancias Cortas de LACCIR (Latin American and Caribbean Collaborative ICT Research). Esta investigación es financiada en parte por LACCIR proyecto S1009LAC006.

REFERENCES

- [1] C. Harris and M. Stephens, "A combined edge and corner detector," *4th Alvey Vision Conference*, 1988.
- [2] D. Lowe, "Object recognition from local scale-invariant features," *Proceedings of the Seventh IEEE International Conference on Computer Vision*, pp. 1150–1157 vol.2, 1999.
- [3] E. Nowak, F. Jurie, and B. Triggs, "Sampling strategies for bag-of-features image classification," *Computer Vision–ECCV 2006*, pp. 490–503, 2006.
- [4] O. Tuzel, F. Porikli, and P. Meer, "Region covariance: A fast descriptor for detection and classification," *Computer Vision–ECCV*, pp. 589–600, 2006.
- [5] D. Jugessur and G. Dudek, "Local appearance for robust object recognition," *cvpr*, 2000.
- [6] J. Batista, "A region covariance embedded in a particle filter for multi-objects tracking," *Update*, 2008.
- [7] X. Li, W. Hu, Z. Zhang, X. Zhang, M. Zhu, and J. Cheng, "Visual tracking via incremental log-euclidean riemannian subspace learning," in *Proc. CVPR*, 2008.
- [8] O. Tuzel, F. Porikli, and P. Meer, "Pedestrian detection via classification on Riemannian manifolds." *IEEE transactions on pattern analysis and machine intelligence*, vol. 30, no. 10, pp. 1713–27, octubre 2008.
- [9] J. Yao, J.-m. Odobez, and C. Parc, "Fast Human Detection from Videos Using Covariance Features," *Learning*, 2008.
- [10] P. Viola and M. Jones, "Robust real-time object detection," *International Journal of Computer Vision*, 2002.
- [11] N. Dalal, B. Triggs, and C. Schmid, "Human detection using oriented histograms of flow and appearance," *Computer Vision–ECCV*, p. 428â441, 2006.
- [12] M. Andriluka, S. Roth, and B. Schiele, "Pictorial structures revisited: People detection and articulated pose estimation," *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1014–1021, junio 2009.
- [13] L. Itti, C. Koch, and E. Niebur, "A Model of Saliency-Based Visual Attention for Rapid Scene Analysis," *Colorectal disease : the official journal of the Association of Coloproctology of Great Britain and Ireland*, vol. 4, no. 2, pp. 147–149, marzo 2002.
- [14] D. Walther and C. Koch, "Modeling attention to salient proto-objects." *Neural networks : the official journal of the International Neural Network Society*, vol. 19, no. 9, pp. 1395–407, noviembre 2006.
- [15] P. L. Rosin, "A simple method for detecting salient regions," *Pattern Recognition*, vol. 42, no. 11, pp. 2363–2371, noviembre 2009.
- [16] F. Porikli and O. Tuzel, "Fast Construction of Covariance Matrices for Arbitrary Size Image Windows," *2006 International Conference on Image Processing*, pp. 1581–1584, octubre 2006.
- [17] P. Viola and M. Jones, "Rapid Object Detection Using a Boosted Cascade of Simple Features," *Proc. IEEE Conf. Computer Vision and Pattern Recognition*, vol. 01, pp. vol1pp511–518, 2001.
- [18] P. Cortez, D. Mery, and L. E. Sucar, "Object Tracking based on Covariance Descriptors and On-Line Naive Bayes Nearest Neighbor Classifier," *Submitted to the Pacific-Rim Symposium on Image and Video Technology*, 2010.