

IMAGENES SIMILARES DE FACEBOOK

Carolina Guevara¹, Jos Fajardo², Hemerson Morn³, Felix Drouet⁴

¹ Universidad de Guayaquil

rosa.guevarap, jose.fajardod, hemerson.moranp, felix.drouetv @ug.edu.ec

Abstract—El uso de la coincidencia gráfica aproximada para la minería de subgrafos frecuentes ha sido identificado en diferentes aplicaciones como una necesidad. Para satisfacer esta necesidad, varios algoritmos de minería se han desarrollado y se han utilizado en diferentes dominios de la ciencia. En este trabajo, se presenta un nuevo marco para la clasificación de imágenes basado en el gráfico, donde se utiliza un algoritmo para la minería de subgrafos frecuentes. En este momento, la correspondencia aproximada gráfica se aproxima nos permiten desarrollar clasificadores robustos en presencia de distorsiones en este tipo de datos. La experimentación se realizó sobre una colección de imágenes reales, donde cada imagen se representa como un gráfico. Los resultados de esta experimentación muestran la utilidad de los algoritmos aproximados y el marco propuesto.

Palabras clave: minería de grafos aproximados; cotejo aproximado de grafos; representación de imágenes; clasificación de imágenes; selección de características.

I. INTRODUCCION

En los últimos años se ha incrementado la necesidad de convertir grandes volúmenes de datos en información útil. Los objetos en muchas de estas bases de datos están o pudieran estar representados como grafos. Como respuesta a esta necesidad, varios autores han desarrollado técnicas y métodos para procesar estas bases de datos (JIMNEZ et al., 2010). Un ejemplo de estas técnicas es el descubrimiento de patrones frecuentes (YUN y RYU, 2011).

La minería de patrones frecuentes, especialmente la detección de subgrafos frecuentes en colecciones de grafos es un problema en tareas de minería de grafos (GAGO-ALONSO et al., 2009; NIJSSEN and KOK, 2004; YAN and HUAN, 2002). En minería de subgrafos frecuentes, existen dos enfoques para evaluar la similitud de grafos, conocido como cotejo exacto de grafos y cotejo aproximado de grafos.

El cotejo exacto de grafos ha sido satisfactoriamente utilizado en diversas aplicaciones (EICHINGER and BHM, 2010; GAGO-ALONSO et al., 2010; JIANG et al., 2010a); sin embargo, existen problemas concretos donde un cotejo exacto no es aplicable con resultados satisfactorios (HOLDER et al., 1992). En ocasiones, los subgrafos muestran ligeras diferencias en los datos. Un ejemplo de estas diferencias se puede ver en el procesamiento de imágenes, donde estas diferencias pueden estar debido a ruido y distorsiones, o simplemente pueden presentar ligeras diferencias espaciales entre instancias de objetos iguales.

Esto significa que se debe tolerar cierto nivel de distorsiones geométricas, variaciones semánticas o desajustes entre vértices o aristas, mientras se realiza la búsqueda de los subgrafos frecuentes. Por este motivo, varios autores han expresado la necesidad del uso del cotejo aproximado entre grafos para la minería de subgrafos frecuentes en colecciones de grafos (BORGELT and BERTHOLD, 2002; HOSSAIN and ANGRYK, 2007; KETKAR et al., 2006; KOYUTRK et al., 2004). Estos autores defienden la idea de que se pudieran detectar subgrafos frecuentes con mayor interés para aplicaciones y usuarios.

En respuesta a esta necesidad, varios algoritmos han sido desarrollados para la minería de subgrafos frecuentes utilizando cotejo aproximado entre grafos para diferentes dominios de la ciencia como: análisis de estructuras bioquímicas (CHEN et al., 2007; JIA et al., 2011; XIAO et al., 2008; ZHANG and YANG, 2008; ZOU et al., 2009), redes genéticas regulatorias (SONG and CHEN, 2006); análisis de circuitos, redes sociales, y análisis de vínculos (HOLDER et al., 1992).

Big data y la recuperación de imágenes web por su contenido. Desde una perspectiva big data, en el ámbito web muchas veces es necesario recuperar una imagen de una biblioteca distribuida. En el análisis de una determinada página web son muchos los escenarios que se pueden encontrar: es posible requerir la obtención de la fuente de las imágenes que aparecen en dicha página, también es posible analizar si sus imágenes han sido publicadas en terceras páginas, o reconocer a una persona en las fotografías de dicha web. En todos estos casos el problema central es analizar el contenido de las imágenes de interés localizadas en una web para buscar y encontrar una imagen en una biblioteca de imágenes determinadas. Como se puede intuir, cuando se está hablando de grandes cantidades de datos la búsqueda en miles o millones de imágenes será una búsqueda computacionalmente cara y difícil de ejecutar en equipos informáticos básicos.

La pregunta clave en este contexto es si la tecnología evoluciona y crece a la misma velocidad que lo están haciendo estos datos. Para el campo en el que se centra este artículo, el análisis de imágenes, la respuesta es no. Trabajos como los de Guo y Dyer (2005) o White et al. (2010), explican cómo la infraestructura (recursos de almacenamiento y computación necesarios) es difícilmente adquirible para llevar a cabo aplicaciones de análisis de imagen a gran escala, por lo tanto existen pocos investigadores que se aventuren en esta área. Esto a su vez provoca que el número de trabajos dedicados

al procesamiento de grandes volúmenes de imágenes o videos sea relativamente escaso y monopolizado por grandes grupos de investigación. El tratamiento de la imagen, por su relevancia y complejidad, debe tener un espacio propio en el mundo big data. Si prosiguen las tendencias actuales y la sentencia de que los datos son el nuevo petróleo del siglo XXI acuada por

Por otro lado, la minería de subgrafos frecuentes se ha utilizado satisfactoriamente en clasificación de imágenes (BAHADIR and SELIM, 2010; JIANG and COENEN, 2008; JIANG et al., 2010b). No obstante, casi todos los enfoques han estado basados en el uso del cotejo exacto. En este trabajo, se propone un esquema basado en grafos para la clasificación de imágenes. Este esquema utiliza los subgrafos frecuentes como características obtenidas mediante algoritmos de minería de subgrafos frecuentes en una colección de imágenes reales. Utilizando este esquema se evalúan los algoritmos mediante la clasificación de imágenes.

II. TRABAJOS RELACIONADOS

1 Clasificación de imágenes utilizando minería de subgrafos frecuentes aproximados Tipo de artículo: Artículo original Temática: Inteligencia artificial, Procesamiento de imágenes, Reconocimiento de patrones

Resumen: Se aborda el análisis web desde el punto de vista de las imágenes, empleando tecnologías big data. Las imágenes cada vez tienen más peso en la web por lo que cualquier análisis que se realice debe considerar este tipo de información. Los grandes volúmenes de imágenes existentes hacen necesaria la utilización de grandes infraestructuras de computación para realizar este tipo de trabajos, así como tecnologías de visión artificial específicas. Se muestran tecnologías big data que pueden ser utilizadas dentro del campo del análisis de imágenes a gran escala. Además, se propone una arquitectura que permite recuperar imágenes de una biblioteca de imágenes de forma eficiente y con un bajo coste computacional. Esta arquitectura puede servir como base para los análisis web e investigaciones que requieran un estudio detallado de las imágenes similares, sin la necesidad de disponer de hardware específico para ello.

2 Tecnologías big data para análisis y recuperación de imágenes web Rodríguez-Vaamonde, Sergio; Torre-Bastida, Ana-Isabel; Garrote, Estibaliz (2014). Tecnologías big data para análisis y recuperación de imágenes web. El profesional de la información, v. 23, n. 6, noviembre-diciembre, pp. 567-574.

Resumen: El uso del cotejo aproximado de grafos para la minería de subgrafos frecuentes se ha identificado en diferentes aplicaciones como una necesidad. Con el fin de enfrentar este reto, varios algoritmos de minería han sido desarrollados y han sido utilizados en varios dominios de la ciencia. En este trabajo, se presenta un nuevo esquema para la clasificación de imágenes basado en grafos donde se utiliza un algoritmo para la minería de subgrafos frecuentes. Esta vez el enfoque del cotejo aproximado de grafos nos permite desarrollar clasificadores robustos ante distorsiones presentes

en este tipo de datos. La experimentación se realiza sobre una colección real de imágenes representadas en forma de grafos. Los resultados de esta experimentación muestran la utilidad del uso de los algoritmos aproximados y del esquema propuesto.

A. Big data y la recuperación de imágenes web por su contenido

Desde una perspectiva big data, en el ámbito web muchas veces es necesario recuperar una imagen de una biblioteca distribuida. En el análisis de una determinada página web son muchos los escenarios que se pueden encontrar: es posible requerir la obtención de la fuente de las imágenes que aparecen en dicha página, también es posible analizar si sus imágenes han sido publicadas en terceras páginas, o reconocer a una persona en las fotografías de dicha web.

En todos estos casos el problema central es analizar el contenido de las imágenes de interés localizadas en una web para buscar y encontrar una imagen en una biblioteca de imágenes determinadas. Como se puede intuir, cuando se está hablando de grandes cantidades de datos la búsqueda en miles o millones de imágenes ser una búsqueda computacionalmente cara y difícil de ejecutar en equipos informáticos básicos.

La pregunta clave en este contexto es si la tecnología evoluciona y crece a la misma velocidad que lo están haciendo estos datos. Para el campo en el que se centra este artículo, el análisis de imágenes, la respuesta es no. Trabajos como los de Guo y Dyer (2005) o White et al. (2010), explican como la infraestructura (recursos de almacenamiento y computación necesarios) es difícilmente adquirible para llevar a cabo aplicaciones de análisis de imagen a gran escala, por lo tanto existen pocos investigadores que se aventuren en esta área.

Esto a su vez provoca que el número de trabajos dedicados al procesamiento de grandes volúmenes de imágenes o videos sea relativamente escaso y monopolizado por grandes grupos de investigación. El tratamiento de la imagen, por su relevancia y complejidad, debe tener un espacio propio en el mundo big data. Si prosiguen las tendencias actuales y la sentencia de que los datos son el nuevo petróleo del siglo XXI acuada por Andreas Weigend es cierta, el potencial que se pueda extraer de estos datos depender irremediablemente de las tecnologías y algoritmos desarrollados para ello, y en el caso de las imágenes está claro que se necesita potenciar ambos.

Para ello presentamos en primer lugar un resumen de las principales tecnologías big data existentes y su aplicabilidad a los datos en formato imagen. Abarcar el cien por ciento de las tecnologías queda fuera del alcance de este estudio, por lo que nos centraremos en el caso particular de la recuperación rápida de imágenes en base a su contenido y presentaremos una arquitectura que contiene los ingredientes necesarios para hacer búsquedas rápidas de imágenes que permitan un amplio abanico de análisis web sin necesidad de una gran inversión en infraestructura.

B. Tecnologías big data aplicables a imágenes

Las necesidades que deben cumplir las tecnologías del big data se basan en el procesamiento eficiente de grandes

cantidades de datos con un tiempo reducido o tolerable.

La complejidad que aade el hecho de que los datos se encuentren en formato imagen es otra variable a considerar. En general las tecnologas de mayor relevancia para datos no estructurados como imgenes son las bases de datos NoSQL (Leavitt, 2010) y los modelos de programacin Map-Reduce (Bajcsy et al., 2013), ambas relacionadas con el procesamiento de datos en lotes. Por otro lado estn los CEP (complex event processing), los IMDG (in-memory data grids), o los sistemas de computacin distribuida, para el procesamiento de datos en tiempo real. En la tabla siguiente se muestra una taxonoma de estas tecnologas en forma de cuadrante ordenado por volumen de datos y tiempos de rendimiento de las tecnologas big data. En los siguientes apartados se detallan las tecnologas ms relevantes en los mbitos anteriores, centrndose en aquellas especialmente tiles para el procesamiento de imgenes a gran escala.

TABLE I
TAXONOMIA DE TECNOLOGIAS DE BIG DATA APLICABLES AL PROCESAMIENTO DE IMAGENES

VOLUMEN DE DATOS	POCOS	% MUCHOS
Tiempo real	Analitica stream eventos complejos BD en memoria	% Analitica tiempo real % Grid datos en memoria % Plataforma especializada
Lotes	Analitica operacin OLPT / OLAP BD relacional	% Analitica lotes %Map-Reduce % BD NoSQL

III. DATOS

TABLE II
DATOS ESTADISTICOS DEL PROYECTO.

MEDIA	MEDIANA	VARIANZA	% DESV.ESTANDAR
25.5	25.5	222.83	14.93

A. Arquitectura software de recuperacin rpida de imgenes similares

Para poder hacer este tipo de anlisis y recuperacin, se propone la utilizacin de tecnologas big data que permitan un acceso eficiente a toda la informacin disponible en imgenes. La propuesta se puede resumir en la figura siguiente. El primer problema es el almacenamiento eficiente de las imgenes de entrada. Estas imgenes pueden entrar al sistema de diferentes formas, en funcin del anlisis al que se est dedicando esta arquitectura. En el caso de anlisis de imgenes web, la entrada sera por un robot automtico que obtuviese las imgenes de las paginas web.

IV. METODOLOGIA

Dado un conjunto de imgenes representados en forma de grafos pre-etiquetados propuesto por Riesen y Bunke (RIESEN and BUNKE, 2008) se utilizan los algoritmos para la MSFA con el objetivo de obtener todos los subgrafos frecuentes de esta coleccin. Luego, estos subgrafos son utilizados para construir los vectores de caractersticas de las imgenes

Distancia	Frmula
Euclidea	$d(x, y) = \sqrt{\sum_{i=1}^n (x_i - y_i)^2}$
Coseno	$d(x, y) = 1 - \cos(\theta) = 1 - \frac{x \cdot y}{\sqrt{(x \cdot x)(y \cdot y)}}$
Chi Cuadrado	$d(x, y) = \frac{1}{2} \sum_{i=1}^n \frac{ x_i - y_i ^2}{x_i + y_i}$

originales. Finalmente, se emplea un clasificador usando estos vectores como datos para realizar la clasificacin de dichas imgenes.

A. Clasificacin de imgenes

Teniendo en cuenta los subgrafos frecuentes obtenidos, se construyen los vectores de caractersticas, entonces una imagen es representada como un vector de caractersticas, donde es el nmero total de subgrafos identificados. Por tanto, se construye una matriz donde el nmero de las filas corresponde al nmero de grafos (imagenes) en la coleccin y el nmero de las columnas corresponde a la cantidad de subgrafos frecuentes (caracteristicas). Cada valor de caracterstica puede ser asignado utilizando una configuracin binaria o una configuracin de similitud. En la configuracin binaria, una celda de la matriz es si la caracterstica ocurre en la imagen de la coleccin y en otro caso . Una celda de la matriz en la configuracin de similitud es el mayor valor de similitud de una ocurrencia de la caracterstica en la imagen. El valor de la similitud de cada caracterstica se obtiene mediante utilizando la funcin de similitud que le corresponda a cada algoritmo en particular.

Para la clasificacin mediante el clasificador SVM (de sus siglas en ingles, Support Vector Machine) se utiliz el paquete libSVM1. En el caso de APGM, se usa la matriz indizada por las etiquetas de los vrices, las cuales representan el tipo de vertice.

B. Clasificacin de imgenes

Teniendo en cuenta los subgrafos frecuentes obtenidos, se construyen los vectores de caractersticas, entonces una imagen es representada como un vector de caractersticas, donde es el numero total de subgrafos identificados. Por tanto, se construye una matriz donde el nmero de las filas corresponde al nmero de grafos (imagenes) en la coleccin y el nmero de las columnas corresponde a la cantidad de subgrafos frecuentes (caracteristicas).

Cada valor de caracterstica puede ser asignado utilizando una configuracin binaria o una configuracin de similitud. En la configuracin binaria, una celda de la matriz es si la caracterstica ocurre en la imagen de la coleccin y en otro caso . Una celda de la matriz en la configuracin de similitud es el mayor valor de similitud de una ocurrencia de la caracterstica en la imagen. El valor de la similitud de cada caracterstica se obtiene mediante utilizando la funcin de similitud que le corresponda a cada algoritmo en particular.

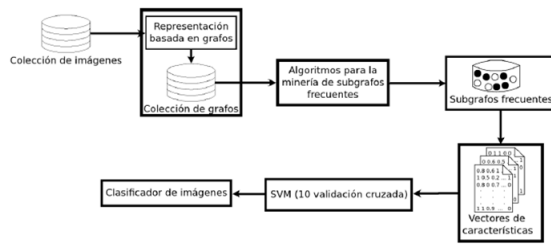


Figura 1. Esquema de clasificación basada en grafos.

C. Esquema para la clasificación de imágenes

En esta sección se presenta el esquema propuesto para mostrar la utilidad de los métodos aproximados en tareas de clasificación de imágenes. Los detalles de la clasificación de imágenes son introducidos con el esquema propuesto.

de búsqueda de una imagen concreta, se dispondrá de una base de datos de imágenes y habrá que introducir manualmente la nueva imagen de consulta. En cualquier caso el problema principal consiste en almacenar las imágenes y, además de la propia imagen, es crucial almacenar sus descriptores visuales. Ya que esta no es una información estática y es dependiente del análisis o búsqueda a realizar, es necesario disponer de un almacén de datos lo suficientemente flexible. Además se debe poder guardar una gran colección de datos, como son el origen de las imágenes (web, fecha de acceso, etc.), anotaciones manuales o comentarios. Para ello, cualquier base de datos NoSQL de las presentadas en la tabla 1 es una opción válida ya que en general permiten disponer de un esquema flexible y son capaces de tener réplicas o nodos distribuidos. Esta última característica la hace ideal para el análisis de imágenes web en cualquier punto geográfico, ya que permite replicas distribuidas en todo el mundo. Sobre este almacén de datos NoSQL es necesario construir el sistema de recuperación de figuras. No es factible computar para cada imagen de la colección la distancia a la imagen de consulta. Por ello es fundamental utilizar algún sistema de búsqueda aproximada de las más cercanas. El algoritmo más utilizado en el campo del análisis de imágenes (Kulis; Grauman, 2009) es locality-sensitive hashing (LSH) (Slaney; Casey, 2008), por lo que será el algoritmo base para la recuperación de figuras. Este algoritmo permite generar una firma numérica (o hash) para cada descriptor o conjunto de descriptores de imagen, de tal forma que aquellos vectores que tengan una distancia euclídea muy baja, y por tanto sean vectores muy similares, posean la misma firma numérica. Este tipo de algoritmos es muy útil para encontrar entradas similares dentro de grandes colecciones de datos, por ejemplo buscando páginas web similares (Slaney; Casey, 2008), por lo que es lógica su aplicación a los descriptores visuales de imágenes. Una vez se tienen las firmas para todas las imágenes de la colección, la búsqueda de las similares es sencilla: dada una imagen de consulta, se genera su hash. Con cada firma se busca en toda

la base de datos las que posean la misma firma y todas ellas serán las imágenes más similares. Para hacer esta comparación, se puede pensar en que existe el mismo problema de búsqueda que antes, pero nada más lejos de la realidad. Las bases de datos NoSQL actuales para big data, disponen de técnicas de indexación y búsqueda

rápida de un número concreto, como puede ser el algoritmo de búsqueda en árbol binario de la base de datos NoSQL MongoDB o el uso de cualquier tecnología IMDG de las propuestas. Por ello la búsqueda ya no se circunscribe a calcular una distancia entre vectores sino a usar una arquitectura de índices para encontrar un número concreto. Tras este paso ya se dispone de un conjunto de imágenes de la biblioteca similares a la de la entrada. En función de la analítica web que se está ejecutando, quizá conocer este número es suficiente. En muchos casos la recuperación de figuras tiene como objetivo aquella que más se parece a la de entrada. En este caso, es obligatorio calcular la distancia concreta, pero ya que se ha obtenido un conjunto de imágenes parecidas, se puede calcular la distancia sobre ese conjunto de unas decenas de imágenes similares en unos pocos segundos, en vez de sobre el total de la biblioteca. Para este último paso, también se va a aprovechar el almacén de datos NoSQL distribuido propuesto en el inicio. Ya que las figuras pueden estar almacenadas en localizaciones diferentes y que cada cálculo de la distancia entre la imagen de consulta y la similar es independiente, es posible usar el paradigma Map-Reduce expuesto con anterioridad. Este modelo permitirá en la función Map el cálculo de la distancia entre cada imagen similar y la de consulta, ejecutándose en cada nodo de la red distribuida de almacenamiento. Por otro lado, el método Reduce se encargará de ordenar todas las distancias y podrá generar el ranking final de imágenes similares para la analítica.

V. RESULTADOS

Mediante los resultados experimentales, determinaremos la utilidad del uso de la minería de subgrafos frecuentes aproximados en tareas de clasificación de imágenes. Esperamos realizar un código útil y eficiente, es decir que trabaje de la manera especificada, reflejando sobre todo los conocimientos adquiridos durante el semestre.

VI. BIBLIOGRAFIA

- Ajovalasit, A., Petrucci, G. and Scafidi, M., RGB photoelasticity applied to the analysis of membrane residual stress in glass. *Measurement science and technology*, vol 23, pp. 1-4, 2012
- Azzam, R.M.A. The intertwined history of polarimetry and ellipsometry. *Thin Solid Films*, Volume 519, pp. 2584- 2588, 2011.
- Kasimayan T. and Ramesh, K., Digital reflection photoelasticity using conventional reflection polariscope. *Optics and Lasers in Engineering*, vol 34, pp. 45-51, 2010
- Kasimayan, T., Ramesh, K., Digital reflection photoelasticity using conventional reflection polariscope. *ScienceDirect*, vol 34, pp. 45-51, 2010.
- Wang, Z., Bovik, A. C. H. R. and Sheikh, E., Simoncelli, P., Image quality assessment: from error visibility to structural

similarity. IEEE Transactions on Image Processing, vol 13, pp. 600–612, 2004.

R. Dosselmann, X. and Yang, A., comprehensive assessment of the structural similarity index. Signal, Image and Video Processing, vol 5, pp. 8191, 2011.

Maldonado, M., Sanchez, G. and Branch, J., Registration of range images using a histogram based metric. Dyna, vol 79, (176), pp. 27-34, 2012.