

# Reconocimiento de objetos con información de color y profundidad

Curso: Introducción al Reconocimiento de Patrones  
Guillermo Carbajal – Alvaro Gómez  
Diciembre 2011

## Resumen

El reconocimiento de objetos mediante imágenes de color se viene investigando desde hace varios años con diversos enfoques. Uno de los enfoques más utilizados es el reconocimiento a partir de características locales con ciertas propiedades de invarianza frente a cambios de iluminación y punto de vista. En dicho enfoque, se aprende un cierto vocabulario de características y cada imagen se caracteriza por la frecuencia y ubicación en que aparecen las palabras de ese vocabulario. El reconocimiento de nuevas imágenes se puede realizar mediante la selección del vecino más cercano en el conjunto de entrenamiento o mediante un clasificador entrenado con el conjunto de entrenamiento.

En los últimos dos años se ha popularizado el uso de sensores que permiten obtener mapas de profundidad además de imágenes RGB (sensores RGB-D). A partir de los mapas de profundidad se pueden extraer características que aportan valiosa información sobre la forma de los objetos. La utilización en conjunto de características de color y de forma pueden determinar un mejor desempeño en el reconocimiento de instancias de objetos y especialmente en el reconocimiento de categorías de objetos.

En este trabajo se analiza el reconocimiento de categorías de objetos combinando información de color y de profundidad.

## I. Introducción

El reconocimiento de objetos se puede separar en dos grupos:

- reconocimiento de instancias  
implica reconocer un objeto previamente conocido en una toma distinta que puede tener un nuevo punto de vista, puede incluir otros objetos extraños y/o puede presentar oclusiones.
- reconocimiento de categorías  
consiste en reconocer que un cierto objeto no necesariamente visto antes es de una cierta categoría (por ejemplo “es un auto”, “es una persona”)

El reconocimiento de objetos mediante imágenes color se viene investigando desde hace varios años. Algunos de los enfoques que se han desarrollado a lo largo de los años son:

- Detección de líneas, contornos y/o superficies para luego matchear contra modelos 2D o 3D
- Adquisición de imágenes desde diversas posiciones y orientaciones para representarlas en un espacio vectorial y realizar una descomposición en una base con los valores propios más importantes (ej. Eigenfaces [2])
- Extracción de un conjunto de características locales que tengan propiedades de invarianza frente a cambios de iluminación y punto de vista. Se reconoce una imagen comparando estas características locales con las características de las imágenes de la base conocida (debe haber una suficiente cantidad de correspondencias y esas correspondencias deben ser coherentes con una transformación que alinee las imágenes). (Ej. Matching con SIFT [3])

Cuando la cantidad de imágenes consideradas crece, no es posible la comparación 1 a N. Las características locales son mapeadas a un conjunto de “palabras visuales”. Estas “palabras visuales” se aprenden por ejemplo mediante k-means sobre un conjunto de entrenamiento. El reconocimiento de una nueva imagen se realiza comparando contra la base la frecuencia de aparición de “palabras visuales”. Esto da un ranking de candidatos que se puede refinar teniendo en cuenta la coherencia espacial de correspondencias [11,12].

En los últimos dos años se ha popularizado el uso de sensores que permiten obtener mapas de profundidad además de imágenes RGB (sensores RGB-D) [19]. A partir de los mapas de profundidad se pueden extraer características que aportan valiosa información sobre la forma de los objetos. La utilización en conjunto de características de color y de forma pueden determinar un mejor desempeño en el reconocimiento de instancias de objetos y especialmente en el reconocimiento de categorías de objetos.

En este trabajo se estudia el reconocimiento de categorías y se realizan pruebas sobre una extensa base de imágenes RGB-D de objetos domésticos. Se comparan los resultados de clasificación utilizando características RGB, características sobre la profundidad y la combinación de las mismas.

El resto del artículo se organiza de la siguiente manera: La sección II describe la base de imágenes utilizada. La sección III presenta el esquema de reconocimiento utilizado. La sección IV presenta las herramientas de software utilizadas para las corridas sobre la base y la pequeña aplicación desarrollada. La sección V describe los experimentos realizados, los resultados obtenidos y su comparación con algunos resultados de la bibliografía. Finalmente, en la sección VI se presentan las conclusiones del trabajo.

## II Base de imágenes

Se utiliza la base de datos disponible en <http://www.cs.washington.edu/rgbd-dataset/> [5].

La figura 1 muestra algunos de los objetos de la base

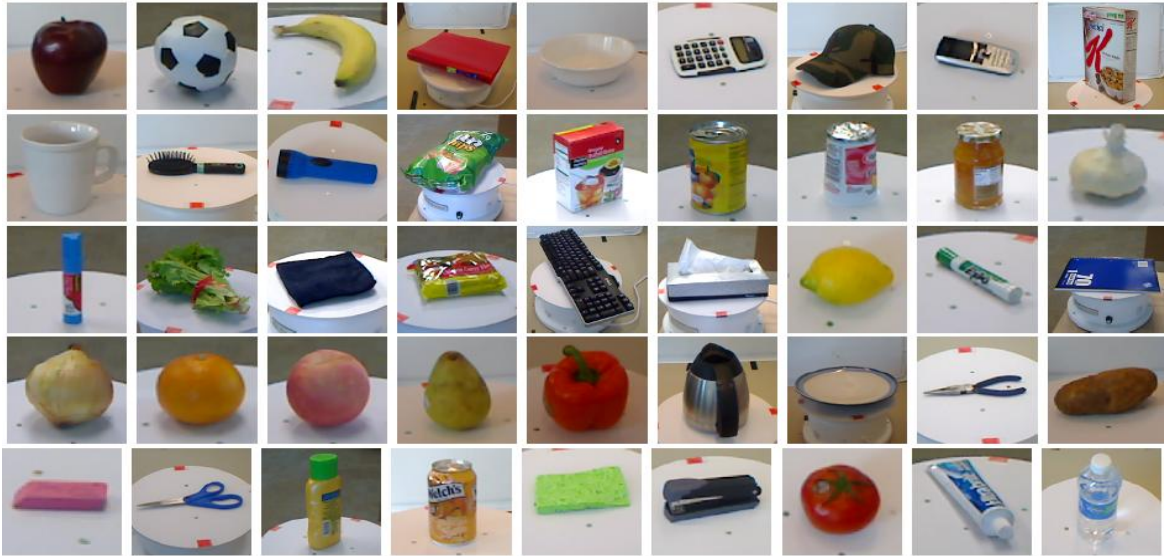


Fig. 1 Algunos de los objetos de la base. Imagen tomada de [5]

Esta es una base de 300 objetos domésticos con múltiples imágenes RGB-D desde diversas posiciones tomadas con un sensor tipo Kinect. Las imágenes han sido tomadas colocando los objetos sobre una mesa giratoria, cubriendo los 360° de azimut en torno al objeto y desde 3 diferentes elevaciones (aproximadamente 30°, 45° y 60°). En la figura 2 se muestra un ejemplo.



Fig. 2 Varias vistas de una taza extraídas del video tomado desde una elevación de 30°. Se cuentan con aproximadamente 600 vistas por instancia.

Para cada elevación se cuenta con un promedio de 200 imágenes por instancia.

Los objetos están organizados en 51 categorías y cada categoría cuenta con un conjunto de diferentes instancias. El número de instancias por categoría varía en un rango de 3 a 12. En la figura 3 se muestran algunas de las categorías.



Fig. 3 Cada fila de la tabla muestran instancias de una de las categorías de la base. La base cuenta con un total de 51 categorías y el número de instancias por categoría varía en un rango de 3 a 12.

- De izquierda a derecha se muestran las instancias 1, 2, 3, 4 y 5 de algunas categorías elegidas para ejemplificar

Los autores de la base han realizado pruebas de reconocimiento de instancias y categorías de objetos sobre la base utilizando la información RGB, de profundidad y la combinación de ambas [5]. Para la información de color extraen un muestreo denso de patches SIFT a dos escalas diferentes, sobre toda la imagen y también sobre una división 2x2 de la imagen. Para la información de profundidad utilizan Spin Images [17] sobre los puntos de la nube. En los dos casos la información se integra mediante la técnica de Efficient Matching Kernels [18] y se reduce de dimensión mediante PCA. Con este esquema, los autores de la base realizan pruebas de reconocimiento de instancias de objetos y de categorías y reportan los resultados de la figura 4.

Classifier	Shape	Vision	All
	Category		
LinSVM	53.1 ± 1.7	74.3 ± 3.3	81.9 ± 2.8
kSVM	64.7 ± 2.2	74.5 ± 3.1	83.8 ± 3.5
RF	66.8 ± 2.5	74.7 ± 3.6	79.6 ± 4.0
	Instance (Alternating contiguous frames)		
LinSVM	32.4 ± 0.5	90.9 ± 0.5	90.2 ± 0.6
kSVM	51.2 ± 0.8	91.0 ± 0.5	90.6 ± 0.6
RF	52.7 ± 1.0	90.1 ± 0.8	90.5 ± 0.4
	Instance (Leave-sequence-out)		
LinSVM	32.3	59.3	73.9
kSVM	46.2	60.7	74.8
RF	45.5	59.9	73.1

Fig. 8. Category and instance recognition performance of various classifiers on the RGB-D Object Dataset using shape features, visual features, and with all features. LinSVM is linear SVM, kSVM is gaussian kernel SVM, RF is random forest.

Fig. 4 Resultados de reconocimiento de instancias y categorías reportado en [5]

Para el reconocimiento de instancias, Lai et al [5]. realizan dos pruebas diferentes variando los conjuntos de entrenamiento y test. En el caso que denominan "Leave-sequence-out", se entrena con las imágenes de los video a 30° y 60° de elevación y se testea con las imágenes del video a 45° de elevación. En el caso que denominan "Alternating-contiguous-frames" los conjuntos de entrenamiento y test tienen imágenes intercaladas de los tres videos.

Para el reconocimiento de categorías, se deja una instancia al azar para test de cada categoría y se entrena con las imágenes del resto de las categorías.

En la figura 4 puede verse que la combinación de información de color y de forma se vuelve importante en el reconocimiento de categoría donde las diferentes instancias de una categoría suelen tener variaciones importantes en su aspecto visual pero comparten formas similares. La información de forma también aporta en el reconocimiento de instancia "Leave-sequence-out" donde el aspecto visual de los conjuntos de entrenamiento y test tiene diferencias por provenir de videos con puntos de vista diferentes.

### **III Esquema de reconocimiento**

Para el reconocimiento de categorías se utilizó el método conocido como “Bag of Features” o “Bag of Words” [11,12,20] (sección III.a). Las imágenes de la base se describieron utilizando información de intensidad (SIFT), color (opponent-SIFT) y profundidad (FPFH) . Los dos primeros se describen en la sección III.b y el tercero en la sección III.c.

La clasificación se realizó utilizando Support Vector Machines (SVM) lineal y con kernel chi cuadrado. Se realizaron pruebas de clasificación utilizando los descriptores “puros” y combinando la información de profundidad con la rgb. La combinación se realizó siguiendo las siguientes dos estrategias diferentes: concatenando descriptores y mezclando clasificadores expertos. (sección III.d)

#### **III.a Reconocimiento mediante “Bag of Features” o “Bag of Words”**

Si se quiere saber de que trata un determinado texto en general basta con conocer algunas de las palabras del mismo, a partir de ciertas palabras claves suele inferirse si se trata de una noticia, una publicidad o una ley, por ejemplo. Uno de los enfoques más simples y utilizados en el reconocimiento de objetos consiste en suponer que es posible describir un objeto y por lo tanto una escena utilizando un conjunto de palabras. De manera análoga a cómo sucede en los textos, la frecuencia con que aparecen las palabras del vocabulario determinarán la categoría del objeto.

Las imágenes rgb-d describen los distintos objetos que aparecen en ellas utilizando información de color y profundidad pero no utilizan letras. Menos aún palabras. Para emplear este enfoque es necesario

1. Elegir una manera adecuada de describir las imágenes.

Se debe traducir la información de color y profundidad en descriptores que permitan formar un vocabulario. Algunas propiedades deseables de los descriptores son las siguientes:

- robustos a cambios en las condiciones de adquisición: iluminación, perspectiva, rotación, etc.
- objetos “similares” deben describirse de manera “similar

En este trabajo se utilizaron tres descriptores locales de las imágenes, uno de ellos utiliza información de intensidad (SIFT), otro de color (opponentSIFT) y el tercero de profundidad (FPFH).

2. Construir un diccionario

Cualquiera sea el descriptor utilizado, la cantidad de posibles valores que éstos pueden tomar es infinita. A fin de caracterizar una imagen por la frecuencia con que aparecen los descriptores locales es que se restringe la cantidad de posibles valores que pueden adquirir. El diccionario se construye de la siguiente manera

- Para un conjunto de imágenes representativo de los objetos presentes en la base se extraen los descriptores.
- Se elige un número de palabras para el diccionario y se realiza cuantización vectorial, en este caso se utilizó k-means como método de cuantización.

3. Asignar a los descriptores locales palabras pertenecientes al diccionario

Una vez calculados los descriptores locales se les debe asignar palabras que formen parte del diccionario generado en el paso anterior. El método más simple consiste en asignar a cada descriptor la palabra del diccionario más cercana. Este fue el método utilizado.

Finalmente se cuenta con un conjunto de histogramas ( uno por cada par imagen-descriptor) que en este caso serán utilizados como insumos de un clasificador SVM.

### III.b Características sobre las imágenes RGB. SIFT y Opponent-SIFT

El algoritmo SIFT fue publicado originalmente por David Lowe en 1999 [3] y es ampliamente utilizado para extraer características descriptivas de las imágenes. Estas características son invariantes a cambios de escala, traslación, rotación y parcialmente invariantes a cambios de iluminación y afinidades en la imagen.

Para este trabajo se utilizó la variante densa de SIFT que consiste en calcular el descriptor para una grilla uniforme de la imagen (ver figura 5). Esta variante, a diferencia de otras que primero buscan puntos claves en la imagen, asegura contar con un importante número de descriptores por imagen sin importar las características de la imagen.

Cada uno de los puntos de la grilla es el centro de un patch conformado por 4x4 subregiones. Lo que hace el algoritmo SIFT es calcular el gradiente de los píxeles pertenecientes a una misma subregión y forma con las orientaciones del gradiente un histograma de 8 bins para cada subregión. Luego concatena las orientaciones de los 16 píxeles obteniendo un descriptor del patch de dimensión 128.

Tanto la separación entre los puntos de la grilla como el tamaño de cada una de las subregiones son parámetros del algoritmo. Cuánto menor sea la separación entre los puntos de la grilla mayor será el número de patches por imagen. En este trabajo se fijó una separación de tres píxeles entre los puntos de la grilla y se utilizaron cuatro tamaños de subregiones por imagen, 4, 6, 8 y 10 píxeles.

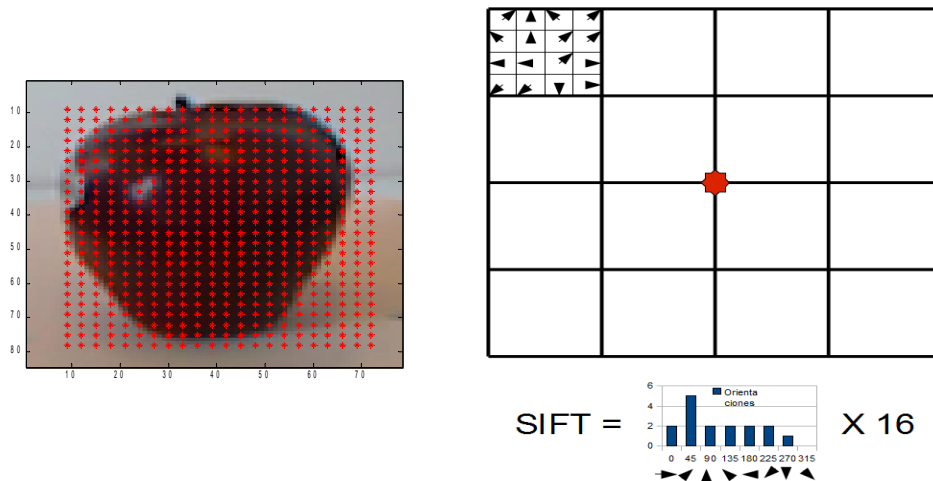


Fig. 5 A la izquierda se muestra un ejemplo del muestreo denso de las imágenes. Los puntos rojos son los centros de los patches dónde se calcula SIFT. A la derecha se esquematiza el cálculo del descriptor para un patch con subregiones de tamaño 4. Cada subregión produce un histograma de 8 bins, luego se concatenan los 16 histogramas para formar el descriptor de dimensión 128.

Como las imágenes con las que se trabajó tienen información de color también se probó una alternativa a SIFT que utiliza información de color llamada Opponent SIFT [22]. El vector descriptor se construye aplicando SIFT a cada uno de los tres canales del espacio de opuestos y luego concatenando los resultados. El espacio de opuestos se define como:

$$\begin{pmatrix} O_1 \\ O_2 \\ O_3 \end{pmatrix} = \begin{pmatrix} (R-G)/\sqrt{2} \\ (R+G-2B)/\sqrt{6} \\ (R+G+B)/\sqrt{3} \end{pmatrix}$$

Los descriptores Opponent SIFT son vectores de dimensión 384 en los que la información de intensidad está contenida en el tercer canal mientras que la de color en los dos primeros. Los canales uno y dos poseen

además la propiedad de ser invariantes a cambio en la intensidad de la luz ya que la resta de los canales hace que la variación en uno de ellos se cancele con la de los otros.

### **Agregado de información espacial**

Al utilizar el método de bag of features se pierde la información espacial presente en la ubicación de los patches, esto puede dar lugar a ambigüedades. Por ejemplo, las tres formas distintas de distribuir los trozos de manzana que se muestran a continuación producen el mismo histograma y por lo tanto a los efectos de un clasificador son iguales. Una forma de mitigar esta limitación es realizar histogramas en varios niveles.

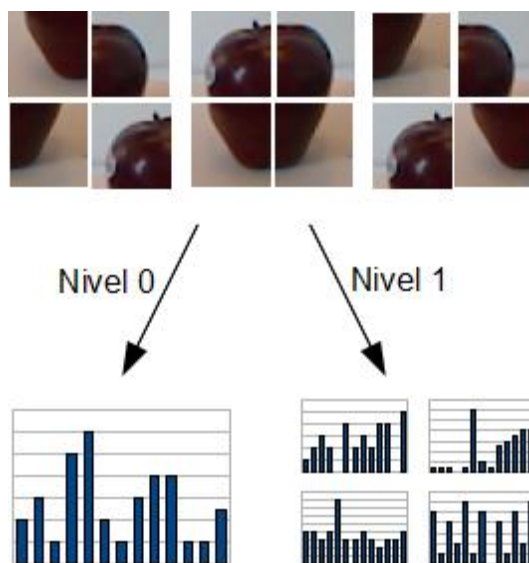


Fig 6. Descripción de la imagen mediante histogramas que conservan información espacial.

En la figura 6 se esquematiza una posible representación de la imagen en dos niveles. El primer nivel corresponde al histograma de la imagen entera mientras que en el segundo se divide la imagen en 4 y se realiza un histograma para cada nivel. El vector que caracteriza la imagen es la concatenación de todos los histogramas. En el ejemplo, si el diccionario utilizado tiene K palabras entonces la imagen se representa mediante un vector de dimensión  $K \times (1+4)$ .

### **III.c Características 3D – Fast Point Feature Histogram**

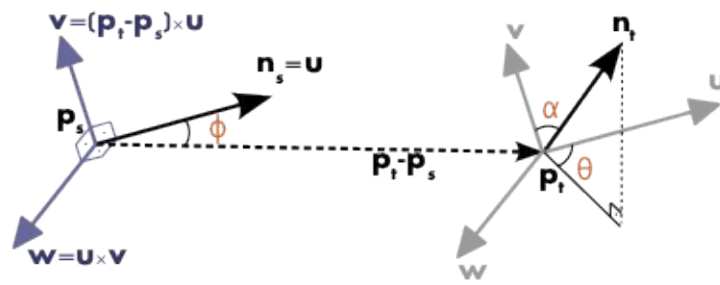
El Point Feature Histogram (PFH) [4, 21] codifica las propiedades geométricas de una nube en el entorno de un punto. Se basa en la relación entre los vecinos del punto en cuestión y las normales estimadas en los mismos, de manera de representar las variaciones de la superficie muestreada.

El PFH se calcula tomando la relación entre todos los pares de puntos en el vecindario del punto de interés. Entre cada par de puntos se calculan 4 valores: la distancia entre el par de puntos y las cantidades como se muestra en la figura 7.

El descriptor PFH para el punto de interés es el histograma de las cuartetos de todos los pares de puntos en la vecindad.

Para este trabajo se usó una variante de PFH de bajo costo computacional llamada Fast Point Feature Histogram (FPFH) [4]. El número de descriptores que se calculan por imagen es variable. Se fijó un máximo de 300 descriptores FPFH por imagen.





$$\begin{aligned}
 \mathbf{u} &= \mathbf{n}_s & \alpha &= \mathbf{v} \cdot \mathbf{n}_t \\
 \mathbf{v} &= \mathbf{u} \times \frac{(\mathbf{p}_t - \mathbf{p}_s)}{\|\mathbf{p}_t - \mathbf{p}_s\|_2} & \phi &= \mathbf{u} \cdot \frac{(\mathbf{p}_t - \mathbf{p}_s)}{d} \\
 \mathbf{w} &= \mathbf{u} \times \mathbf{v} & \theta &= \arctan(\mathbf{w} \cdot \mathbf{n}_t, \mathbf{u} \cdot \mathbf{n}_t)
 \end{aligned}$$

Fig. 7 Cálculo de la relación en el par de puntos  $P_s$  y  $P_t$  con normales estimadas  $n_s$  y  $n_t$ . Tomado de [7]

### III.d Clasificación con SVM

Los descriptores de color y profundidad se utilizaron para entrenar clasificadores SVM. Se realizaron pruebas tanto con SVM lineal como utilizando un kernel del tipo chi cuadrado. En el apéndice 1 se describen los fundamentos del funcionamiento de SVM. A continuación se describe como se utilizó la técnica de clasificación SVM en este trabajo.

Para cada una de las 51 categorías se eligieron hasta 4 instancias para entrenamiento y se reservó una para test del modelo SVM obtenido durante el entrenamiento.

Se evaluó el desempeño de los siguientes métodos de clasificación:

1. clasificación utilizando descriptores "puros" : se probó con SIFT, FPFH y OppositeSIFT. Para cada categoría se entrenó un clasificador uno contra todos. Al llegar una imagen no conocida hasta el momento se calculan los descriptores, se generan los histogramas y se la evalúa con todos los clasificadores.

Para cada categoría el clasificador devuelve un score que es la distancia al hiperplano de separación calculado en SVM. Este score da información de que tan probable es que la imagen pertenezca a su categoría. Un score de 1 o más indica que es muy probable que pertenezca, entre 0 y 1 que es probable y un score negativo indica que es muy poco probable que el patch pertenezca a la categoría.

A la imagen se le asocia la categoría que produjo un score más grande. A todas las imágenes se les asocia una categoría.

2. clasificación utilizando concatenación de vectores SIFT y FPFH: como insumo al clasificador SVM se usó un vector conformado por la concatenación del descriptor SIFT y del descriptor FPFH.
3. combinación de clasificadores utilizando decisión del experto. Para emplear este método se dividieron los patches correspondientes a las instancias de test en dos subconjuntos. Cada subconjunto contiene aproximadamente la mitad de los patches de cada instancia de test. El procedimiento de combinación es el siguiente:
  - Se crean los modelos SVM para ambos descriptores.
  - Se testea el rendimiento de los clasificadores utilizando uno de los subconjuntos reservados para test. El porcentaje de acierto para cada descriptor se guarda en un vector de dimensión igual al número de categorías (en la figura 8 se muestran bajo el símbolo de porcentaje)
  - Se calculan los scores para los descriptores que pertenecen al otro subconjunto de test utilizando SIFT y FPFH. (en la figura 8 son los vectores que aparecen debajo de la palabra "scores")

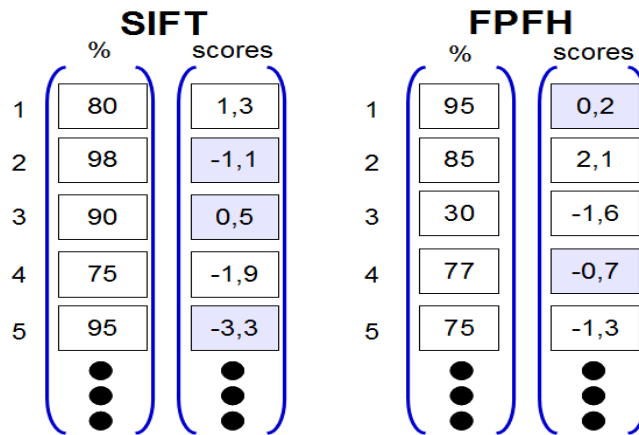


Fig. 8 Esquema que muestra los vectores que interviene en la decisión por método experto.

- Para cada categoría se elige el descriptor con mejor desempeño utilizando el primer conjunto . En el ejemplo se elije FPFH para la categoría 1, SIFT para la 2 y 3, FPFH para la cuatro, etc. Esto produce un nuevo vector de características de dimensión igual al número de categorías. Finalmente se asocia al descriptor la categoría con mayor score. En el ejemplo se asociaría al patch la categoría 3.

## **IV Software**

Para el desarrollo de software de ese proyecto se realizaron programas en Matlab y C utilizando las siguientes bibliotecas:

### **Point Cloud Library [7]**

Point Cloud Library (PCL) es un proyecto abierto destinado al procesamiento de nubes de puntos 3D. PCL cuenta con algoritmos que forman parte del estado del arte y permite realizar tareas como filtrado, estimación de características, reconstrucción y visualización de superficies, registrado y segmentación además de algunas herramientas de mapeo y reconocimiento de objetos. Es libre tanto para uso comercial como para investigación. Se encuentra financiado por Willow Garage, NVidia, Google, Toyota, and Trimble.

En este trabajo se utiliza la implementación de Fast Point Feature Histogram y la lectura y visualización de nubes de puntos.

### **VLFeat [8]**

VLFeat es una biblioteca libre que implementa algunos algoritmos clásicos de visión por computadora. En este trabajo se utilizaron las implementaciones de SIFT, OppositeSIFT, k-means y SVM presentes en la biblioteca. Está escrita en C y cuenta con interfaces para MATLAB que facilitan su utilización. Soporta Windows, Mac OS X y linux. La versión utilizada fue la 0.9.13.

### **Corridas sobre la base de imágenes RGB-D dataset**

Las corridas se realizaron tomando como base el ejemplo “phow\_caltech101” disponible en la biblioteca VLFeat [8]. En dicho ejemplo se realiza entrenamiento y clasificación en base a características SIFT sobre las imágenes de la base Caltech 101 [10].

## **Aplicación**

Para probar resultados más allá de las corridas sobre la base de imágenes, se implementó una pequeña aplicación en Matlab con las siguientes funcionalidades:

- Adquisición desde sensor Kinect
- Armado de base de imágenes
- Entrenamiento con la base creada
- Clasificación utilizando características SIFT, FPFH, combinación de ambas
- Clasificación mediante modelo precargado

La figura 9 muestra una vista de la interfaz gráfica de la aplicación.

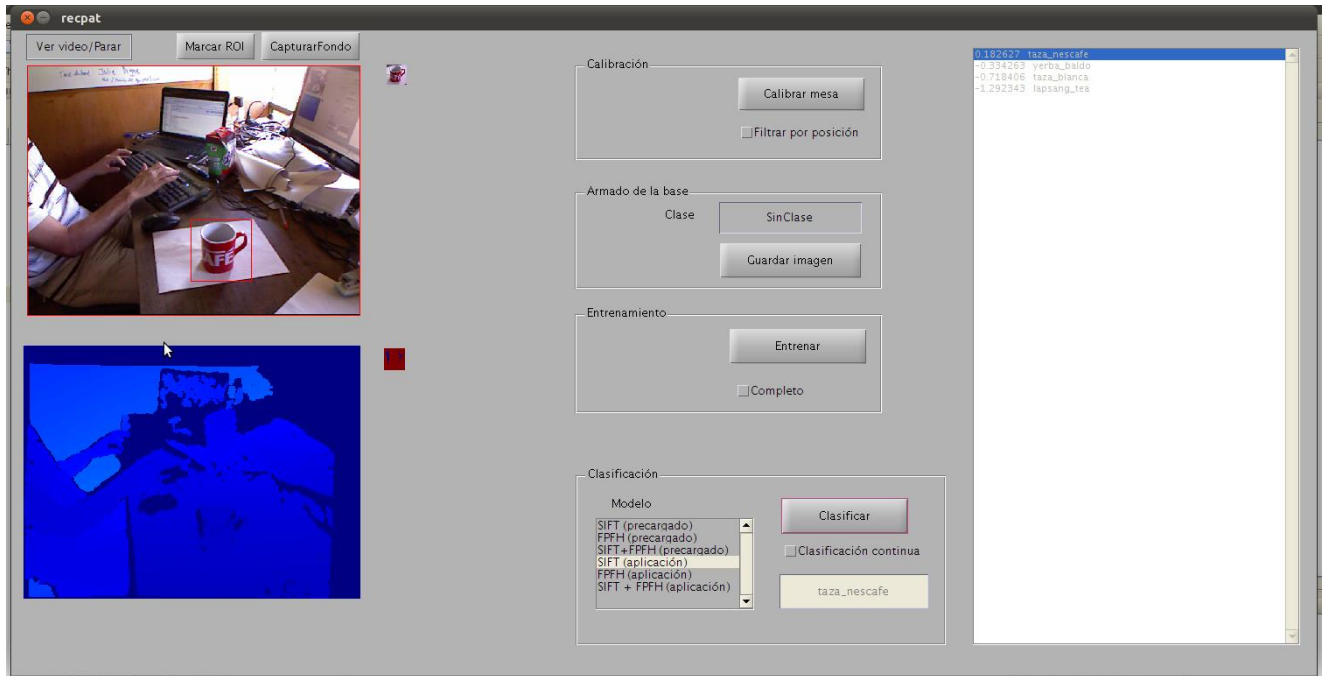


Fig. 9 Interfaz gráfica de la aplicación

## V Experimentos y resultados

Se realizaron pruebas tendientes a evaluar el desempeño de las distintas variantes de clasificación y la influencia de los distintos elementos que conforman el esquema de clasificación: tipo de descriptor, número de palabras del diccionario, información espacial de los histogramas, método de clasificación.

Las pruebas se realizaron con las 51 categorías de la base y se utilizaron hasta cinco instancias por clase. Una de las instancias se usó para test mientras que la restantes se utilizaron para entrenamiento.

- **V.a Tipo de descriptor**

Se evaluó el desempeño de SIFT, Opponent SIFT y FPFH individualmente. Para la prueba se dejó la primera instancia de cada clase para test (7019 imágenes) y se entrenó con las restantes (25122 imágenes).

- **V.a.1. SIFT**

Se dividió cada una de las imágenes en cuatro regiones y se aplicó SIFT a cada una de ellas. El descriptor SIFT de la imagen es la concatenación de los cuatro histogramas. El número de palabras utilizado fue 600. La matriz de confusión obtenida al dejar la primera instancia de cada categoría para test se muestran en la tabla 2.

- **V.a.2. Opponent SIFT**

Se realizó la misma prueba que en V.a.1 pero utilizando OpponentSIFT en lugar de SIFT. En la tabla 6 se observan los resultados. A continuación se comparan los desempeños.

Descriptor	Desempeño
SIFT	75,94%
OpponentSIFT	74,61%

- **V.a.3. FPFH**

El número de palabras utilizado para formar el diccionario fue 300. Los histogramas se realizaron utilizando la imagen completa. El resultado de la prueba se muestra en la tabla 1.

- **V.b Número de palabras del diccionario**

Para evaluar la influencia del diccionario se utilizó el descriptor SIFT. Se realizaron pruebas utilizando 100, 300 y 600 palabras. Para la prueba se dejó la primera instancia de cada clase para test (7019 imágenes) y se entrenó con las restantes (25122 imágenes).

Número de Palabras	Desempeño
100 (2x2)	71,70%
300 (2x2)	75,72%
600 (2x2)	75,94%

- **V.c Información espacial en los histogramas**

Para evaluar la influencia del diccionario se utilizó el descriptor SIFT. Se dejó la primera instancia de cada clase para test (7019 imágenes) y se entrenó con las restantes (25122 imágenes). Se ejecutaron las siguientes pruebas con un diccionario de 600 palabras:

- utilizando histograma de imagen completa.
- dividiendo la imagen en cuatro regiones y concatenando los histogramas
- concatenando el histograma de imagen completa con los histogramas de las subregiones.

División Espacial	Desempeño
imagen completa	72,83%
600 (2x2)	75,94%
600 (1 + 2x2)	76,81%

Se observa que el rendimiento de SIFT mejora un poco aunque la mejora no parece ser significativa. Eso se debe a que los objetos abarcan una gran parte de la imagen.

- **V.d Método de combinación: concatenación de características y combinación de clasificadores**

Número de Palabras	SIFT	FPFH	Concatenación	Experto
100 (2x2)	71,70%	54,82%	78,57%	76,05%
300 (2x2)	75,72%	54,82%	80,98%	79,40%
600 (2x2)	75,94%	54,82%	81,42%	78,20%

- **V.e SVM lineal – SVM chi2**

Para evaluar la influencia del kernel en la clasificación se entrenó un modelo SVM lineal y uno que utiliza kernel del tipo chi cuadrado. El kernel chi2 ha sido reportado como una de las mejores opciones para histogramas y su mapeo puede calcularse en forma eficiente [14, 23]. Se utilizó como insumo al clasificador la concatenación de los descriptores SIFT y FPFH.

Tipo de kernel	Rendimiento
lineal	72,40%
Chi cuadrado	80,99%

Los resultados obtenidos utilizando el kernel son significativamente mejores que los que se logran con SVM lineal. Las matrices de confusión se muestran en las tablas 4 y 5 (svm lineal y utilizando kernel respectivamente) .

- **V.f Comparación con la publicación de referencia**

Si se compara el rendimiento de los descriptores de forma se observa que los resultados obtenidos por los autores de la base es significativamente mayor que los que se obtuvieron en este trabajo. Los autores de la base reportan una media de 64,7% utilizando SPIN mientras que en este trabajo se logró un 54,82% en una sola corrida.

En lo que respecta a los descriptores que utilizan información rgb los resultados aquí obtenidos son levemente mejores, en el entorno de un 3%.

Cuando se combinan ambas fuentes de información los resultados que se obtienen también son del orden. Los autores de la base reportan un 83,8% con una desviación estándar de 3,5% utilizando kernel SVM y aquí se obtuvo una media de 82,1%.

Los resultados se obtuvieron promediando 10 corridas. En cada corrida se reservó una instancia por categoría para test mediante sorteo. La incertidumbre que se indica equivale a dos desviaciones estándar.

	Autores de la base	Grupo de Recpat 2011
Descriptor de forma	64,7 +/- 2,2	54,82
Descriptor rgb	74,5 +/- 3,1	64,7 +/- 2,5
Descriptor forma + descriptor rgb	83,8 +/- 3,5	82,1 +/- 3,9















## **VI Conclusiones**

Combinar la información de color con la de profundidad contribuye a mejorar el desempeño del reconocimiento de categorías. Un claro ejemplo de ello es la categoría bol (ver figura 3), su forma característica hace que el aprovechamiento de dicha información sea relevante para la detección de la categoría.

Si bien el desempeño que se obtuvo utilizando el descriptor de forma es menor que el que reportan los autores de la base, vale mencionar que el menor desempeño del descriptor de forma aquí utilizado no necesariamente responde a la menor bondad de FPFH respecto a SPIN. Los autores de la base dividen el espacio ocupado por el objeto a reconocer en una grilla de 3x3x3 y calculan 1000 descriptores SPIN por región. Luego realizan PCA a los descriptores de cada región y se quedan con las 100 primeros componentes. Finalmente concatenan los 27 descriptores y obtienen un descriptor de forma de tamaño 2700. En este trabajo se utilizaron tan solo 300 descriptores por imagen, es de esperar que un mayor número de descriptores por imagen y/o el aprovechamiento de la información espacial mejore el desempeño del descriptor de forma.

Los resultados que aquí se obtuvieron al combinar forma con rgb son muy similares a los reportados por los autores de la base. Eso quiere decir que pese a que no se utiliza un gran número de descriptores por imagen ni información espacial el descriptor FPFH logra en gran medida extraer los detalles necesarios para complementar la información de intensidad de la imagen.

# Apéndice 1

## SVM

SVM (Support Vector Machines) es una herramienta de aprendizaje supervisado sumamente utilizada debido a su alta capacidad de discriminación entre clases de datos. Lo que hace SVM es encontrar el hiperplano que mejor separa un conjunto de datos. Gracias al truco de los kernels, este funcionamiento básico se generaliza y se logran bordes de decisión no lineales en el espacio original. A continuación se describe brevemente el funcionamiento de SVM cuando el problema de clasificación es linealmente separable, cuando no lo es y se utiliza un clasificador lineal y cuando se utiliza un borde de decisión no lineal .

### 1. Problema separable linealmente

Intuitivamente, parece razonable encontrar el hiperplano que maximizar el margen, es decir la distancia mínima de un patrón al hiperplano. Cuanto mayor es el margen, mayor es la tolerancia al ruido y por lo tanto mayor la capacidad de generalización. SVM realiza algo parecido a esto. SVM encuentra el hiperplano donde la distancia al hiperplano desde los puntos más cercanos al mismo es lo más grande posible.

Se busca el hiperplano separador  $\mathbf{w}^t \mathbf{x} + b = 0$  que maximiza el margen. La distancia de un punto  $\mathbf{x}_i$  al hiperplano es  $|\mathbf{w}^t \mathbf{x}_i + b| / \|\mathbf{w}\|$  .

Si se considera el hiperplano canónico, es decir aquel que cumple  $|\mathbf{w}^t \mathbf{x}_i + b| = 1$  para los puntos más cercanos entonces la distancia desde estos puntos al hiperplano es  $\|\mathbf{w}\|^{-1}$  y por lo tanto el margen es:

$$m = 2 \times \|\mathbf{w}\|^{-1}$$

Maximizar el margen es por tanto equivalente al siguiente problema de optimización:

$$\min_{\mathbf{w}, b} J(\mathbf{w}) = 1/2 \|\mathbf{w}\|^2 \text{ sujeto a } t_i(\mathbf{w}^t \mathbf{x}_i + b) \geq 1 \forall i .$$

Para resolver este problema se introduce un multiplicador de Lagrange  $\lambda_i \geq 0$  por cada restricción y se obtiene un problema sin restricciones en el que se minimiza el lagrangiano

$$L(\mathbf{w}, b, \lambda) = 1/2 \|\mathbf{w}\|^2 - \sum_{i=1}^n \lambda_i [t_i(\mathbf{w}^t \mathbf{x}_i + b) - 1]$$

Derivando con respecto a  $\|\mathbf{w}\|$  e igualando a 0 se obtiene:  $\mathbf{w} = \sum_{i=1}^n \lambda_i t_i \mathbf{x}_i$

Derivando con respecto a  $b$  e igualando a 0 se llega a:  $0 = \sum_{i=1}^n \lambda_i t_i$

Sustituyendo estas expresiones en la expresión del Lagrangiano se obtiene el problema dual.

$$\max_{\lambda} \tilde{L}(\lambda) = \sum_{i=1}^n \lambda_i - 1/2 \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j t_i t_j \mathbf{x}_i \cdot \mathbf{x}_j$$

$$\text{con las restricciones: } \lambda_i \geq 0 \quad \text{y} \quad \sum_{i=1}^n \lambda_i t_i = 0$$

El problema dual satisface las condiciones de KKT:

$$\lambda_i \geq 0$$

$$t_i(\mathbf{w}^t \mathbf{x}_i + b) \geq 1 \quad \forall i$$

$$\lambda_i [t_i(\mathbf{w}^t \mathbf{x}_i + b) - 1] = 0$$

Por lo tanto para cada patrón  $\mathbf{x}_i$  se tienen dos casos:

- Si  $\lambda_i = 0$  estos puntos no contribuyen a definir el hiperplano separador.
- Si  $t_i(\mathbf{w}^t \mathbf{x}_i + b) = 1$  estos puntos definen el hiperplano separador y se denominan vectores de soporte. Sólo los vectores de soporte contribuyen a definir el hiperplano separador óptimo.

El vector  $\mathbf{w}$  se define como  $\mathbf{w} = \sum_{i=1}^n \lambda_i t_i \mathbf{x}_i$

Mientras que el  $b$  se obtiene a partir de  $t_i(\mathbf{w}^t \mathbf{x}_i + b) = 1$

Sólo los vectores de soporte son necesarios para clasificar. La función de clasificación es

$$f(\mathbf{x}) = \sum_{i=1}^n \lambda_i t_i \mathbf{x}_i \cdot \mathbf{x} + b$$

## 2. Problema no separable linealmente

A continuación se introducen las variables  $\xi_i \geq 0$ . Las restricciones son ahora de la forma  $t_i(\mathbf{w}^t \mathbf{x}_i + b) \geq 1 - \xi_i$  con:

- $\xi_i = 0$  si los puntos están bien clasificados y además se encuentran fuera del margen.
- $0 \leq \xi_i \leq 1$  Si los puntos están bien clasificados pero caen dentro del margen.
- $\xi_i = 1$  si los puntos están mal clasificados

El objetivo será ahora maximizar el margen penalizando a su vez los ejemplos mal clasificados. El problema de



optimización se escribe ahora cómo

$$\min_{\mathbf{w}, b, \xi} J(\mathbf{w}, \xi) = C \sum_{i=1}^n \xi_i + 1/2 \|\mathbf{w}\|^2$$

$$\text{con las restricciones } t_i(\mathbf{w}^t \mathbf{x}_i + b) \geq 1 - \xi_i \text{ y } \xi_i \geq 0 .$$

El parámetro C controla el peso relativo que se le da al error de clasificación en entrenamiento y a la complejidad (margen). Un valor grande de C favorece modelos con menos error mientras que C pequeños favorecen modelos más simples. Para resolver este problema se introducen, al igual que en el caso anterior los multiplicadores de Lagrange, en este caso  $\lambda_i$  y  $\mu_i$

$$L(\mathbf{w}, b, \lambda, \mu) = 1/2 \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \lambda_i [t_i(\mathbf{w}^t \mathbf{x}_i + b) - 1 + \xi_i] - \sum_{i=1}^n \mu_i \xi_i$$

Las condiciones de KKT son ahora:

$$\lambda_i \geq 0$$

$$t_i(\mathbf{w}^t \mathbf{x}_i + b) - 1 + \xi_i \geq 0$$

$$\lambda_i [t_i(\mathbf{w}^t \mathbf{x}_i + b) - 1 + \xi_i] = 0$$

$$\mu_i \geq 0$$

$$\xi_i \geq 0$$

$$\mu_i \xi_i = 0$$

$$\text{Derivando con respecto a } \mathbf{w} \text{ e igualando a 0: } \mathbf{w} = \sum_{i=1}^n \lambda_i t_i \mathbf{x}_i$$

$$\text{Derivando con respecto a } b \text{ e igualando a 0: } 0 = \sum_{i=1}^n \lambda_i t_i$$

$$\text{Derivando con respecto a } \xi_i \text{ e igualando a 0: } \lambda_i = C - \mu_i$$

Sustituyendo estas expresiones en L se obtiene el problema dual.

$$\max_{\lambda} \tilde{L}(\lambda) = \sum_{i=1}^n \lambda_i - 1/2 \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j t_i t_j \mathbf{x}_i \mathbf{x}_j$$

$$\text{con las restricciones: } 0 \leq \lambda_i \leq C \text{ y } \sum_{i=1}^n \lambda_i t_i = 0$$

La única diferencia con el caso anterior reside en la primera restricción.

Nuevamente para cada patrón  $\mathbf{x}_i$  se tienen dos casos:

- Si  $\lambda_i=0$  estos puntos están bien clasificados y fuera del margen.
- Si  $t_i(\mathbf{w}^t \mathbf{x}_i + b) = 1 - \xi_i$  estos puntos son vectores de soporte.

Los vectores de soporte se dividen en dos casos:

- Si  $\lambda_i < C$  entonces  $\mu_i > 0$  y  $\xi_i = 0$ . Los puntos están sobre el margen.
- Si  $\lambda_i = C$  entonces  $\mu_i = 0$  y  $\xi_i > 0$ . Los puntos están dentro del margen. Se encuentran bien clasificados si  $\xi_i \leq 1$  y mal clasificados si  $\xi_i > 1$

Nuevamente el vector  $\mathbf{w}$  se define cómo  $\mathbf{w} = \sum_{i=1}^n \lambda_i t_i \mathbf{x}_i$

Mientras que el  $b$  se obtiene a partir de  $t_i(\mathbf{w}^t \mathbf{x}_i + b) = 1$

### 3. Problema no lineal

Para resolver problemas no separables linealmente se puede adoptar una estrategia que consiste en dos etapas

1. Realizar una proyección lineal de los datos sobre un espacio de alta dimensión.
2. Buscar el hiperplano separador óptimo en el nuevo espacio

Al proyectar a un espacio de alta dimensión se tiene mayor probabilidad de que las clases sean separables. El inconveniente del alto costo computacional que implica trabajar en un espacio de alta dimensionalidad se evita gracias al popular truco de los kernels. Gracias a este elegante truco se realiza la proyección al espacio de alta dimensionalidad sólo de manera implícita. Análogamente al caso anterior el vector  $\mathbf{w}$  se define cómo

$\mathbf{w} = \sum_{i=1}^n \lambda_i t_i \Phi(\mathbf{x}_i)$  siendo  $\Phi(\mathbf{x}_i)$  la proyección del vector  $\mathbf{x}_i$  a la alta dimensión. Para clasificar un nuevo patrón basta con evaluar

$$f(\mathbf{x}) = \sum_{i=1}^n \lambda_i t_i \Phi(\mathbf{x}_i) \Phi(\mathbf{x}) + b = \sum_{i=1}^n \lambda_i t_i k(\mathbf{x}_i, \mathbf{x}) + b$$

## Bibliografía

- [1] Richard Szeliski  
Computer Vision, Algorithms and Applications  
Springer-Verlag London Limited 2011  
ISBN 978-1-84882-934-3
- [2] M. Turk and A. Pentland  
"Face recognition using eigenfaces". *Proc. IEEE Conference on Computer Vision and Pattern Recognition, 1991*. pp. 586–591.
- [3] Lowe, D. G.,  
"Object recognition from local scale-invariant features", International Conference on Computer Vision, Corfu, Greece, September 1999.
- [4] Rusu, Radu Bogdan  
Semantic 3D Object Maps for Everyday Manipulation in Human Living Environments  
PhD Thesis - Computer Science department, Technische Universitaet Muenchen, Germany  
October 2009  
<http://files.rbrusu.com/publications/RusuPhDThesis.pdf>  
Acceso en diciembre de 2011
- [5] A Large-Scale Hierarchical Multi-View RGB-D Object Dataset  
Kevin Lai, Liefeng Bo, Xiaofeng Ren, and Dieter Fox  
IEEE International Conference on Robotics and Automation (ICRA), May 2011.
- [6] K. Lai, L. Bo, X. Ren, and D. Fox  
Sparse Distance Learning for Object Recognition Combining RGB and Depth Information  
IEEE International Conference on Robotics and Automation (ICRA), May 2011.  
Best Vision Paper Award.
- [7] Point Cloud Library  
<http://pointclouds.org/>  
Acceso en diciembre de 2011
- [8] Andrea Vedaldi, Brian Fulkerson  
VL-Feat  
<http://www.vlfeat.org/index.html>  
Acceso en diciembre de 2011
- [9] Sebastien Paris  
Scenes/Objects classification toolbox  
<http://www.mathworks.com/matlabcentral/fileexchange/29800-scenesobjects-classification-toolbox>  
Acceso en diciembre de 2011
- [10] L. Fei-Fei, R. Fergus and P. Perona.  
*Learning generative visual models from few training examples: an incremental Bayesian approach tested on 101 object categories*. IEEE. CVPR 2004, Workshop on Generative-Model Based Vision. 2004  
[http://www.vision.caltech.edu/Image\\_Datasets/Caltech101/](http://www.vision.caltech.edu/Image_Datasets/Caltech101/)  
Acceso en diciembre de 2011

- [11] Svetlana Lazebnik, Cordelia Schmid, and Jean Ponce. [Beyond Bags of Features: Spatial Pyramid Matching for Recognizing Natural Scene Categories](#). Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, New York, June 2006, vol. II, pp. 2169-2178.
- [12] K. Grauman and T. Darrell. Pyramid match kernels: Discriminative classification with sets of image features. In Proc. ICCV, 2005.
- [13] G. Csurka, C. Dance, L.X. Fan, J. Willamowski, and C. Bray (2004). "Visual categorization with bags of keypoints". *Proc. of ECCV International Workshop on Statistical Learning in Computer Vision*.
- [14] Perronnin, F., Sanchez, J, Liu, Yan.. (2010). Large-Scale Image Categorization with Explicit Data Embedding. Learning, 2297-2304.
- [15] C. Wallraven, B. Caputo, and A. Graf. Recognition with local features: the kernel recipe. In Proc. ICCV, volume 1, pages 257–264, 2003.
- [16] J. Willamowski, D. Arregui, G. Csurka, C. R. Dance, and L. Fan. Categorizing nine visual classes using local appearance descriptors. In ICPR Workshop on Learning for Adaptable Visual Systems, 2004.
- [17] Using Spin Images for Efficient Object Recognition in Cluttered 3D Scenes  
Andrew E. Johnson, Member, IEEE, and Martial Hebert, Member, IEEE
- [18] Liefeng Bo and Cristian Sminchisescu, Efficient Match Kernels between Sets of Features for Visual Recognition, Advances in Neural Information Processing Systems (NIPS), December, 2009.
- [19] <http://www.primesense.com>
- [20] <http://people.csail.mit.edu/fergus/iccv2005/bagwords.html>
- [21] [http://www.pointclouds.org/documentation/tutorials/pfh\\_estimation.php#pfh-estimation](http://www.pointclouds.org/documentation/tutorials/pfh_estimation.php#pfh-estimation)
- [22] Koen E. A. van de Sande, Theo Gevers and Cees G. M. Snoek, "Evaluating Color Descriptors for Object and Scene Recognition", IEEE Transactions on Pattern Analysis and Machine Intelligence, volume 32 (9), pages 1582-1596, 2010
- [23] A.Vedaldi and A.Zisserman, Efficient Additive Kernels via Explicit Feature Maps in *Proceedings of the IEEE Conf. on Computer Vision and Pattern Recognition(CVPR)*, 2010.