



Presentación:

Instituto Tecnológico de las Américas

Grupo:

#8

Materia:

Programación 3

Sección:

G4

Día/hora:

Martes y jueves / 2:00 PM – 4:00 PM

Maestro:

Willis Polanco

Fecha de entrega:

Viernes 17 de abril del 2021

Tema:

Análisis y diseño del proyecto

El reconocimiento facial:

El reconocimiento facial (RF) automatizado por computadora es un área de investigación en pleno desarrollo. Identificar a una persona por medio de una imagen de su rostro implica emular el proceso cognitivo que realiza un ser humano al reconocer a sus pares. Aunque existen trabajos en el campo de la psicofísica y la neurociencia, aún no se sabe con certeza cómo funcionan estos procesos internamente en el cerebro humano.

Independientemente de la técnica utilizada para la solución, el RF requiere de tres etapas: 1) la detección del rostro en una imagen, 2) la extracción de características y 3) la identificación y/o verificación de la cara mediante la clasificación de las características. En la identificación, el sistema informa la identidad de la persona, mientras que en la verificación confirma o rechaza esta. Cada una de estas áreas ha sido objeto de investigaciones, dando origen a diferentes técnicas para resolver cada problema.

Detección del rostro

La detección implica encontrar las áreas dentro de una imagen que contienen un rostro. Básicamente se trata de descartar todo lo que sea fondo, y así obtener la ubicación y tamaño exacto de la cara.

Uno de los primeros métodos para resolver este problema fue el de plantillas deformables. Cada característica del rostro (ojos y boca en este caso) es representada por medio de una plantilla. La idea general es tener una función de energía que relacione las intensidades de la imagen (valles y picos) con las propiedades de la plantilla. Una vez ubicada la plantilla se ajustan sus parámetros para minimizar la función de energía.

Otros enfoques se vuelcan a la clasificación de patrones por medio de ejemplos [SP98]. Como sólo hay que detectar caras (sin identificar sujetos específicamente) se crean patrones para zonas que representan una cara y zonas que no. Generalmente la solución para encontrar patrones que no son rostros, es por medio de la modalidad denominada “*bootstrap*”, en la cual se comienza con patrones que son rostros, se pone en marcha el clasificador, y todos los falsos positivos se agregan como patrones de *no-rostro*. En una siguiente pasada, si aparecen nuevos falsos positivos, éstos se siguen agregando como patrones.

Cabe destacar que el método de clasificación es independiente de las imágenes utilizadas como patrones. En [SP98] se utiliza una red neuronal artificial como clasificador, específicamente un perceptrón multicapa.

Los últimos trabajos se basan en el método propuesto en [VJ04], que alcanza las tasas de reconocimiento de los mejores sistemas, pero con la ventaja de obtener una alta velocidad de detección. Una de sus contribuciones son los clasificadores en cascada, que permiten deshacerse de las regiones de fondo lo antes posible. Varios clasificadores son utilizados en serie, los primeros descartan la mayoría de regiones que no son propensas a contener un rostro. Las regiones más prometedoras pasan a los clasificadores más complejos y que requieren de mayor tiempo de cálculo.

1.1.2 Extracción de características

La extracción de características se refiere a la obtención de propiedades o parámetros particulares de cada rostro para luego poder ser clasificados. Se pueden tomar tres enfoques diferentes: 1) un enfoque holístico, basándose en la imagen del rostro como un todo, 2) un enfoque mediante características

1.1. Estado del arte

locales, dando mayor importancia a las diferentes partes del rostro, o 3) un enfoque híbrido basado en la idea de que el sistema de percepción humana combina características locales y globales para el reconocimiento [ZCPR03].

En el enfoque holístico, la técnica más influyente es la de las *eigenfaces* [TP91]. Se utiliza un análisis de componentes principales para representar la imagen completa del rostro en un espacio de dimensiones reducido. Esta técnica no necesita de grandes cantidades de fotos por sujeto, tiene una baja sensibilidad al ruido, y aunque presenta problemas para manejar las variaciones de iluminación, pueden realizarse pequeños cambios para mejorar la efectividad en dichas condiciones [BHK97].

Por otro lado, se encuentran los modelos basados en las características locales. Uno de los métodos de mayor éxito es el de *Elastic Bunch Graph Matching* [WFKM97]. Aquí, los rostros son representados en forma de grafos. Los puntos importantes en la cara que conforman los nodos del grafo (como los ojos o la nariz), son descritos por un conjunto de onditas de Gabor de diferentes escalas y rotaciones, logrando robustez para cambios en iluminación, traslación, distorsión, rotación y escalado [ZJN07].

Los modelos híbridos son adaptaciones del método de las *eigenfaces* aplicado a características locales [PMS94]. Así, se combinan ventajas de los métodos holísticos y las características locales.

1.1.3 Reconocimiento

El reconocimiento consiste en clasificar las características extraídas de cada rostro. Esta clasificación puede ser realizada de manera *supervisada*, en la cual un patrón de entrada es identificado como miembro de una clase predefinida, o de manera *no supervisada*, donde el patrón es asignado a una clase desconocida. Aquí, cada clase es un sujeto, por lo tanto, al clasificar las características se está indicando a qué sujeto pertenecen.

Para el diseño de clasificadores se pueden distinguir tres aproximaciones basadas en: 1) concepto de similaridad, 2)

aproximación probabilística y 3) optimización de un criterio de error [JDM00].

Las aproximaciones basadas en el concepto de similaridad son las más simples e intuitivas, donde los patrones similares son asignados a la misma clase. Para ello, se establece una métrica que define la similaridad y luego se clasifica por medio de plantillas o mínima distancia usando uno o varios prototipos por clase. La técnica de *eigenfaces* original [TP91] aplica la regla del vecino más cercano, utilizando como métrica la distancia Euclídea, donde

1.2. Propuesta del proyecto

cada prototipo es la media de los patrones de entrenamiento.

En el enfoque probabilístico, se utilizan los conceptos de la teoría de la decisión estadística para establecer los bordes de decisión de las diferentes clases. Se asume que las características, que representan a un patrón, tienen una función de densidad de probabilidad condicionada a la clase. Así, un vector patrón \mathbf{x} perteneciente a la clase ω_i es visto como una observación aleatoria desde la función de probabilidad $p(\mathbf{x}|\omega_i)$. Las reglas de decisión más conocidas son la *regla de decisión de Bayes* y la *regla de máxima probabilidad*.

Por último, la tercer aproximación se basa en construir los bordes de decisión optimizando algún criterio de error. El objetivo es minimizar el error de clasificación entre la respuesta deseada y la salida del clasificador. Un ejemplo de este tipo de clasificadores son las *redes neuronales* [JMM96], las cuales pueden considerarse como sistemas distribuidos y paralelos que consisten de pequeñas unidades de procesamiento masivamente conectadas. Están conformadas por redes de grafos ponderados donde los nodos son las neuronas artificiales y los bordes (con pesos) son las conexiones entre las neuronas de entrada y salida de la red. Para lograr una clasificación

adecuada son entrenadas con un algoritmo de entrenamiento a partir de un conjunto de datos. Los tipos de redes más comúnmente usados son las redes hacia adelante, como por ejemplo, el perceptrón multicapa.

1.2 Propuesta del proyecto

La idea central de este proyecto consiste en realizar el diseño y la implementación de un sistema integral para identificación de personas mediante el reconocimiento del rostro. El ambiente de aplicación del sistema será el Centro de Investigación y Desarrollo en Señales, Sistemas e Inteligencia Computacional (*sinc(i)*) de la FICH, UNL y el Centro de Investigación y Desarrollo de Ingeniería en Sistemas de Información (CIDISI) de la UTN-FRSF.

Para ello, se desarrolló una interfaz de usuario y se instaló una cámara web en una PC de escritorio en ambos centros de investigación para la captura de imágenes. Con el conjunto de imágenes obtenidas se realizó una clasificación manual de las mismas conformando una base de rostros para las pruebas del sistema.

Ya con la base de rostros, se desarrolló el sistema en base a las tres etapas básicas de todo sistema de reconocimiento facial:

1.3. Objetivos

1. Localización del rostro: el cual ubica la cara en la imagen aplicando el algoritmo propuesto en [ML06], basado en una segmentación de piel para detectar posibles regiones caras, y luego mediante combinación de plantillas decide por una región, recortándola como posible rostro.
2. Extracción de características: en esta etapa se obtienen las propiedades más relevantes del rostro mediante el método

eigenfaces, obteniendo como resultado un vector característico que representa la imagen de entrada en un espacio dimensional reducido para su posterior clasificación.

3. Clasificación: aquí se diseñó una red neuronal para el reconocimiento de patrones, más específicamente un perceptrón multicapa. La red se entrenó mediante el algoritmo de retropropagación del error a partir de un conjunto de patrones de ejemplo. Con esto, se logró que la red “aprenda” sobre el conjunto de potenciales usuarios.

Cabe destacar dos fases del proceso de construcción: una de entrenamiento y otra de clasificación. En la primera se ajustan varios parámetros del sistema, siendo el paso más importante el entrenamiento de la red para el conjunto de datos especificado. Una vez logrado esto, se pasa a la fase de clasificación verificando el desempeño mediante la presentación de imágenes rostros de usuarios tanto conocidos como desconocidos.

1.3 Objetivos

Objetivos generales

- Desarrollar un prototipo de herramienta computacional de código abierto que permita realizar reconocimiento facial, basado en la utilización de técnicas del área de Procesamiento de Imágenes e Inteligencia Computacional.
- Aplicar los conocimientos adquiridos en el transcurso de la carrera a un proyecto que realice un aporte ingenieril a la comunidad en general.

-
- Colaborar en proyectos de investigación que requieran la participación de profesionales informáticos en el área de procesamiento de imágenes digitales e inteligencia artificial.

1.4. Alcances

Objetivos específicos

- Implementar los módulos del localizador, normalizador de caras y extractor de características a partir de técnicas de procesamiento digital de imágenes.
- Diseñar e implementar el módulo del reconocedor mediante técnicas de la Inteligencia Computacional.
- Desarrollar un prototipo que permita la prueba del sistema de RF. El software ofrecerá las siguientes funcionalidades.
 - Tomar fotografías en vivo desde una cámara.
 - Registrar y almacenar individuos en el sistema, asociándole una o más imágenes capturadas.
 - Soportar el RF mediante la comparación de una fotografía actual con las registradas en el sistema. Se proveerán dos modos de funcionamiento: a) hacer corresponder la imagen tomada a la persona con la que más se ajuste de las conocidas en el sistema, y b) verificar la identidad del usuario por otro método (por ej. el ingreso del nombre y/o DNI por teclado).
 - Exportar e importar la información del individuo almacenada en el sistema. Esto permitiría transportar su información para ser utilizada en otra ubicación.

1.4 Alcances

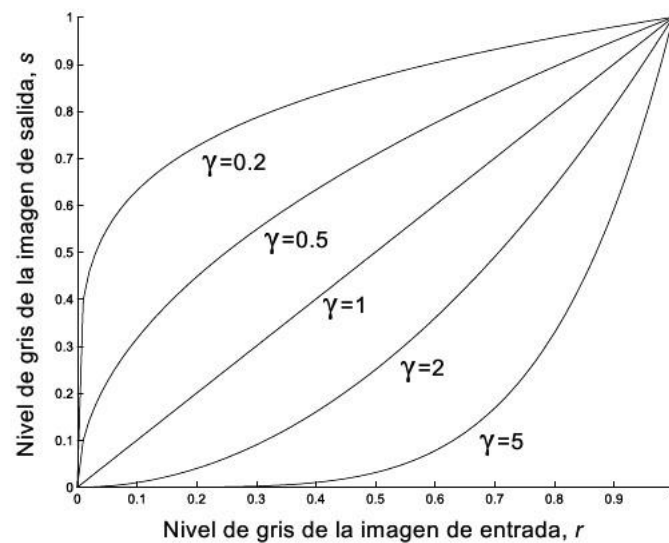
- La cantidad de usuarios se limitará a 11 en esta etapa de desarrollo del prototipo.
- No se hará rechazo de impostores, sino que se supondrá que todos los usuarios del sistema se encuentran registrados.
- No incluirá el desarrollo del módulo de administración de usuarios.

Marco teórico

2.1.1 Imagen digital

Una imagen en escala de grises puede ser definida como una función bidimensional $f(x,y)$, donde x e y son coordenadas *espaciales*, y el valor de f en un par de coordenadas dado es denominado *intensidad* o *nivel de gris*. Cuando x , y y f son cantidades discretas, la imagen es llamada *imagen digital* [GW02]. De este modo, ésta puede ser representada como una matriz de $M \times N$ elementos llamados *píxeles*.

$$f(x,y) = \begin{bmatrix} f(0,0) & f(0,1) & \dots & f(0,N-1) \\ f(1,0) & f(1,1) & \dots & f(1,N-1) \\ \vdots & \vdots & \ddots & \vdots \\ f(M-1,0) & f(M-1,1) & \dots & f(M-1,N-1) \end{bmatrix} \quad (2.1)$$



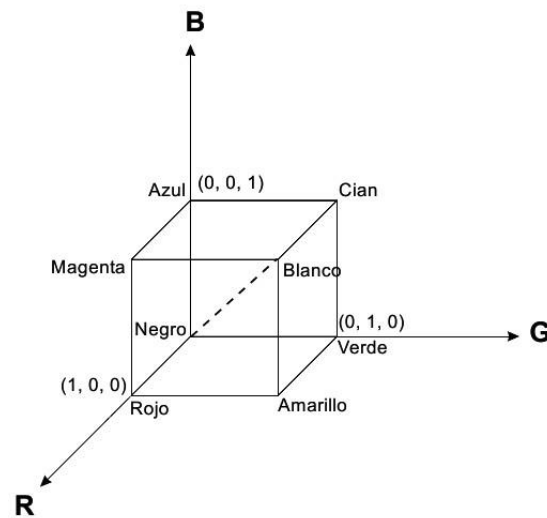
2.1.2 Corrección gamma

En el área de procesamiento de imágenes, uno se encuentra con la necesidad de mejorar una imagen con el objetivo de obtener otra más aceptable para una aplicación específica. Dependiendo del problema en sí, existen diversas técnicas para lograrlo, entre las cuales nos encontramos con las *transformaciones de potencia*. Estas transformaciones son útiles para la manipulación de contraste y tienen la forma básica

$$s = cr^\gamma,$$

donde r es el valor de entrada, s es el de salida y c y γ son constantes positivas. En la Figura 2.1 se muestran las transformaciones para diferentes valores de γ . Para $\gamma < 1$, ésta amplifica los valores oscuros de entrada y contrae los claros, mientras que lo opuesto ocurre para $\gamma > 1$.

Una variedad de dispositivos usados para la captura de imágenes responden de acuerdo a la ley de potencia. Por convención, el exponente en la ecuación es referido como *gamma*, y el proceso usado para la corrección de este fenómeno es llamado *corrección gamma*.



2.1.3 Modelos de color

El uso del color en el procesamiento de imágenes es de importancia ya que suele simplificar la identificación y extracción de objetos en una escena. Un modelo de color especifica un sistema de coordenadas y un subespacio en el cual un color es definido por un punto. Muchos modelos están orientados al hardware como lo son el RGB (rojo, verde y azul), CMY (cian, magenta y amarillo) y el CMYK (cian, magenta, amarillo y negro); otros están orientados hacia aplicaciones donde se pretende manipular el color, tales como HSI (tono, saturación e intensidad), HSV (tono, saturación y valor) y el YCbCr (luminancia, crominancia azul y crominancia roja).

A continuación se dan las bases de los modelos empleados en este trabajo.

Modelo de color RGB

Este modelo se basa sobre un sistema de coordenadas cartesianas y el subespacio en el cual se define el color es el cubo normalizado entre $[0,1]$ (Ver Figura 2.2). Los valores rojo (R), verde (G) y azul (B) están sobre los vértices $(1,0,0)$, $(0,1,0)$ y $(0,0,1)$ respectivamente; el origen corresponde al negro, en el punto más lejano a éste se ubica el blanco, y el resto de los vértices corresponden a los colores secundarios: cian, magenta y amarillo. El resto de colores está compuesto por la combinación de los colores primarios R, G y B, y son puntos dentro del cubo definido por un vector extendido desde el origen.

Las imágenes representadas en el modelo RGB consisten de tres planos, uno por cada color primario. El número de bits usados para representar cada píxel es llamado *profundidad*.

Modelo de color YCbCr

Este es un modelo de color usado en sistemas de video y fotografía digital, donde Y es la componente de luminancia y Cb y Cr son las componentes de cromaticidad azul y roja.

YCbCr es el equivalente digital al modelo YUV (Sistema de TV). Concentra la mayor parte de la información de la imagen en la luminancia y menos en la cromaticidad. El resultado es que los elementos de YCbCr están menos correlacionados y pueden ser codificados por separado.

El paso de un sistema RGB a un YCbCr es obtenido mediante las ecuaciones:

$$\begin{aligned} Y &= 0,3R + 0,6G + 0,1B \\ Cb &= \frac{B-Y}{2} + 0,5 \\ Cr &= \frac{R-Y}{2} + 0,5. \end{aligned}$$

2.1.4 Procesamiento morfológico

El procesamiento morfológico es una herramienta que nos permite extraer componentes de una imagen que son útiles para la representación y descripción de regiones. Dos operaciones de interés en esta área son la dilatación y erosión.

Dilatación

Consideremos a las imágenes binarias como conjuntos de elementos en Z^2 , donde cada elemento de un conjunto corresponde a un par de coordenadas (x,y) en la imagen. Sean A y B dos conjuntos en Z^2 , la dilatación de A por B es denotada por $A \oplus B$ y se define como

$$A \oplus B = \{z | (\hat{B})_z \cap A \neq \emptyset\},$$

donde \hat{B} es la reflexión de B alrededor de su origen. A la imagen B se la suele llamar *elemento estructurante*. Esta ecuación nos dice que la dilatación de A por B es el conjunto de todos los desplazamientos z , tal que \hat{B} y A estén solapadas al menos por un elemento. Los desplazamientos son realizados en base al origen de \hat{B} .

Una de las aplicaciones de la dilatación es para “relleno de huecos” de tamaño igual o menor que el elemento estructurante.

Erosión

Al igual que en la dilatación, para conjuntos A y B en Z^2 , la erosión de A por B se denota por $A \ominus B$ y es definida como

$$A \ominus B = \{z | (B)_z \subseteq A\}.$$

Esto nos indica que la erosión corresponde a todos los puntos z tales que B , trasladado por z , es un subconjunto de A . Esta operación es la opuesta a la dilatación y su mayor uso es para eliminar detalles irrelevantes de tamaño menor o igual que B .

Extensión a imágenes en escala de grises

Aquí, las imágenes son representadas como $f(x,y)$ y $b(x,y)$, imagen de entrada y elemento estructurante respectivamente, donde se asume que ambas son funciones discretas.

La dilatación de f por b se define como

$$f \oplus b(s,t) = \max \{f(s-x, t-y) + b(x,y) | (s-x), (t-y) \in D_f; (x,y) \in D_b\}, \quad (2.6)$$

donde D_f y D_b son el dominio de f y b , respectivamente. Las condiciones de que $(s-x)$ y $(t-y)$ deben estar en el dominio de f , y que x e y tienen que estar en el dominio de b , indican solapamiento de las dos funciones en al menos un elemento, al igual que en el caso binario.

El efecto logrado es una imagen de salida que tiende a ser más brillante que la de entrada, si todos los valores de los elementos estructurantes son positivos, y la eliminación o reducción de detalles oscuros dependiendo de los valores y forma de b .

La erosión es definida como

$$f \ominus b(s,t) = \min \{f(s+x, t+y) - b(x,y) | (s+x), (t+y) \in D_f; (x,y) \in D_b\} \quad (2.7),$$

donde D_f y D_b son el dominio de f y b , respectivamente. Las condiciones de que $(s+x)$ y $(t+y)$ tienen que estar en el dominio de f , y x e y en el dominio de b , expresan que el elemento estructurante tiene que estar completamente contenido en la imagen erosionada, es decir, f .

Como en el caso binario, la imagen de salida tiende a ser más oscura que la imagen de entrada, si todos los valores del elemento estructurante son positivos. Los detalles brillosos en la imagen de entrada son reducidos si éstos son más pequeños en área que el elemento estructurante, con un grado de reducción dependiendo de los valores de grises circundantes a los detalles brillosos y por la forma y amplitud de los valores de b .

2.2 Análisis de las componentes principales aplicado a imágenes faciales

El Análisis de las Componentes Principales (PCA, del inglés *Principal Component Analysis*) es una técnica estadística que realiza una transformación lineal y ortogonal de los datos a un nuevo sistema de coordenadas tal que la máxima varianza queda proyectada sobre la primer coordenada, la segunda mayor varianza sobre la segunda coordenada, y así sucesivamente. Transforma un conjunto de variables correlacionadas en un conjunto de variables no correlacionadas, y simplifica la transformación encontrando las componentes más cercanas a las variables originales pero ordenadas en forma decreciente al orden de su varianza. Puede ser utilizada para reducir la dimensionalidad del conjunto de datos reteniendo aquellas características del conjunto que mayor contribución hacen a su varianza, manteniendo las primeras componentes principales de más bajo orden. Tales componentes frecuentemente contienen la mayor información del conjunto de datos.

Para el caso de imágenes faciales, se pretende extraer las características más relevantes en una imagen facial, las cuales

pueden o no estar relacionadas a las características faciales como boca, nariz, ojos, etc.. En términos matemáticos, se quiere encontrar las componentes principales de la distribución de caras, o lo que es lo mismo, los eigenvectores de la matriz de covarianza del conjunto de imágenes caras [TP91]. Estos eigenvectores pueden ser vistos como el conjunto de vectores que caracterizan la variación entre las imágenes faciales. A partir de esto, cada cara puede ser representada exactamente como una combinación lineal de los eigenvectores. Además, cada rostro se puede aproximar usando aquellos que tienen asociados los eigenvalores más altos, es decir, los que responden a la mayor varianza del conjunto de imágenes faciales.

2.2.1 Cálculo de eigenfaces

Sea $I(x,y)$ una imagen facial en escala de grises de $N \times N$. Una imagen puede ser representada como un vector de dimensión N^2 , o equivalentemente, un punto en el espacio N^2 -dimensional. Por lo tanto, un conjunto de imágenes es mapeado a una colección de puntos en un espacio inmenso, considerando el tamaño de las imágenes con las que se suele trabajar. De esta manera, un conjunto de imágenes faciales, dado que todas tienen un aspecto similar, se mapean en un subespacio del espacio de imágenes y así pueden ser descritas por un subespacio de menor dimensión. La idea es encontrar los vectores que mejor describan la distribución de imágenes cara dentro del espacio completo de imágenes. Estos vectores de longitud N^2 , llamados *eigenfaces*, definen el subespacio de imágenes rostro, el cual se denomina *espacio de caras*.

Sea $\Gamma_1, \Gamma_2, \Gamma_3, \dots, \Gamma_M$ un conjunto de imágenes faciales. La cara media del conjunto es definida por

$$\Psi = \frac{1}{M} \sum_{n=1}^M \Gamma_n.$$

Cada cara difiere de ésta por el vector $\Phi_i = \Gamma_i - \Psi$. A partir de este conjunto de vectores, se busca un conjunto de M vectores ortonormales, \mathbf{u}_n , los cuales mejor describan los datos. El k -ésimo vector, \mathbf{u}_k , es seleccionado tal que

$$\lambda_k = \frac{1}{M} \sum_{n=1}^M (\mathbf{u}_k^T \Phi_n)^2$$

es un máximo, sujeto a que los vectores \mathbf{u}_k sean ortonormales. Los vectores \mathbf{u}_k y los escalares λ_k son los eigenvectores y eigenvalores, respectivamente, de la matriz de covarianza

$$\begin{aligned} C &= \frac{1}{M} \sum_{n=1}^M \Phi_n \Phi_n^T \\ &= AA^T, \end{aligned}$$

donde $A = [\Phi_1 \Phi_2 \Phi_3 \dots \Phi_M]$. La matriz C resultante es de $N^2 \times N^2$ y obtener N^2 eigenvectores y eigenvalores es un tarea costosa para el tamaño de imágenes comúnmente procesadas.

Una consideración a tener en cuenta es que si el número de puntos en el espacio de imágenes es menor que la dimensión del mismo ($M < N^2$), habrá solamente $M - 1$ eigenvectores significantes, el resto tendrían asociados eigenvalores iguales a cero. En este caso, se tienen M imágenes caras de dimensión N^2 y se pueden encontrar los eigenvectores N^2 -dimensionales resolviendo los eigenvectores para una matriz de $M \times M$ y luego realizar una combinación lineal de las imágenes faciales Φ_l .

Consideremos los eigenvectores \mathbf{v}_l de $A^T A$ tal que

$$A^T A \mathbf{v}_l = \mu_l \mathbf{v}_l.$$

Premultiplicando ambos lados por A , se tiene

$$AA^T A \mathbf{v}_l = \mu_l A \mathbf{v}_l.$$

De aquí se ve que $A \mathbf{v}_l$ son los eigenvectores de la matriz de covarianza C . Entonces, se construye la matriz de $M \times M$, $L = A^T A$, donde $L_{mn} = \Phi_m^T \Phi_n$, y se encuentran los M eigenvectores \mathbf{v}_l , de L . Estos

vectores determinan la combinaci3n lineal de las M im3genes faciales para obtener los eigenvectores

$$\mathbf{u}_l = A\mathbf{v}_l$$

$$= \sum_{k=1}^M \mathbf{v}_{lk} \Phi_k.$$

De esta forma se reducen enormemente los c3lculos del orden del n3mero de p3xeles en las im3genes (N^2) al orden del n3mero de im3genes en el conjunto (M). Los eigenvalores asociados permiten ordenar los eigenvectores de acuerdo al aporte que hacen a la variaci3n entre im3genes.

Proyecci3n y reconstrucci3n de una imagen facial

Una imagen facial Γ es proyectada sobre el espacio de caras por medio de la operaci3n

$$x_k = \mathbf{u}_k^T (\Gamma - \Psi), \quad \text{para } k = 1, \dots, M. \quad (2.14)$$

Los pesos x_k forman un vector $\mathbf{x}^T = [x_1 \ x_2 \ x_3 \ \dots \ x_M]$ que describen el aporte de cada eigenface en representar la imagen cara, tratando a las eigenfaces como un conjunto base para las im3genes faciales.

La reconstrucci3n exacta de una cara se realiza a trav3s de una combinaci3n lineal de las eigenfaces ponderadas con los respectivos pesos x_k

M

$$\Gamma = \sum_{k=1}^M x_k \mathbf{u}_k.$$

$k=1$

2.3 Reconocimiento de patrones

Un *patr3n* es un *arreglo de caracter3sticas* que describen una sen3al, y es representado por un vector multidimensional. Una *clase* es un conjunto de patrones que comparten una propiedad en com3n. Las clases son denotadas por $\omega_1, \omega_2, \dots, \omega_W$, donde W es la cantidad de clases. Las t3cnicas de reconocimiento de patrones son

aquellas que asignan los patrones a sus respectivas clases teniendo en cuenta las características de éstos.

El reconocimiento de patrones se divide en dos áreas: *reconocimiento basado en la teoría de la decisión* y *reconocimiento estructural*. El primero trata con patrones representados por características cuantitativas, y el segundo con las relaciones estructurales de las características de los patrones.

Reconocimiento basado en la teoría de la decisión

Estos métodos se basan en *funciones de decisión*. Siendo $\mathbf{x} = (x_1, x_2, \dots, x_n)^T$ un vector patrón y existiendo W clases $\omega_1, \omega_2, \dots, \omega_W$, la idea es encontrar W funciones $d_1(\mathbf{x}), d_2(\mathbf{x}), \dots, d_W(\mathbf{x})$ tal que si un patrón \mathbf{x} pertenece a la clase ω_i , entonces

$$d_i(\mathbf{x}) > d_j(\mathbf{x}), \quad j = 1, 2, \dots, W; \quad j \neq i$$

Es decir, un patrón \mathbf{x} desconocido es evaluado en las W funciones de decisión, si el valor más alto es obtenido en la función d_i , entonces se dice que \mathbf{x} pertenece a la clase ω_i .

2.3.1 Redes neuronales

Uno de los métodos basados en la teoría de la decisión es la red neuronal. *Una red neuronal es un procesador masivamente distribuido que tiene una propensión natural para almacenar el conocimiento experimental y ponerlo a disposición para su uso* [Hay99].

La idea de una red neuronal surge con el objetivo de modelar el comportamiento del cerebro humano al realizar una tarea, como es el reconocimiento de las personas, empleando una interconexión masiva de pequeñas unidades de procesamiento las cuales se denominan *neuronas* o *nodos*. Las redes neuronales se asemejan al cerebro en dos aspectos: 1) el conocimiento es adquirido

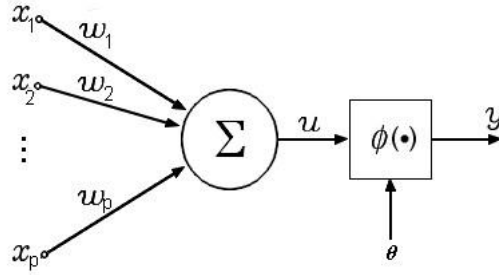


Figura 2.3: Arquitectura de un perceptrón simple.

por un proceso de aprendizaje y 2) los pesos de conexión entre neuronas, conocidos como *pesos sinápticos*, se usan para almacenar el conocimiento. El algoritmo de aprendizaje modifica los pesos sinápticos hasta lograr el conocimiento requerido.

La manera en la cual las neuronas de una red están estructuradas está ligada al algoritmo de aprendizaje usado para el entrenamiento de la red. Entre las arquitecturas posibles para modelos neuronales están las denominadas *redes hacia adelante*, como lo son el perceptrón simple, perceptrón multicapa y redes de función de base radial, las *redes recurrentes* y otros *modelos híbridos*.

A continuación se explicará el funcionamiento de un perceptrón simple, la arquitectura utilizada en este trabajo, el perceptrón multicapa, y se hará una introducción a las redes de función de base radial.

Perceptrón simple

El perceptrón es la forma más simple de red neuronal, consistiendo de una única neurona. Su arquitectura se ilustra en la Figura 2.3. Está constituida por un combinador lineal que suma las señales de entrada x_j ponderadas por los pesos sinápticos w_j correspondientes. La señal de entrada es denotada por

$$\mathbf{x} = (x_1, x_2, \dots, x_p)^T$$

y el conjunto de pesos sinápticos está dado por el vector

$$\mathbf{w} = (w_1, w_2, \dots, w_p)^T.$$

La salida del combinador lineal, representada por u , es la entrada a una *función de activación* que limita la amplitud de la salida de la neurona. Típicamente, el rango de amplitud normalizado de la salida es entre $[0, 1]$ o $[-1, 1]$.

La arquitectura también incluye un umbral θ que tiene como objetivo fijar un valor mínimo que debe superar la salida del combinador lineal para activar la neurona. Este umbral puede ser agregado al vector \mathbf{w} como $w_0 = \theta$, considerando una entrada adicional $x_0 = 1$.

En términos matemáticos, una neurona está representada por medio de las ecuaciones

p

$$u = \sum_{j=0}^p w_j x_j$$

$j=0$

y

$$y = \varphi(u),$$

donde φ representa la función de activación e y es la salida de la neurona.

Existen tres tipos básicos de funciones de activación para determinar la salida:

a) *Función umbral*. Aplica la siguiente ecuación

$$y = \begin{cases} 1 & \text{si } u \geq 0 \\ 0 & \text{si } u < 0 \end{cases}$$

b) *Función bipolar*. Se la define como

$$y = \begin{cases} 1 & \text{si } u \geq 0 \\ -1 & \text{si } u < 0 \end{cases} \quad (2.20) \quad \text{c) Función}$$

sigmoidal. Es la más comúnmente usada en el modelado de redes neuronales, también llamada *función logística*. Se define como una función estrictamente incremental de la siguiente manera

$$\phi(u) = \frac{1}{1 + e^{-u}}$$

En la Figura 2.4 se muestran las tres funciones descritas. Cuando la salida de la función de activación es cercana o igual a su valor máximo, el perceptrón es activado; mientras que cuando está cercano a su valor mínimo, se dice que está desactivado. Por lo tanto, se puede ver al perceptrón simple como un clasificador de dos clases.

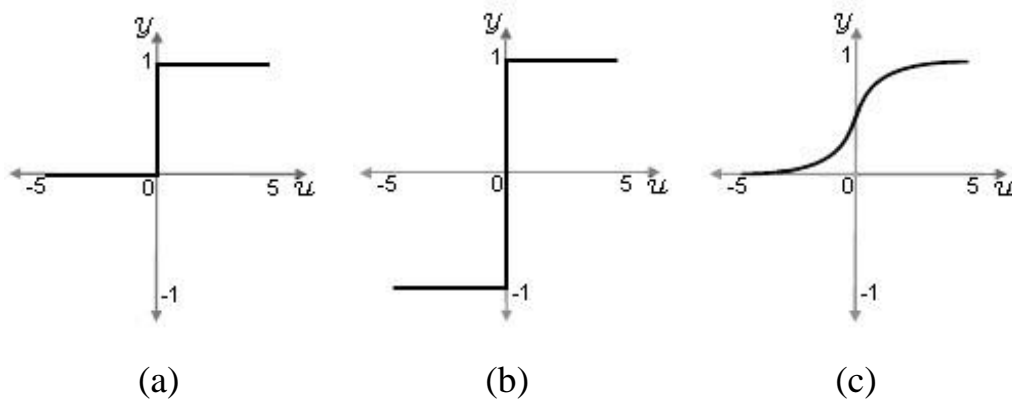


Figura 2.4: Funciones de activación. a) Umbral, b) Bipolar, c) Sigmoidal.

Como se mencionó anteriormente, los pesos sinápticos se utilizan para almacenar el conocimiento de la neurona, siendo fijos o adaptados mediante el aprendizaje. Este aprendizaje implica que los parámetros de la red sean adaptados a través de un proceso continuo de estimulación, por el ambiente en el cual se encuentra. La estimulación es llevada a cabo por medio de los patrones de entrenamiento, y cuando algún criterio de finalización es

satisfecho, el algoritmo es detenido y se dice que la red est´a entrenada y lista para su funcionamiento.

Hay una variedad de algoritmos y reglas de aprendizaje para los diferentes dise˜nos de redes. Para el caso del perceptr´on, una de ellas es la denominada *regla delta*. Esta se basa en ajustar los pesos con el objetivo de minimizar la diferencia entre la respuesta deseada y la obtenida por el perceptr´on. La respuesta deseada es denotada por $d(n)$ y la salida producida por un est´ımulo $\mathbf{x}(n)$, por $y(n)$. Entonces, el error medido en el paso n del entrenamiento est´a dado por

$$e(n) = d(n) - y(n).$$

El ajuste de los pesos se realiza minimizando el criterio de error cuadr´atico instant´aneo

$$E(\mathbf{w}) = \frac{1}{2}e^2(n),$$

cuyo gradiente es

$$\nabla E(\mathbf{w}) = -e(n)\phi^0(n)\mathbf{x}(n).$$

De esta manera, el valor actualizado est´a dado por

$$w_j(n+1) = w_j(n) + \eta e(n)\phi^0(n)x_j(n),$$

donde η es una constante positiva que determina la tasa de aprendizaje. Este par´ametro tiene un impacto importante en el desempe˜no del algoritmo. Si es peque˜no, el proceso de aprendizaje avanza lentamente pero la convergencia a una soluci´on estable puede tardar demasiado tiempo; en cambio si η es grande el aprendizaje es m´as acelerado pero hay peligro de que el proceso de aprendizaje diverja y se obtenga una soluci´on inestable. El perceptr´on es capaz de lograr un ajuste de los pesos con error cero solamente cuando los patrones son *separables linealmente*, es decir, los patrones pertenecen a dos clases diferentes que pueden ser separados geom´etricamente por un hiperplano.

Perceptrón multicapa

El perceptrón multicapa (MLP, del inglés *Multi-Layer Perceptron*) representa una generalización del perceptrón simple. Se constituye por una capa de entrada, una o varias capas ocultas y una capa de salida. Un esquema de una red con una capa oculta se muestra en la Figura 2.5. La capa de entrada simplemente copia el patrón de entrada, las capas ocultas y de salida la conforman un conjunto de neuronas. La señal de entrada se propaga por la red de capa en capa, de izquierda a derecha. Además, se puede observar que cada neurona en alguna capa está conectada a todas las neuronas de la capa previa y la siguiente.

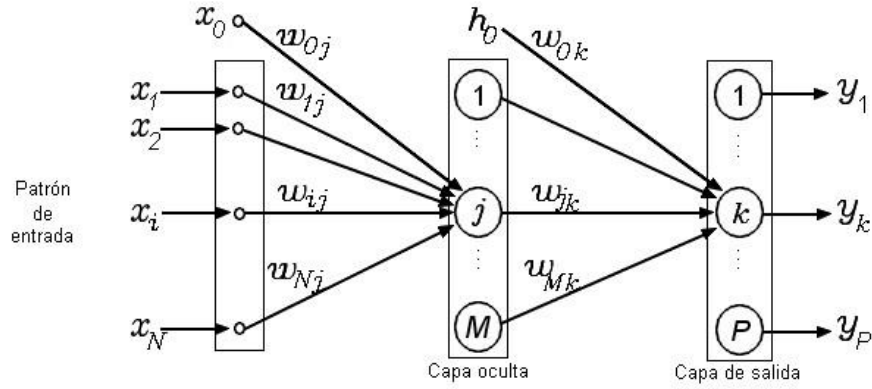
El MLP es entrenado en una manera supervisada con algoritmo de *retropropagación del error*, el cual está basado en regla delta, mencionado anteriormente para el caso del perceptrón simple [Hay99]. Este algoritmo comienza con una fase de propagación hacia adelante, donde un patrón de entrenamiento es aplicado en la capa de entrada obteniendo la salida de cada neurona que constituye la red. Específicamente, la salida de la neurona j es calculada de la misma manera que para el caso del perceptrón simple,

$$y_j(n) = \varphi_j(u_j(n)), \quad (2.26)$$

donde φ_j es la función de activación y $u_j(n)$ la salida del combinador lineal definida por

$$u_j(n) = \sum_{i=0}^p w_{ij}(n) y_i(n), \quad (2.27)$$

donde p es el número de entradas aplicado a la neurona j , $w_{ij}(n)$ es el peso sináptico que conecta la neurona i a la neurona j , y $y_i(n)$ es la entrada de la neurona j , o equivalentemente, la salida de la neurona i .



Siguiendo con el esquema de la Figura 2.5, consideremos los nodos de la capa de salida. El error medido en la neurona k en el paso n es definido por

$$e_k(n) = d_k(n) - y_k(n).$$

Entonces, partiendo del error cuadrático instantáneo en la neurona k definido como $\frac{1}{2}e_k^2(n)$, la suma instantánea de errores cuadráticos es obtenida por sumar $\frac{1}{2}e_k^2(n)$ sobre todas las neuronas, expresado matemáticamente como

$$E(n) = \frac{1}{2} \sum_{k=1}^P e_k^2(n),$$

donde P es la cantidad de neuronas en la capa de salida. La suma instantánea de errores cuadráticos $E(n)$ es una función de todos los pesos sinápticos de la red, por lo tanto, el objetivo del proceso de aprendizaje es ajustar los pesos minimizando $E(n)$.

La regla delta establece que la corrección aplicada a los pesos sinápticos en los nodos de la capa de salida es

$$\begin{aligned} \Delta\omega_{jk}(n) &= -\eta \frac{\partial E(n)}{\partial \omega_{jk}(n)} \\ &= \eta e_k(n) \phi'_k(u_k(n)) y_j(n) \end{aligned}$$

siendo η la tasa de aprendizaje. Esto es posible dado que se conocen las respuestas deseadas de cada neurona, mientras que no ocurre lo mismo para los nodos de las capas ocultas. De todos modos, es posible inferir una regla de aprendizaje basada en modificar estos pesos para tratar de disminuir el error en la capa de salida, mediante la propagación hacia atrás de las cantidades e_k . Así, la actualización de los pesos en las capas ocultas se realiza mediante el gradiente descendiente de la función criterio error cuadrático, pero calculado con respecto a los pesos w_{ij} como

$$\Delta\omega_{ij}(n) = -\eta \frac{\partial E(n)}{\partial \omega_{ij}(n)}.$$

A través de la regla de la cadena, se puede calcular la derivada parcial de (2.31) según

$$\frac{\partial E(n)}{\partial \omega_{ij}(n)} = \frac{\partial E(n)}{\partial y_j(n)} \frac{\partial y_j(n)}{\partial u_j(n)} \frac{\partial u_j(n)}{\partial \omega_{ij}(n)},$$

donde las últimas dos derivadas parciales son calculadas como

$$\frac{\partial y_j(n)}{\partial u_j(n)} = \phi'_j(u_j(n))$$

y

$$\frac{\partial u_j(n)}{\partial \omega_{ij}(n)} = y_i;$$

y por último, se calcula la derivada parcial $\partial E(n)/\partial y_j(n)$ como

$$\frac{\partial E(n)}{\partial y_j(n)} = \sum_{k=1}^P e_k(n) \frac{\partial e_k(n)}{\partial y_j(n)}.$$

Usando la regla de la cadena para la derivada parcial $\partial e_k(n)/\partial y_j(n)$ se reescribe la ecuación y se obtiene

$$\begin{aligned} \frac{\partial E(n)}{\partial y_j(n)} &= \sum_{k=1}^P e_k(n) \frac{\partial e_k(n)}{\partial u_k(n)} \frac{\partial u_k(n)}{\partial y_j(n)} \\ &= - \sum_{k=1}^P e_k(n) \phi'_k(u_k(n)) \omega_{jk}(n). \end{aligned}$$

Reemplazando en (2.31) se obtiene la regla de actualización de pesos en las capas ocultas

$$\Delta\omega_{ij}(n) = \eta \left[\sum_{k=1}^P e_k(n) \phi'_k(u_k(n)) \omega_{jk}(n) \right] \phi'_j(u_j(n)) y_i.$$

La actualización de los pesos es realizada por cada patrón de entrenamiento que se presente a la red. La convergencia del algoritmo de aprendizaje, a una solución adecuada, se comprueba al finalizar la presentación del conjunto completo de entrenamiento, verificando el valor de alguna función error prefijada o mediante alguna técnica de estimación del error.

Con respecto a la tasa de aprendizaje η , aquí se presenta el mismo caso que para el perceptrón simple. Un método para incrementar la tasa de aprendizaje y evitar una solución inestable es modificar la regla delta incluyendo un término denominado *momento*

$$\Delta\omega_{ij}(n) = \alpha \Delta\omega_{ij}(n-1) - \eta \frac{\partial E(n)}{\partial \omega_{ij}(n)},$$

donde α es frecuentemente un número positivo llamado *constante de momento*. Cabe aclarar que los subíndices en la ecuación hacen referencia tanto a pesos de neuronas ocultas como de salidas.

Red de función de base radial

Una red de función de base radial (RBF, del inglés *Radial Basis Function*) está constituida por tres capas cuando es usada para clasificación, una capa de entrada, una oculta y una de salida [HH93]. Un esquema de una red RBF se observa en la Figura 2.6. Al igual que en una red MLP, la capa de entrada copia el patrón de entrada. La capa oculta y de salida la conforman un conjunto de nodos. Las salidas de los nodos en la capa de salida son una combinación lineal ponderada de las salidas de los nodos de la capa oculta.

En términos matemáticos, la salida del nodo k en la capa de salida está dada por

M

$$y_k = \sum_{j=0}^M w_{jk} \phi_j,$$

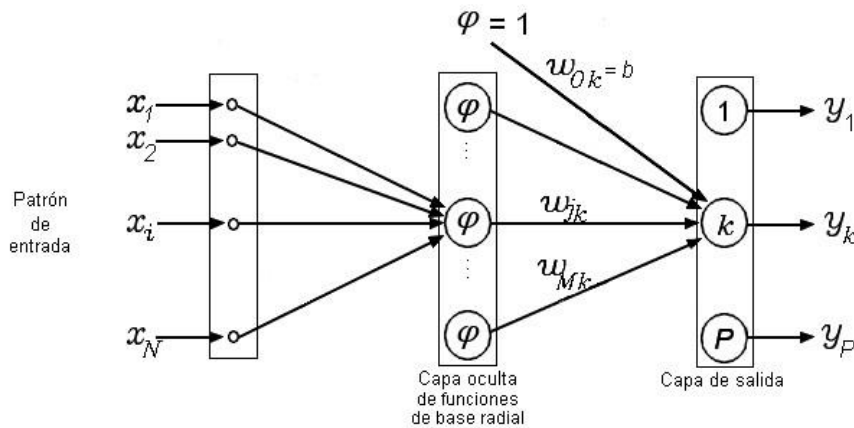
$j=0$

donde w_{jk} es el peso sináptico que conecta el nodo j de la capa oculta al nodo k de la capa de salida, y ϕ_j para $j = 1, 2, \dots, M$ es un conjunto de funciones conocidas como *funciones de base radial*. Estas funciones son las salidas de los nodos de la capa oculta, y comúnmente tienen una forma Gaussiana dada por

$$\phi_j = e^{\frac{-u_j^2}{2\sigma_j^2}},$$

donde

$$u_j = \sum_{i=1}^N (x_i - w_{ij})^2$$



Aquí, ϕ_j es la salida del nodo j en la capa oculta, x_i es la entrada a la red y w_{ij} corresponde a los pesos sinápticos del nodo j . Además, se deben notar que los pesos w_{ij} son el centro de la función Gaussiana, y que σ_j^2 define el ancho de la misma. La salida del nodo j está en el rango de 0 a 1, de modo que cuanto más cerca está la entrada de

la red al centro de la función Gaussiana, mayor será la respuesta del nodo j . El nombre *función de base radial* se debe a que ϕ_j es radialmente simétrica.

Existen varias aproximaciones para el entrenamiento de una red RBF. Algunas de éstas, suelen dividir el aprendizaje en dos etapas. En la primer etapa, el aprendizaje se realiza en la capa oculta mediante un método no supervisado, como por ejemplo, un algoritmo de clustering; en la segunda etapa, el aprendizaje se realiza en la capa de salida usando un método supervisado. Una vez obtenida una solución inicial, se puede aplicar un algoritmo de entrenamiento supervisado en ambas capas, para optimizar los parámetros de la red.

Hay una variedad de algoritmos de clustering que se pueden aplicar en la capa oculta. Uno de los más conocidos por su simplicidad y producir buenos resultados, es el algoritmo *K-means*. Por otro lado, para el caso de la capa de salida, se suele usar la regla delta; mientras que para la optimización de parámetros se puede aplicar un método de gradiente.

Desarrollo del software

En este proyecto se desarrolló un sistema que permite identificar a una persona por medio del reconocimiento de una imagen de su rostro. Para ello, se instaló una cámara web en una PC de escritorio para la captura de imágenes. Por medio de ésta, el sistema toma una foto de la persona, la procesa e informa la identidad de la misma.

A la hora de modelar sistemas, una de las herramientas más utilizadas es el Lenguaje Unificado de Modelado (UML, del inglés *Unified Modeling Language*) [BRJ98]. Este estándar es un lenguaje gráfico para visualizar, especificar, construir y documentar un sistema. Provee un vocabulario y reglas que permiten realizar modelos conceptuales de varios aspectos del sistema, con el objetivo de facilitar la comunicación entre el equipo de desarrollo y los usuarios finales, quienes son los que especifican las funciones y el comportamiento deseado del mismo.

A fin de realizar un análisis y documentación del sistema, en las siguientes secciones se explican las funcionalidades necesarias por medio de UML. Se especifican los requerimientos identificados, y se realizan los diagramas de casos de uso, clases y secuencia respectivamente.

3.1 Análisis de requerimientos

Un requerimiento puede ser definido como:

- Una condición o capacidad que un usuario necesita para resolver un problema o lograr un objetivo.
- Una condición o capacidad que debe tener un sistema para satisfacer un contrato, norma, especificación, u otros documentos formales.
- Una representación documentada de una condición o capacidad como las mencionadas anteriormente.

A su vez, estos requerimientos pueden ser diferenciados en *requerimientos funcionales* y *requerimientos no funcionales*. Los

primeros, son declaraciones de las funciones que el sistema debe ser capaz de realizar, de como debe responder ante entradas particulares; mientras que los segundos, son restricciones del sistema, tales como disponibilidad, mantenimiento, seguridad, capacidad de los dispositivos de entrada/salida, rendimiento (como velocidad y tiempo de respuesta), etc..

En base a estas definiciones, se realiza un análisis de requerimientos con el objetivo de identificar tanto las funcionalidades que se esperan del software como sus limitaciones.

3.1.1 Requerimientos funcionales

El sistema debe ser capaz de:

1. Capturar una imagen a través de una cámara web, con la posibilidad de guardarla.
2. Identificar a una persona por medio de una imagen de su rostro.

3.1.2 Requerimientos no funcionales

Como requerimientos no funcionales se especifican los siguientes:

1. El sistema debe funcionar con una cámara web estándar y luz artificial.
2. El sistema debe ejecutarse en tiempo real, lo que implica que las operaciones computacionales deben procesarse en un tiempo aceptable.

3.2 Diagrama de casos de uso

Un caso de uso captura el comportamiento deseado de un sistema en desarrollo, sin tener que detallar como se implementa éste. Es una

descripci3n de un conjunto de secuencias de acciones que el sistema lleva a cabo para producir un resultado de inter3s para un usuario, reflejando como deber3an interactuar ambas partes.

Un diagrama de casos de uso muestra la relaci3n entre los usuarios del sistema y los casos de uso, modelando los aspectos din3micos del mismo.

Para el caso analizado aqu3, se pretende visualizar como interactu3a el usuario con el sistema haciendo uso de las funcionalidades ofrecidas por el mismo. En la Figura 3.1 se puede observar un diagrama de alto nivel del sistema, y en la Figura 3.2 un diagrama m3as detallado.

A continuaci3n se describen los casos de uso correspondientes al diagrama de la Figura 3.1.

Caso de uso: Capturar imagen
Actor: Usuario
Descripci3n: Captura una imagen de la persona con la c3mara web.
Curso normal:
1) Comienza cuando el usuario presiona el bot3n 'Capturar' o 'Identificar'.
2) Se toma una foto de la persona.
Cursos alternativos:
2.A) La c3mara no est3a inicializada.
2.A.1) El sistema inicializa la c3mara web prevvisualizando lo que 3sta captura.

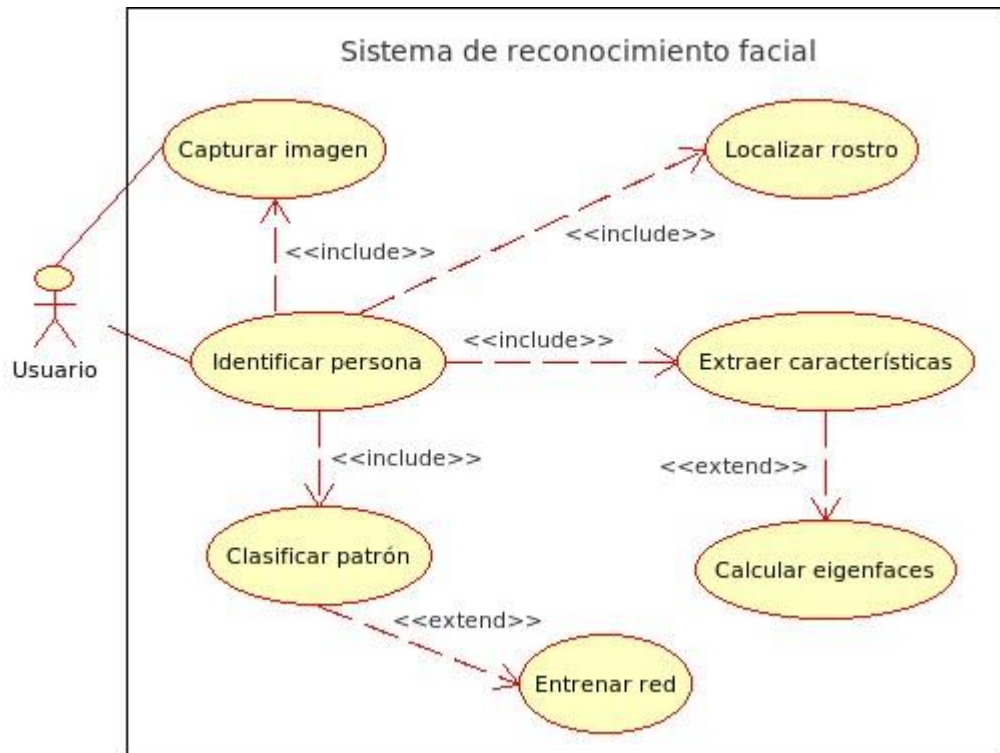


Figura 3.1: Diagrama de casos de uso de alto nivel.

Caso de uso: Identificar persona
Actor: Usuario
Descripción: Identifica a una persona a partir de una imagen de su rostro.
Curso normal:
1) Comienza cuando el usuario presiona el botón 'Identificar'.
2) El sistema captura una imagen del usuario.
3) Se localiza el rostro en la imagen.
4) Se proyecta la imagen facial mediante las eigenfaces y se obtiene un patrón representativo de ésta.
5) Se clasifica el patrón mediante una red neuronal.
6) El sistema informa la identidad de la persona.

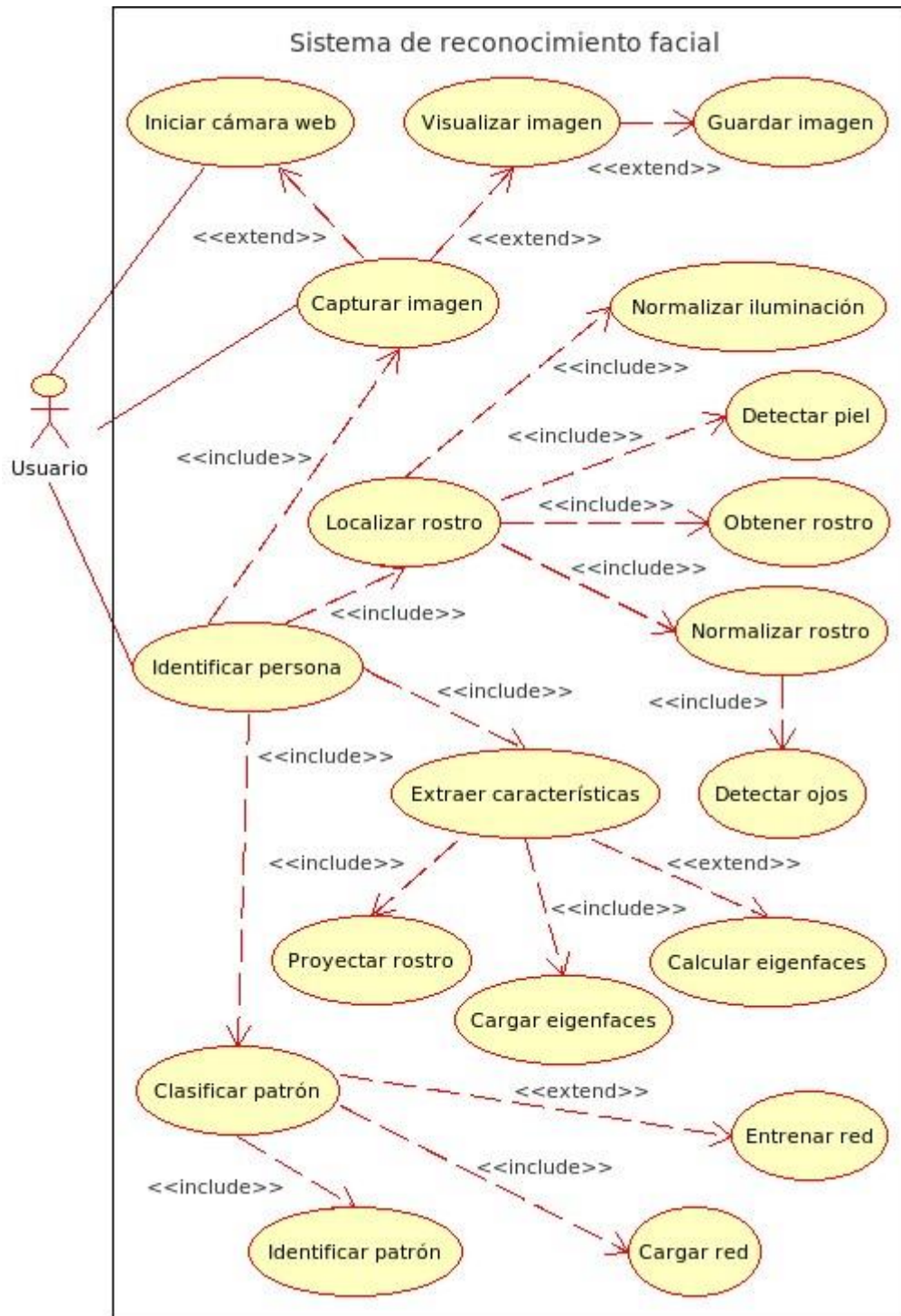


Figura 3.2: Diagrama de casos de uso detallado.

Caso de uso: Localizar rostro

Actor: -
Descripción: Obtiene una subimagen facial normalizada.
Curso normal:
1) Se normaliza la iluminación en la imagen. 2) Se detectan las regiones de piel. 3) Se aplica una plantilla elíptica para ubicar la cara y se recorta una subimagen de tamaño definido. 4) Se detectan los ojos a partir de la imagen obtenida. 5) Se normaliza el rostro a partir de los ojos y se recorta la imagen original nuevamente.
Caso de uso: Extraer características
Actor: -
Descripción: Obtiene un vector representativo de la imagen facial para clasificarlo.
Curso normal:
1) Se cargan las eigenfaces y la imagen media. 2) Se centra la imagen facial sustrayéndole la imagen media. 3) Se proyecta la imagen centrada multiplicándola por cada una de las eigenfaces, generando un punto en el espacio de caras.
Cursos alternativos:
1.A) Las eigenfaces y la imagen media no están calculadas. 1.A.1) Se calculan las eigenfaces y la imagen media. 1.A.2) Se guardan las eigenfaces y la imagen media.
Caso de uso: Calcular eigenfaces
Actor: -
Descripción: Calcula las eigenfaces de todas las imágenes en la base de datos.
Curso normal:
1) Se calcula la imagen media de todas las imágenes en la base de rostros. 2) Se calculan las eigenfaces. 3) Se guardan las eigenfaces y la imagen media.

Caso de uso: Clasificar patrón
Actor: -
Descripción: Clasifica un patrón para identificar a la persona.
Curso normal:
1) Se carga la red. 2) Se aplica el patrón a la red devolviendo como resultado la identificación de la persona.
Cursos alternativos
1.A) La red no está entrenada. 1.A.1) Se entrena la red. 1.A.2) Se guarda la red.
Caso de uso: Entrenar red
Actor: -
Descripción: Entrena la red neuronal para realizar la clasificación de patrones.
Curso normal:
1) Se proyectan todas las imágenes de la base de rostro al espacio de caras obteniendo los patrones correspondientes. 2) Se entrena la red mediante los patrones obtenidos. 3) Se guarda la red entrenada.

3.3 Diagrama de clases

Los diagramas de clases son los más usados en el modelado de sistemas orientados a objetos, los cuales muestran un conjunto de clases y sus relaciones. Se utilizan para modelar la vista de diseño estática de un sistema, lo que implica modelar el vocabulario del mismo, modelar colaboraciones o modelar esquemas.

Así como son importantes para visualizar, especificar y documentar modelos estructurales, también lo son para construir sistemas ejecutables, aplicando ingeniería directa e inversa.

En la Figura 3.3 se puede observar el diagrama correspondiente al sistema de reconocimiento facial. Como se puede notar, se definieron distintas clases que permiten realizar tanto la captura de imágenes como el reconocimiento de una persona.

Las clases *Application*, *MainFrame*, *WebCam* y *SnapshotDialog* son extensiones de clases contenidas en la librería *wxWidgets*¹, las cuales no son mostradas en el diagrama a modo de simplificar el mismo. Esta librería permite crear interfaces de usuario haciendo uso de las bibliotecas ya existentes en el sistema.

La clase *Application* (extensión de *wxApp*) representa la aplicación misma. Es usada para definir y obtener propiedades de ésta, como también para iniciarla.

La clase *MainFrame* (extensión de *wxFrame*) corresponde a la ventana de la aplicación. Gestiona las funcionalidades del menú y los botones incluidos en ella. Inicializa la cámara web, captura una imagen o identifica un usuario, según la acción ejecutada por el usuario. Además, visualiza lo que la cámara web captura.

La clase *WebCam* (extensión de *wxWindow*) se encarga de manejar el dispositivo de captura de imágenes. Detecta si se encuentra conectada una cámara a la PC, inicializa y detiene su controlador, como también captura imágenes a ser visualizadas en la vista previa por la clase *MainFrame*.

La clase *SnapshotDialog* (extensión de *wxDialog*) se utiliza para visualizar la imagen capturada por el usuario en una ventana de diálogo. Permite guardar la misma en el disco rígido o descartarla.

Para el procesamiento básico de las imágenes, se utiliza la clase *CImg*², siendo eficiente y simple de usar. Es libre, de código abierto y es distribuida bajo las licencias CeCILL-C (similar a la LGPL) o CeCILL (similar a la GPL). Permite cargar y guardar imágenes en

¹ <http://www.wxwidgets.org/>

² <http://cimg.sourceforge.net>

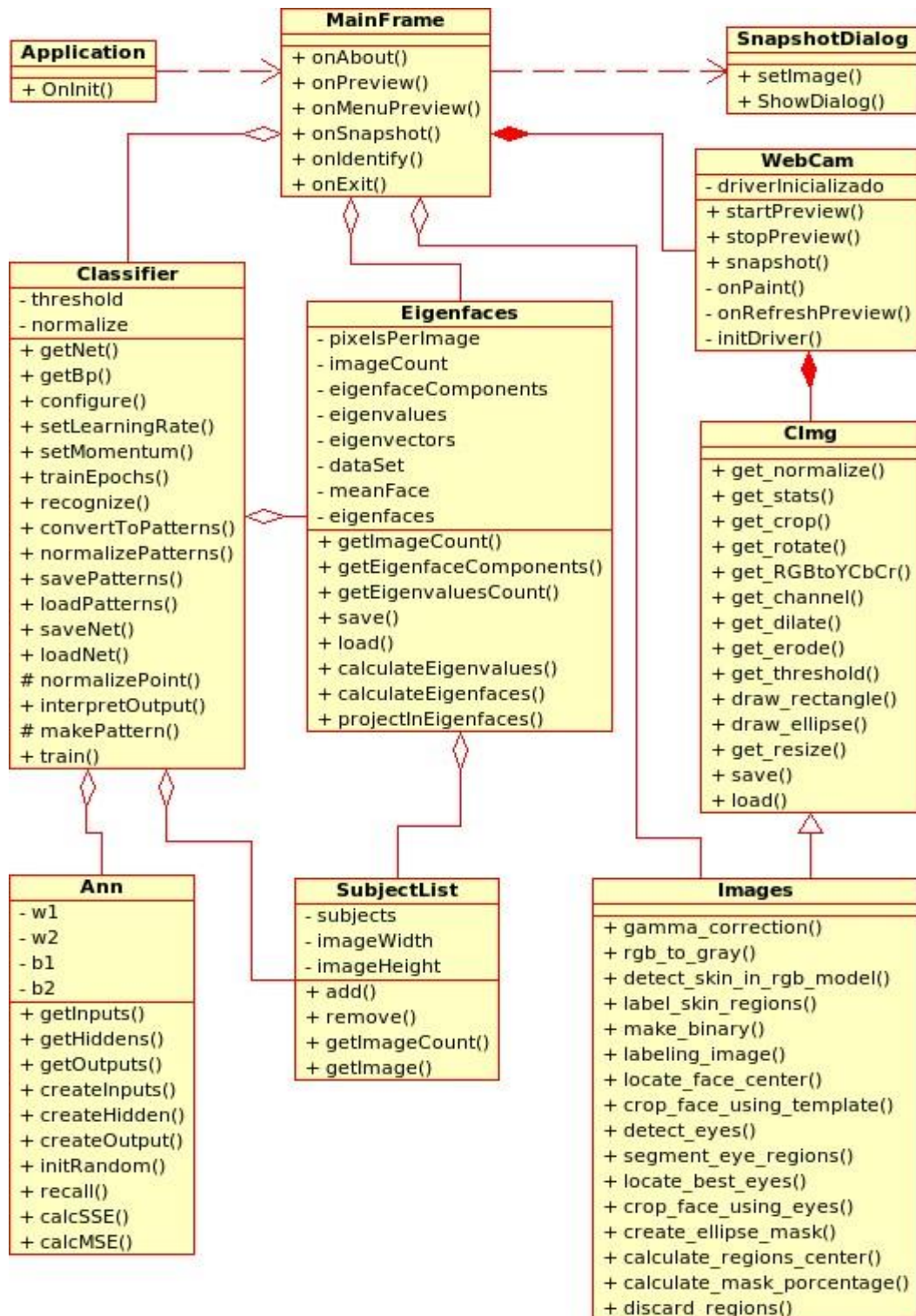
formatos ppm, jpg y otros, manipular imágenes como matrices, aplicar operaciones de dilatación y erosión, rotar y cortar imágenes, efectuar transformaciones a diferentes espacios de color, etc..

La clase *Images* hereda todas las operaciones de *CImg*, y además incluye funciones para normalizar la iluminación, segmentar la piel, localizar el rostro en una imagen y normalizar la geometría del mismo.

La clase *SubjectsList* se encarga de cargar las imágenes de la base de rostros.

La clase *Eigenfaces* contiene las funciones necesarias para generar el espacio de caras. Calcula, guarda y carga la imagen media y las eigenfaces del

Figura 3.3: Diagrama de clases.



3.4. Diagrama de secuencia

conjunto completo de rostros, y proyecta una imagen facial a este espacio.

La clase *Ann* implementa un perceptrón multicapa de tres capas: entrada, oculta y salida. Construye la red según el número de entradas y la cantidad de nodos en la capa oculta y salida e identifica un patrón aplicado a la entrada de la misma.

Por último, la clase *Classifier* se encarga de generar los patrones de entrenamiento haciendo uso de la clase *Eigenfaces*, proyectando cada imagen facial al espacio de caras. Con estos patrones, entrena la red neuronal mediante el algoritmo de retropropagación del error, y permite guardar y cargar tanto la red entrenada como los patrones mismos.

3.4 Diagrama de secuencia

Los diagramas de secuencia son usados para modelar los aspectos dinámicos del sistema. Muestran las interacciones entre los distintos objetos participantes y los mensajes que intercambian. Estos diagramas se forman colocando los objetos y actores que interactúan en el eje horizontal, y en el eje vertical se colocan los mensajes intercambiados en orden temporal, desde arriba hacia abajo. El objeto que inicia la interacción suele ubicarse a la izquierda y los objetos subordinados a la derecha.

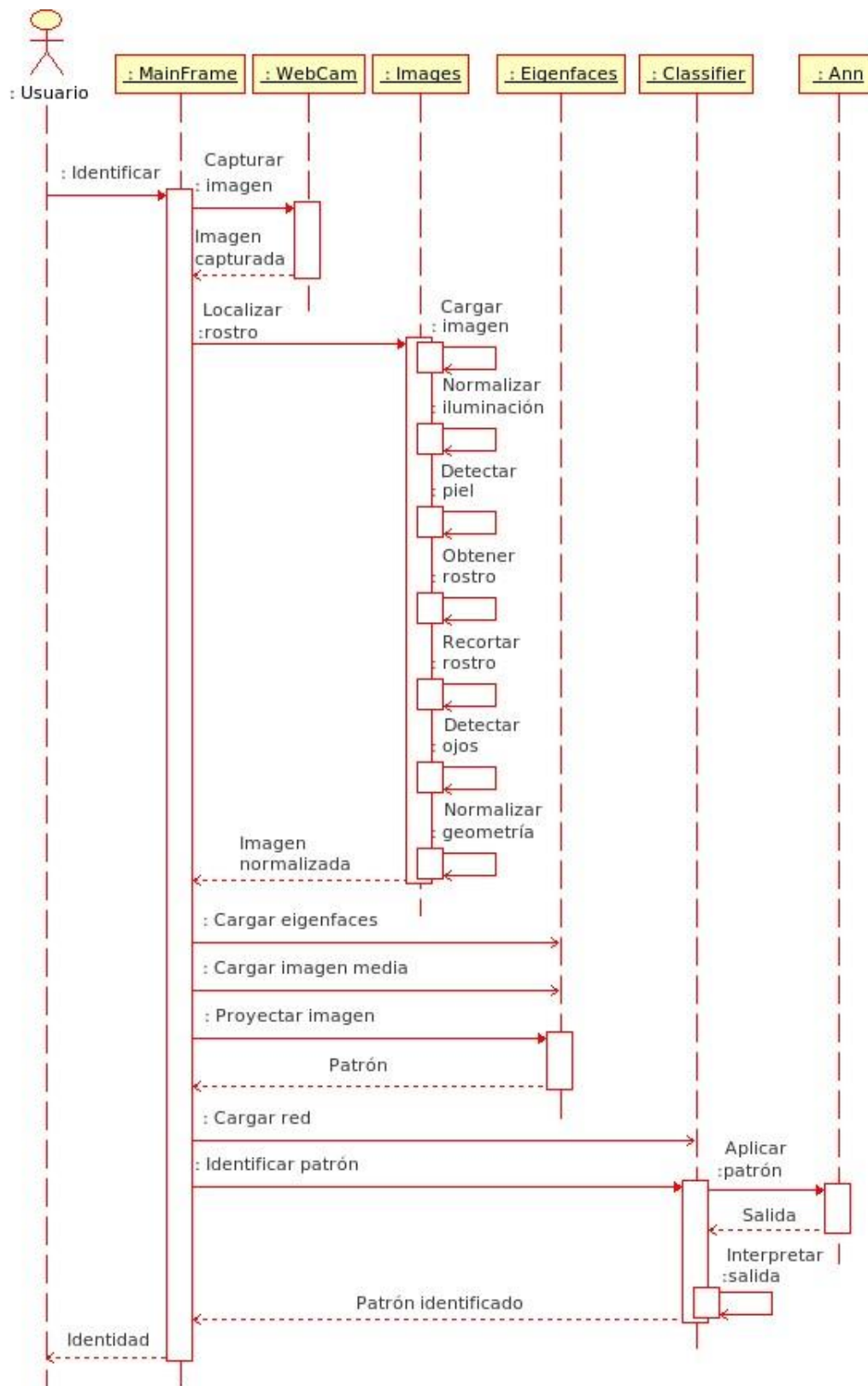
En la Figura 3.4 se puede observar el diagrama correspondiente al momento en que un usuario intenta identificarse ante el sistema. El usuario inicia el reconocimiento al solicitarle a la clase *MainFrame* que lo identifique. Esta instancia a la clase *WebCam* para que tome una foto del usuario. Con la imagen capturada, la clase *Images* la carga, la realza y, a partir de la detección de la piel, ubica el rostro

obteniendo como resultado una subimagen que corresponde a la cara normalizada.

Una vez localizado el rostro, se instancia a la clase Eigenfaces que carga las eigenfaces y la imagen media. Toma la cara normalizada, la transforma en un vector y a este último, lo proyecta al espacio de caras obteniendo un patrón a clasificar. Luego, con la clase Classifier se carga la red neuronal entrenada y se aplica el patrón obtenido a la entrada de la misma. Finalmente, se informa el resultado de la identificación al usuario.

3.4. Diagrama de secuencia

Figura 3.4: Diagrama de secuencia de identificación de persona.



3.5 Interfaz de usuario

A partir de los requerimientos especificados y los distintos diagramas expuestos, se desarrolló la interfaz de usuario implementándose tres funcionalidades: 1) iniciar la cámara web, 2) capturar una imagen y 3) identificar un usuario. Esta se desarrolló en lenguaje C++ utilizando la librería wxWidgets. Como se indicó anteriormente, esta librería permite crear interfaces de usuario basadas en las bibliotecas ya existentes en el sistema, con lo que se obtiene una apariencia muy similar a una aplicación nativa. De esta manera, resulta muy portable entre distintos sistemas operativos, estando disponibles para Windows, OS X, Linux y Unix. Es libre, de código abierto y distribuida bajo una licencia LGPL (similar a la GPL). En la Figura 3.5 se puede observar la interfaz desarrollada. Aquí se pueden notar tres botones en la parte inferior: Vista Previa, Capturar e Identificar. Cada botón realiza una acción como se describe a continuación:

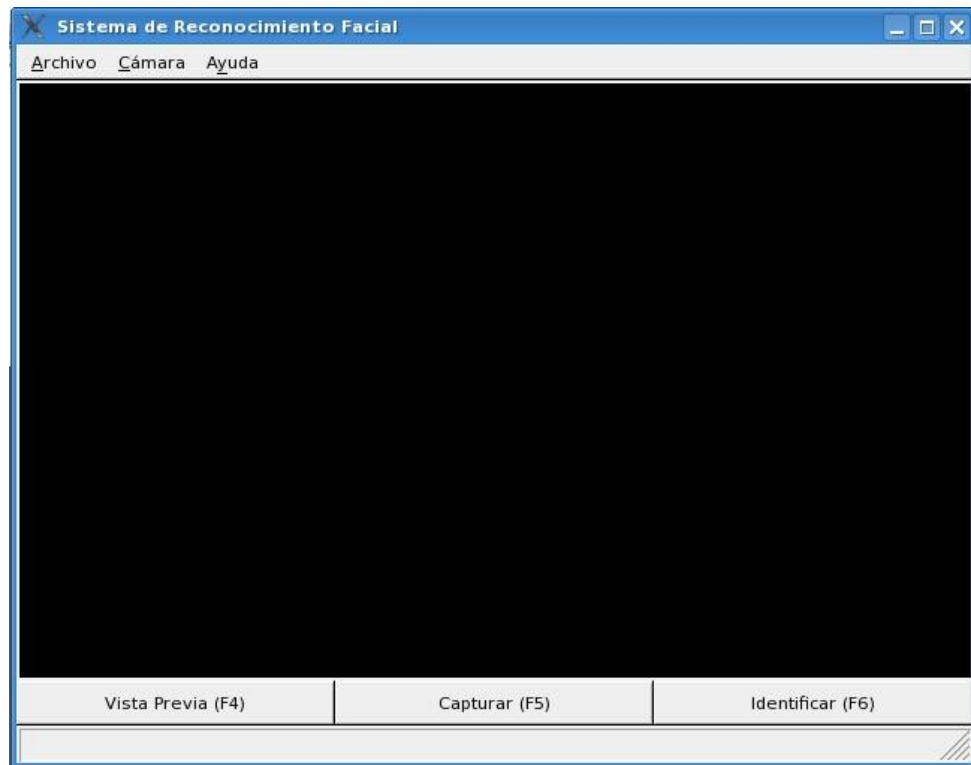


Figura 3.5: Interfaz de usuario.

- Vista Previa: inicializa la cámara web, y reproduce en el recuadro negro lo que ésta captura.
- Capturar: captura una imagen de la vista previa visualizándola en otra ventana, en la cual se puede optar por guardar la imagen tomada (botón Guardar) o descartarla (botón Cancelar). Un ejemplo se muestra en la Figura 3.6.
- Identificar: captura una imagen de la vista previa, aplica las tres etapas para el reconocimiento e indica quién es el usuario que se encuentra frente a la cámara (Ver Figura 3.7).



Figura 3.6: Captura de imagen.

Estas mismas opciones también pueden ser accedidas desde el menú, en la opción Archivo. Además, se debe mencionar que la interfaz funciona sobre los sistemas Windows y Linux.



Figura 3.7: Identificación de usuario.

Las imágenes son capturadas utilizando la cámara web en una resolución de 320 x 240 píxeles. Sobre Linux, se usó el controlador no oficial *spca*³, el cual soporta una variedad de dispositivos. Una lista de éstos puede ser encontrada en su página web. Bajo Windows, se requirió del controlador de la cámara web utilizada.

CAPÍTULO 4

Implementación del sistema

Un sistema de RF consta de tres etapas, como se mencionó anteriormente: (1) detección del rostro, (2) extracción de características y (3) clasificación. En la primera etapa se localiza el rostro en una imagen eliminando todo lo que sea fondo. Una vez obtenida la cara, el segundo paso se encarga de extraer características representativas de ésta, formando así, un vector patrón que se utiliza en la etapa de clasificación para la identificación o verificación de la identidad de una persona con respecto a los usuarios registrados en el sistema. En la Figura 4.1 se puede observar una configuración genérica de un sistema de RF.

A partir de la interfaz de usuario desarrollada, se capturaron un conjunto de imágenes conformando una base de rostros denominada *SINCIDISI*, con la cual se llevaron a cabo las pruebas del sistema. Luego, se desarrollaron

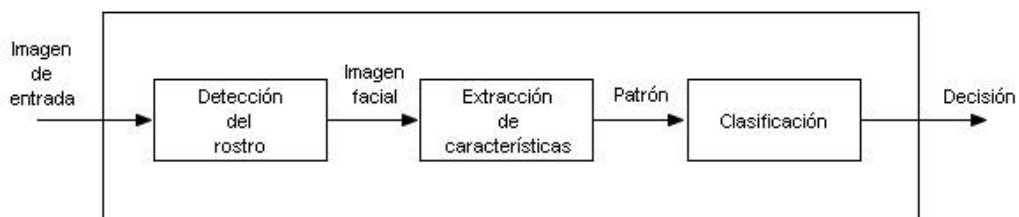


Figura 4.1: Configuración genérica de un sistema de RF.

4.1. Base de datos de rostros SINCIDISI

las tres etapas para el reconocimiento. Previo a la detección del rostro, las imágenes fueron sometidas a un proceso de normalización de la iluminación con el objetivo de realzar las imágenes capturadas. La detección del rostro se basó en el algoritmo propuesto en [ML06], el cual combina detección de piel, aplicación de plantillas y detección de características faciales. Para la extracción de características se proyecta la imagen facial a un espacio de caras, obteniendo un vector patrón que refleja las propiedades más relevantes del rostro. Para la clasificación se diseñó una estructura de red neuronal MLP, constituida por una capa de entrada, una oculta y otra de salida, entrenándola mediante el algoritmo de retropropagación del error. Una vez entrenada, se le presentan los diferentes patrones y devuelve las identidades de los usuarios.

Por otro lado, se debe mencionar que el proceso de construcción se divide en dos fases, una de entrenamiento y una de clasificación. En la primer fase, se captura un conjunto de N imágenes a procesar, denominado *conjunto de imágenes de entrenamiento* $I = \{I_1, \dots, I_N\}$, de un grupo de usuarios. En cada imagen se localiza el rostro obteniendo las imágenes faciales F_i correspondientes, y se calculan los eigenvectores que mejor describen al conjunto facial resultante. Cada rostro F_i es proyectado al nuevo espacio obteniendo un vector característico \mathbf{x}_i . De esta manera, se genera un conjunto de patrones que es utilizado para entrenar la red neuronal. La fase de clasificación toma una imagen de entrada I_q , localiza el rostro y proyecta éste al nuevo espacio, en base a los eigenvectores calculados en la fase de entrenamiento, obteniendo el vector característico \mathbf{x}_q . Luego, este vector es aplicado a la entrada de la red neuronal obteniendo como resultado la identidad del usuario.

4.1 Base de datos de rostros SINCIDISI

Con la interfaz de usuario ya implementada se capturaron varias fotos en un ambiente real, creando así la base de rostros *SINCIDISI*. Para ello, se instaló una cámara web en una PC del centro *sinc(i)*⁴ de la FICH y en una PC del *CIDISF*⁵ de la UTN FRSF (de allí su nombre SINCIDISI), colocándose una fuente de luz artificial para lograr una mejor iluminación en las muestras. Las fotos fueron tomadas con fondo, pose, expresión facial y distancia a la cámara totalmente natural. De esta manera, se obtuvo una cantidad de 448

4.1. Base de datos de rostros SINCIDISI



Figura 4.2: Ejemplos de todos los usuarios de la base de rostros SINCIDISI.

imágenes de 21 personas.

A partir del conjunto de imágenes capturadas, se realizó una selección manual de las mismas a fin de descartar aquellas en las que no se tuvieran variaciones aceptables de pose e iluminación, como

también aquellas con oclusiones parciales del rostro. Además, se eliminaron imágenes en las que se obtuvieron falsas localizaciones del rostro por la presencia de puertas, armarios u otros objetos con un color similar a la piel.

Por otro lado, se notó que la cantidad de fotos para algunos usuarios era escasa y se decidió descartar éstos. A partir de allí, se seleccionaron las imágenes de mejor calidad, en una cantidad igual por cada usuario. De este modo, se obtuvo un conjunto final de 110 imágenes, de un total de 11 personas (10 imágenes por persona). Las fotos fueron almacenadas en formato PPM, con un tamaño de 320 x 240 píxeles. En la Figura 4.2 se exhibe un subconjunto de la base, con ejemplos de todos los usuarios.

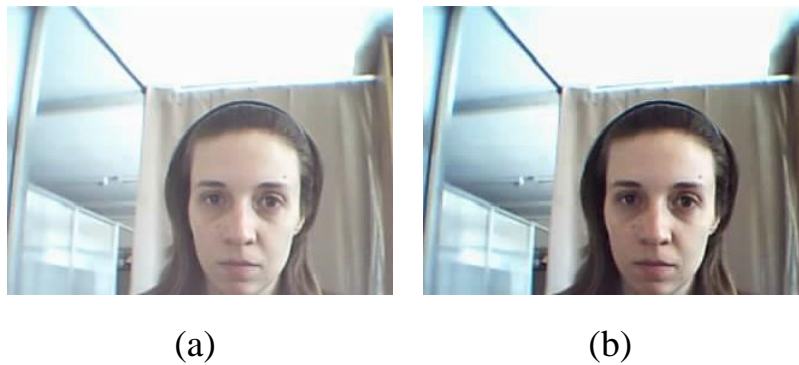


Figura 4.3: Normalización de la iluminación. a) Imagen original, b) Imagen con iluminación normalizada.

4.2 Metodología

4.2.1 Detección del rostro

Un cambio en la distribución de la fuente de luz y en el nivel de iluminación produce un cambio en el color de la piel en la imagen,

degradando el rendimiento de los sistemas de detección de piel [KMB07]. Debido a que el algoritmo aplicado utiliza información sobre el color de piel para ubicar el rostro, es necesario aplicar un proceso de normalización que compense la luminosidad en las imágenes capturadas por la cámara web, ya sean imágenes con poca o mucha iluminación.

Para compensar la luminosidad se aplicó una transformación de potencia

(a cada canal del modelo RGB) dada por

$$f(x,y) = f(x,y)^\gamma,$$

donde γ es adaptada según la cantidad de luz en cada imagen por

$$\gamma = c_1\mu_f + c_2.$$

Aquí, c_1 y c_2 son constantes y μ_f es la media de la imagen a normalizar en escala de grises. Con esto se logra que γ sea directamente proporcional al brillo. En la Figura 4.3 se observa el resultado de aplicar esta transformación.

A partir de la imagen normalizada, como primer paso, el algoritmo localiza el rostro en la imagen. Para esto, realiza una segmentación de piel aplicando

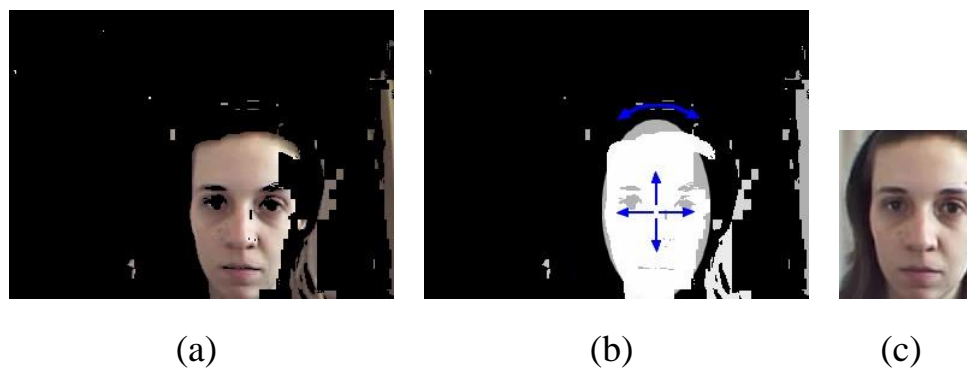


Figura 4.4: Localización del rostro. a) Segmentación de piel, b) Combinación de plantilla sobre las regiones, c) Subimagen con el rostro localizado.

ciertas reglas definidas en [KPS97], basadas en que el color de piel se mantiene en un cierto rango del espacio de color RGB. De esta manera, logra localizar las posibles regiones caras en la imagen. La Figura 4.4(a) muestra el resultado de la segmentación.

Una vez ubicadas las posibles regiones de piel, éstas son etiquetadas mediante un algoritmo recursivo. Luego, se elige aquella región que mejor se asemeje geométricamente a un rostro. Para esto, se aplican plantillas elípticas a las diferentes regiones, haciendo ampliaciones, traslaciones y rotaciones de las mismas para lograr un mejor rendimiento ante variación de pose y distancia a la cámara (Figura 4.4(b)), y se elige aquella región que contiene la mayor cantidad de píxeles que caen dentro de la plantilla. Posteriormente se recorta una subimagen que contiene la región seleccionada, o dicho de otro modo, el rostro del usuario (Figura 4.4(c)).

En segundo lugar, en base a la imagen recortada, el algoritmo normaliza el tamaño y orientación del rostro teniendo en cuenta la distancia entre ojos, posibilitando la comparación de patrones estandarizados geométricamente. Primero detecta las posibles regiones de ojos mediante operaciones morfológicas de dilatación y erosión aplicadas a la componente de luminancia Y del modelo YCbCr, filtrando aquellas que tengan demasiados o muy escasos píxeles como para representar un ojo. En la Figura 4.5(a) y (b) se

observa la componente de luminancia Y y la segmentación resultante respectivamente. Luego, busca el par de ojos candidatos tomando cada región de la mitad izquierda, buscando su par a la derecha que caiga dentro de un rectángulo predefinido (Figura 4.5(c)). Los pares de ojos candidatos resultantes son validados mediante una ponderación de propiedades como posición en la imagen y distancia entre ellos, descartando aquellos que no superan un umbral preestablecido.

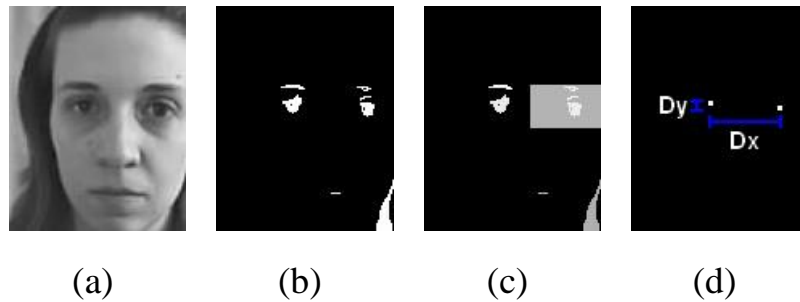


Figura 4.5: Normalización de tamaño. a) Componente Y, b) Posibles regiones ojos, c) Búsqueda del par de ojos, d) Ojos detectados.

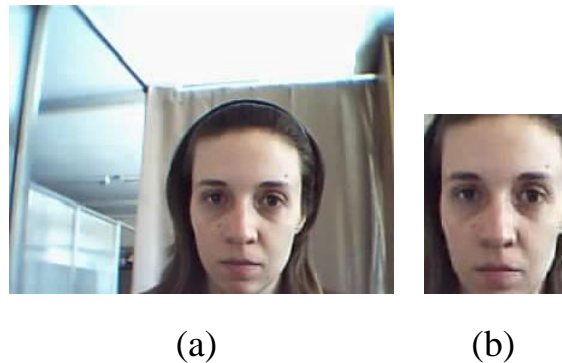


Figura 4.6: Resultado de la detección del rostro. a) Imagen con iluminación normalizada, b) Imagen normalizada.

Con el par de ojos ganador (Figura 4.5(d)), se calcula la distancia ocular horizontal Dx y vertical Dy y el ángulo de inclinación del rostro $\alpha = \arctan(Dy/Dx)$. Se corrige la inclinación del rostro en base al ángulo calculado y luego, en base a Dx y Dy , se recorta la subimagen facial normalizada a partir de la imagen realzada como se observa en la Figura 4.6.

4.2.2 Extracci3n de caracter3sticas

En esta etapa se extraen las caracter3sticas m3as relevantes de las im3genes faciales aplicando el m3etodo *eigenfaces* explicado en la Secci3n 2.2.

Como se explic3o anteriormente, el primer paso es calcular las componentes eigenfaces que describen el conjunto de im3genes faciales aplicando la ecuaci3n (2.13), escrita nuevamente para conveniencia del lector como

$$\begin{aligned} \mathbf{u}_l &= A\mathbf{v}_l \\ (4.3) \\ &= \sum_{k=1}^M \mathbf{v}_{lk} \Phi_k, \end{aligned}$$

donde M es la cantidad de im3genes en el conjunto de entrenamiento. De esta manera, se obtienen M eigenvectores \mathbf{u}_l que conforman el conjunto base de las im3genes faciales.

Dado que no se requiere la reconstrucci3n de las im3genes, se utilizan 3nicamente $M^0 < M$ eigenfaces, las cuales corresponden a las m3as significativas seg3n el orden decreciente de los eigenvalores asociados. Estos M^0 eigenvectores contienen la mayor informaci3n del conjunto de datos.

Una vez calculadas las eigenfaces, tanto en el entrenamiento como en la clasificaci3n, las im3genes faciales son proyectadas al eigenespacio mediante

$$x_k = \mathbf{u}_k^T (\mathbf{I} - \Psi) \quad \text{para } k = 1, \dots, M^0, \quad (4.4)$$

donde I representa una imagen, obteniendo el vector caracter3stico \mathbf{x} correspondiente.

4.2.3 Clasificaci3n

Para el reconocimiento de patrones se disen3o una red neuronal MLP de tres capas (entrada, oculta y salida). Un esquema similar de esta estructura se puede observar en la Figura 2.5.

La entrada a la red se definió según las características extraídas de las imágenes faciales en la etapa anterior, dependiendo de la cantidad de eigenfaces utilizadas para la proyección en el espacio de caras. El número de nodos en la capa oculta se determinó mediante experimentación, donde se observó que cantidad era requerida para obtener un desempeño aceptable. La cantidad de salidas es el número de usuarios a reconocer por el sistema, cada salida representa a uno de éstos. La salida cuyo valor es más alto corresponde al usuario reconocido.

Una vez establecido el diseño, la red se entrenó mediante el algoritmo de retro propagación del error, siguiendo las ecuaciones de la Sección 2.3.1. Esto se realizó en la fase de entrenamiento, en la cual se crearon patrones definiendo la relación entrada-salida para la red. Con dichos patrones, se estimuló la red a fin de minimizar el error entre salida y respuesta deseada. En la fase de

clasificación, la red entrenada es estimulada con el vector característico del usuario que intenta identificarse ante el sistema, y esta retorna la clase a la que pertenece o indica que no pertenece a ninguna de las clases conocidas

Conclusiones y desarrollos futuros

6.1 Conclusiones

En este trabajo se ha implementado un sistema grafico de reconocimiento facial que contempla todas las etapas requeridas, desde la interfaz gráfica de usuario para la captura de la imagen hasta

su clasificación, pasando por la localización del rostro y la extracción de características, identificando al usuario frente a la cámara.

El sistema realizado cumple con los objetivos propuestos logrando un rendimiento aceptable dado el ambiente en el cual se presentó el problema. Este´ presenta una interfaz gráfica sencilla que permite operar fácilmente tanto a un administrador del sistema a la hora de tomar fotos como a un usuario que intenta identificarse. A través de los experimentos se demostró que el clasificador MLP se desempeñó mejor frente a una red RBF y a la técnica del vecino más cercano, logrando un mayor reconocimiento.

6.2 Desarrollos futuros

Con la idea de mejorar el desempeño general del sistema, una mejora estaría dada por un montaje dedicado en el lugar de instalación del mismo. Esto implicaría colocar un fondo fijo de color blanco y un flash como fuente de´ luz para lograr luminosidad constante permitiendo, de esta manera, obtener una mejor calidad en las imágenes capturadas. Por otro lado, una ampliación del sistema estaría dada por el desarrollo de un módulo de administración de usuarios a fin de automatizar el proceso de carga de los mismos. Este incluiría funcionalidades para el registro de nuevos´ usuarios, dar de baja aquellos que ya no se requieran en el sistema, así como también asociar las imágenes capturadas a los usuarios correspondientes.