

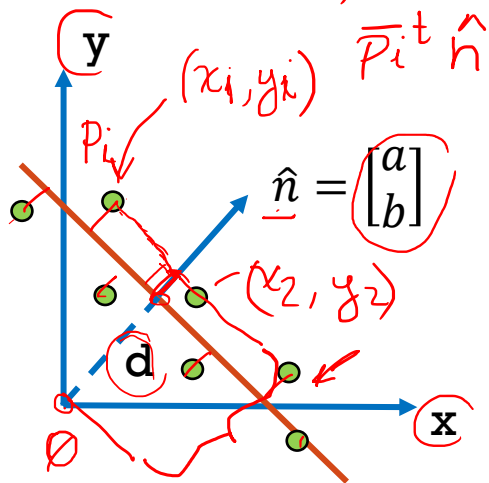
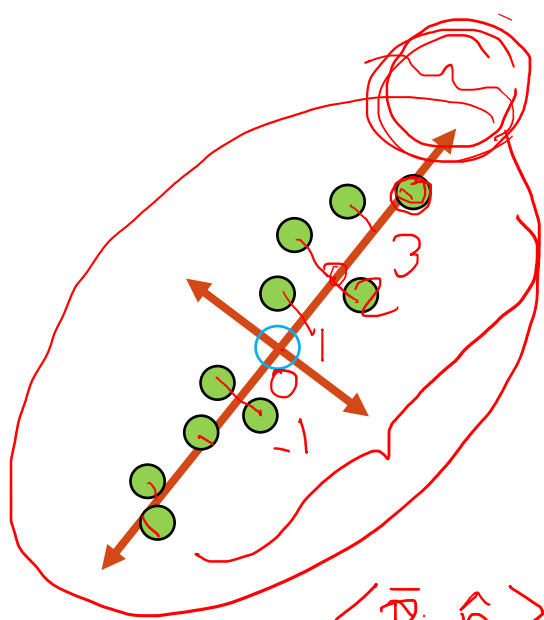
Visión por Computadora I

Ing. Andrés F. Brumovsky
(abrumov@fi.uba.ar)

Laboratorio de Sistemas Embebidos -FIUBA



PCA (PRINCIPAL COMPONENTS ANALYSIS)



Los modelos generativos funcionan mejor en espacios con pocas dimensiones (estimar las funciones de densidad en espacios de muchas dimensiones requieren una gran cantidad de datos)

PCA es una técnica de reducción de dimensiones

1. El primer componente principal es la dirección de máxima varianza (alrededor de la media).
2. Los siguientes componentes principales son ortogonales a los previos y describen la máxima varianza residual

El error de ajuste a una línea lo podemos escribir como

$$E(a, b, d) = \sum_i (ax_i + by_i - d)^2$$

$$\frac{\partial E}{\partial d} = 0 \rightarrow -2 \sum_i (ax_i + by_i - d) = 0$$

$$d = ax_m + by_m$$

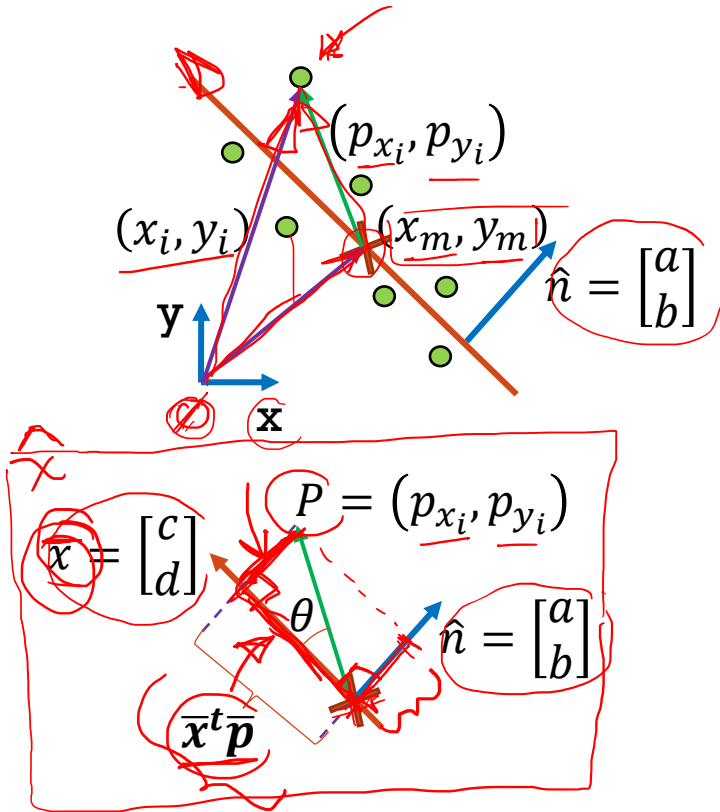
Handwritten notes: $\sum_i ax_i + \sum_i by_i - \sum_i d$, $a \sum_i x_i + b \sum_i y_i - \sum_i d$, $\frac{\sum_i x_i}{n}$, $\frac{\sum_i y_i}{n}$, $n \cdot d$

Substituyendo en la primer ecuación

$$E = \sum_i [a(x_i - x_m) + b(y_i - y_m)]^2 = \|\underline{B}\hat{n}\|_2^2 \rightarrow \underline{B} = \begin{bmatrix} x_1 - x_m & y_1 - y_m \\ \dots & \dots \\ x_n - x_m & y_n - y_m \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix}$$

Donde minimizar $\|\underline{B}\hat{n}\|_2^2$ devuelve los ejes de menor inercia

PCA – INTERPRETACIÓN ALGEBRAICA I



- Por suma de vectores podemos escribir

$$(x_m, y_m) + (p_{x_i}, p_{y_i}) = (x_i, y_i)$$

$$(p_{x_i}, p_{y_i}) = (x_i, y_i) - (x_m, y_m)$$

- Minimizar la proyección sobre la normal $\hat{n}^t \bar{p}$ es equivalente a maximizar la proyección sobre la dirección unitaria en dirección de la línea $\bar{x}^t \bar{p}$
- De la diapositiva anterior sabemos que queremos minimizar la suma de todas estas distancias (para cada punto) al cuadrado ¿Cómo podemos escribir esto de manera algebraica?
- Podemos escribirlo matricialmente de la siguiente manera:

$$E = \bar{x}^t B^t B \bar{x} \text{ sujeto a } \bar{x}^t \bar{x} = 1$$

$$\|B \bar{x}\|_2^2 \rightarrow E$$

$$\langle B \bar{x}, B \bar{x} \rangle$$

$$(B \bar{x})^t B \bar{x}$$

$$\bar{x}^t B^t B \bar{x}$$

- Esto es lo mismo que escribir

$$\sum_i (\bar{x}^t \bar{p}_i)^2$$

- Veamos el caso para solo dos puntos (aunque puede haber m)

$$\begin{aligned} [c \ d] \begin{bmatrix} p_{x_1} & p_{x_2} \\ p_{y_1} & p_{y_2} \end{bmatrix} \begin{bmatrix} p_{x_1} & p_{y_1} \\ p_{x_2} & p_{y_2} \end{bmatrix} \begin{bmatrix} c \\ d \end{bmatrix} &= [c \ d] \begin{bmatrix} p_{x_1}^2 + p_{x_2}^2 & p_{x_1} p_{y_1} + p_{x_2} p_{y_2} \\ p_{x_1} p_{y_1} + p_{x_2} p_{y_2} & p_{y_1}^2 + p_{y_2}^2 \end{bmatrix} \begin{bmatrix} c \\ d \end{bmatrix} \\ &= c \left(c(p_{x_1}^2 + p_{x_2}^2) + d(p_{x_1} p_{y_1} + p_{x_2} p_{y_2}) \right) + d \left(c(p_{x_1} p_{y_1} + p_{x_2} p_{y_2}) + d(p_{y_1}^2 + p_{y_2}^2) \right) \\ &= c^2 p_{x_1}^2 + c^2 p_{x_2}^2 + c d p_{x_1} p_{y_1} + c d p_{x_2} p_{y_2} + c d p_{x_1} p_{y_1} + c d p_{x_2} p_{y_2} + d^2 p_{y_1}^2 + d^2 p_{y_2}^2 \\ &= (c p_{x_1} + d p_{y_1})^2 + (c p_{x_2} + d p_{y_2})^2 \\ &= \sum_i (\bar{x}^t \bar{p}_i)^2 \end{aligned}$$

| \bar{x}^t | B^t | B | \bar{x} |
|-------------|-------|-------------|-----------|
| Línea | | | Línea |
| $c \ d$ | | | $c \ d$ |
| | Pto 1 | $x_1 \ y_1$ | |
| | Pto 2 | $x_2 \ y_2$ | |
| | Pto 3 | $x_3 \ y_3$ | |
| | ... | | |
| | Pto n | | |

PCA - INTERPRETACIÓN ALGEBRAICA II

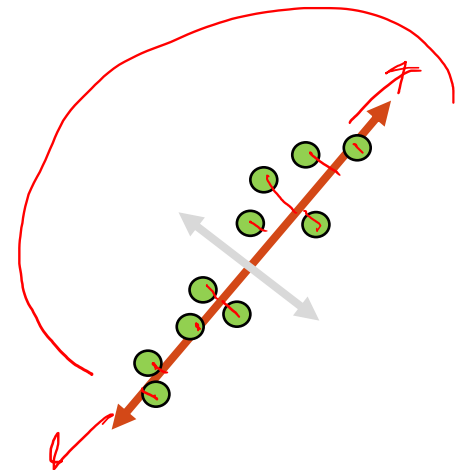
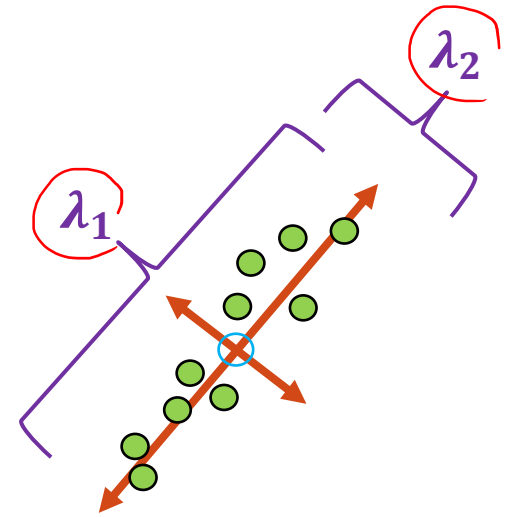
- Dado que lo que queremos hacer es encontrar el máximo de la ecuación anterior pero sujeto a la restricción $x^t x = 1$ (dado que la dirección normal es un vector unitaria) utilizamos el método de multiplicadores de Lagrange.

$$\begin{aligned}
 M &= B^t B \rightarrow E = \bar{x}^t M \bar{x}, & x^t x &= 1 \\
 E' &= \bar{x}^t M \bar{x} - \lambda (\bar{x}^t \bar{x} - 1) & \text{Lagrange} \\
 \frac{\partial E'}{\partial \bar{x}} &= 2M\bar{x} - 2\lambda\bar{x} = 0 \rightarrow M\bar{x} = \lambda\bar{x}
 \end{aligned}$$

- Es decir, x es un autovector de $M = B^t B$
- De la diapositiva anterior también podemos obtener una interpretación algebraica distinta:

$$B^t B = \begin{bmatrix} \sum_{i=1}^n x_i^2 & \sum_{i=1}^n x_i y_i \\ \sum_{i=1}^n x_i y_i & \sum_{i=1}^n y_i^2 \end{bmatrix} \rightarrow \text{matrix de covarianza alrededor del origen}$$

- De donde se deduce que los componentes principales son las direcciones ortogonales de la matriz de covarianza de un conjunto de puntos (los ejes mayores de la elipsoide en el espacio). Es decir, los autovectores de esta matriz de covarianza.
- A su vez los autovalores de esta matriz de covarianza indican la magnitud de las variaciones alrededor de estas direcciones
 - El autovalor λ_M más grande captura la mayor variación alrededor de los vectores de características x
 - El autovalor λ_m más pequeño define las direcciones de menor variación
- Para el grupo de características de la derecha podemos describirlo por su posición a lo largo de la dirección de mayor variación.



EIGENFACES

V → S

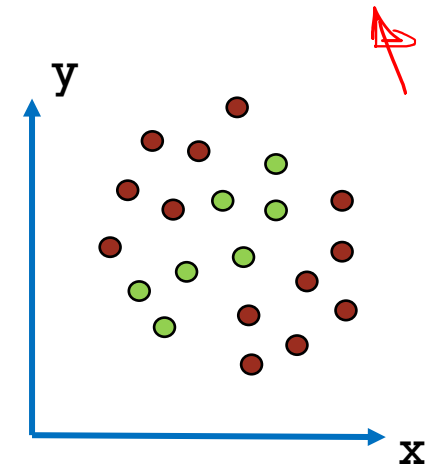
- Supongamos todo el universo de imágenes de 100 x 100 píxeles
 - Si tomamos como característica el valor de intensidad de cada píxel → $\dim(x) = 10000$
 - ¿Cuántas de estas imágenes son solo caras?
- Queremos modelar el subespacio de imágenes de caras
 - Es decir construir un subespacio lineal de menor dimensión que mejor explique la variación en un conjunto de imágenes de caras.
- Supongamos M puntos (x_1, \dots, x_M) en un espacio \mathcal{R}^d con d grande (por ejemplo 10000 o más)
 - Queremos las direcciones u en \mathcal{R}^d que capturen la mayor variación de los x_i , sus coeficientes (proyecciones) serán:

$$u(x_i) = \bar{u}^t (\bar{x}_i - \bar{x}_m)$$

- La varianza de los datos proyectados sobre la dirección \bar{u} se puede escribir entonces como:

$$\begin{aligned} \text{var}(\bar{u}) &= \frac{1}{M} \sum_{i=1}^M \bar{u}^t (\bar{x}_i - \bar{x}_m) (\bar{u}^t (\bar{x}_i - \bar{x}_m))^t \\ &= \bar{u}^t \left[\frac{1}{M} \sum_{i=1}^M (\bar{x}_i - \bar{x}_m) (\bar{x}_i - \bar{x}_m)^t \right] \bar{u} \\ &= \bar{u}^t \Sigma \bar{u} \end{aligned}$$

- Σ es la matriz de covarianza y la dirección de u (vector unitario) que maximiza la varianza es el autovector asociado con el autovalor más grande de Σ .



● Caras
● No caras



EIGENFACES

- Entonces, la dirección que contiene la máxima covarianza de los datos es el autovector correspondiente al autovalor más grande de la matriz de covarianza de los datos.
- Aún más, las k direcciones ortogonales (autovectores) correspondientes a los k autovalores más grandes contienen las direcciones de mayor variación de los datos.
- Entonces podríamos pasar de 10,000 dimensiones (como en el ejemplo) a 40, 30 o inclusive 10 direcciones ortogonales (dimensiones) de mayor variación.
- Ahora, supongamos que tenemos justamente un conjunto de M puntos con N dimensiones.

$$\boxed{B_{M \times N}} \rightarrow (B^t B)_{N \times N}$$

Si $N = 10,000 \Rightarrow (B^t B)$ tiene 100,000,000 datos!

$$\begin{matrix} M \\ 100 \end{matrix} \left[\begin{array}{c} \\ \\ \\ \end{array} \right] \underbrace{\hspace{10em}}_{N = 10,000}$$

- Truco:**

$$\begin{aligned} \rightarrow B^t B \bar{v}_i &= \lambda \bar{v}_i \\ \color{green}{B} B^t B \bar{v}_i &= \lambda \color{green}{B} \bar{v}_i \\ (B B^t) B \bar{v}_i &= \lambda B \bar{v}_i \end{aligned}$$

$(B B^t)_{M \times M}$

- Finalmente, si calculamos los autovectores de $B B^t$ (de dimensión $M \times M$) y los premultiplicamos por la matriz B obtendremos los autovectores de la matriz $B^t B$
- M entonces será la cantidad de caras posibles y entre todo el universo de caras resultará $M \ll N$
- ¿Ahora, cuántos autovectores tiene esta matriz? $\rightarrow M-1$ (los datos tienen restada la media)



EIGENFACES

▪ Mathew A. Turk y Alex P. Pentland – 1991

1. Asumen que la mayoría de las caras se encuentran en un subespacio de dimensión baja definido por las primeras k direcciones de máxima varianza ($k \ll d$)
2. Usan PCA para determinar los autovectores (o “eigenfaces”) u_1, u_2, \dots, u_k que representan el subespacio.
3. Representan todas las imágenes de rostros en el conjunto de datos (dataset) como combinaciones lineales de eigenfaces

▪ Imágenes de entrenamiento $\bar{x}_1, \dots, \bar{x}_M$

▪ Las eigenfaces son en realidad vectores (arreglados con matriz y mostrados como imagen) de 10.000 posiciones (en este ejemplo)

▪ Si quisiéramos hacer el producto interno (proyección) de una cara nueva \bar{x} en este espacio de eigenfaces solo tenemos que superponer la imagen de la cara sobre cualquiera de estas eigenfaces y multiplicar punto a punto.

$$\begin{aligned} \bar{x} &\rightarrow [\bar{u}_1^t(\bar{x} - \bar{u}_0), \dots, \bar{u}_k^t(\bar{x} - \bar{u}_0)] \\ &= [w_1, \dots, w_k] \end{aligned}$$

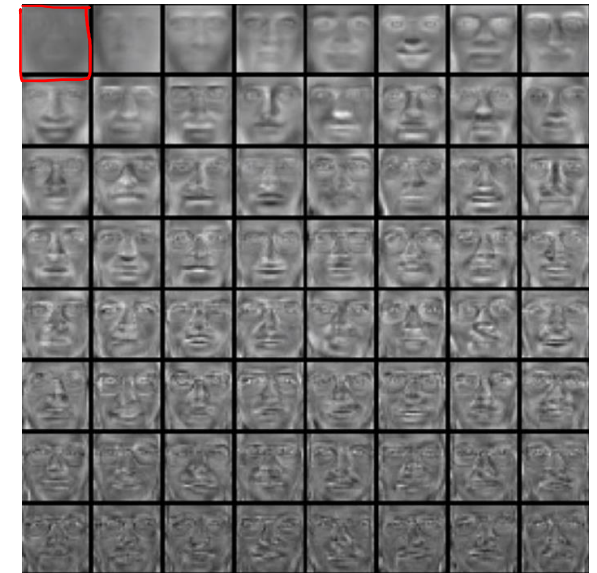
▪ Donde estos pesos w_i son la representación de la cara, de 10.000 valores a k

$$\hat{x} = \bar{u}_0 + w_1 \bar{u}_1 + w_2 \bar{u}_2 + w_3 \bar{u}_3 + w_4 \bar{u}_4 + \dots$$

Dataset



Cara promedio (\bar{u}_0)



Eigenfaces: u_1, \dots, u_k



EIGENFACES



$$\bar{u}_1, \dots, \bar{u}_k$$

$$\bar{u}_0 + 3\sigma_k \bar{u}_k$$

$$\bar{u}_0 - 3\sigma_k \bar{u}_k$$

▪ ¿Cómo lo usamos para reconocimiento?. Dada una nueva imagen \bar{x}

1. Proyectamos en el subespacio: $[w_1, \dots, w_k] = [\bar{u}_1^t(\bar{x} - \bar{u}_0), \dots, \bar{u}_k^t(\bar{x} - \bar{u}_0)]$

2. Opcionalmente podemos verificar el error de reconstrucción $\bar{x} - \hat{x}$ para determinar si es o no una imagen real

3. La clasificamos como la cara de entrenamiento más cercana en el espacio k-dimensional

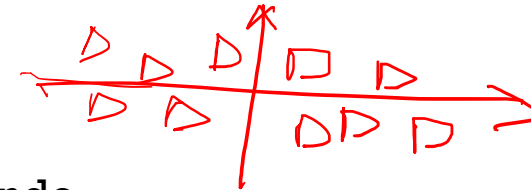
4. Podemos pensar en este procedimiento como en un modelo generativo (se modelan las poblaciones)

▪ Contraste

1. No es robusto frente a desalineaciones

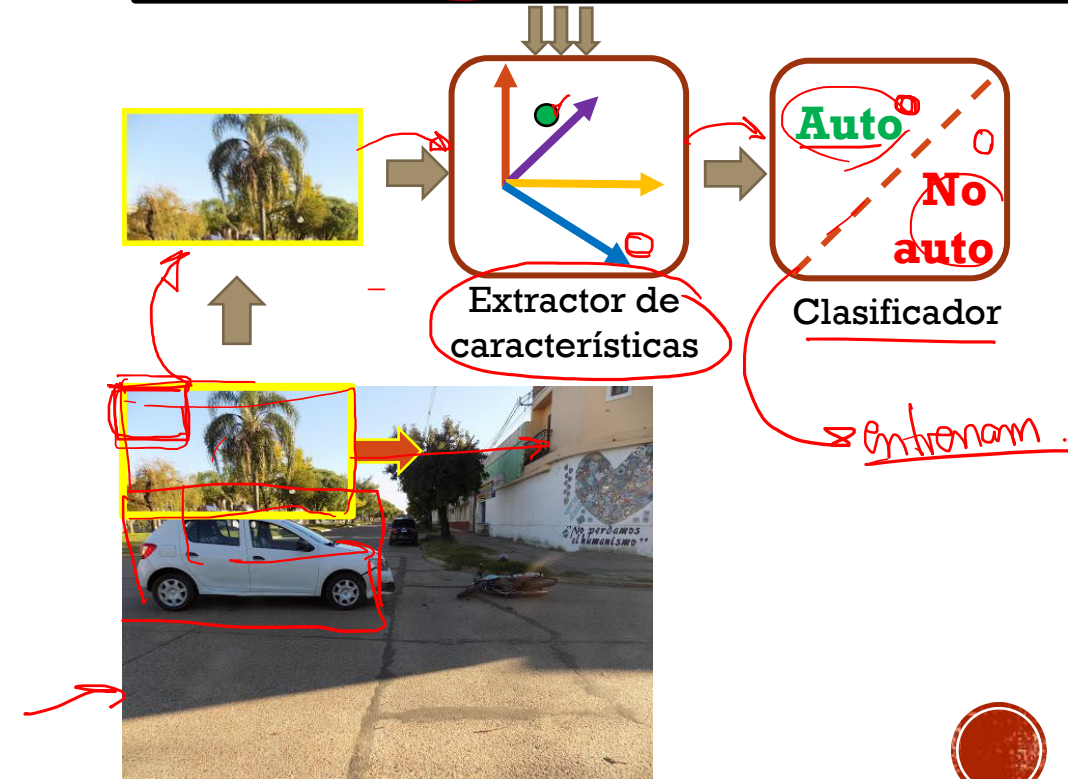
2. No es robusto frente a fuertes variaciones del fondo

3. A veces la dirección de máxima varianza no es la óptima para la clasificación



CLASIFICACIÓN DISCRIMINATIVA

- “No buscamos modelar una clase sino aprender a discriminar entre ellas”
- Distintos métodos de clasificación discriminativa
 - **Nearest Neighbor** – Shakhnarovich, Viola, Darrel (2003) – Berg, Berg, Malik (2005)
 - **SVMs (Support Vector Machine)** – Guton, Vapnik, Heisele, Serre, Poggio (2001)
 - **Boosting** – Viola, Jones (2001) – Torralba et. al. (2004) – Opelt et. al. (2006)
 - **Random Forest** – Breiman (1984) – Shotton, et al CVPR 2008
 - **Redes neuronales** – LeCun, Bottou, Bengio, Haffner (1998) – Rowley, Baluja, Kanade (1998)
- Asumimos
 - Un número fijo conocido de clases (no aparecerán clases nuevas)
 - Amplio número de datos de entrenamiento de cada clase
 - Igual costo de cometer errores de clasificación (entre clases)
 - No sabemos “a priori” las características que definen una clase
- Dividimos el procedimiento en dos partes
 - **Entrenamiento**
 - Conseguir las imágenes modelo: Casos positivos + casos negativos
 - Definir las características a evaluar: LBP, HoG, etc.
 - Enseñar/entrenar un clasificador
 - **Testeo**
 - Generar candidatos (en una nueva imagen)
 - Ventana deslizante, diferentes escalas
 - Puntuar los candidatos

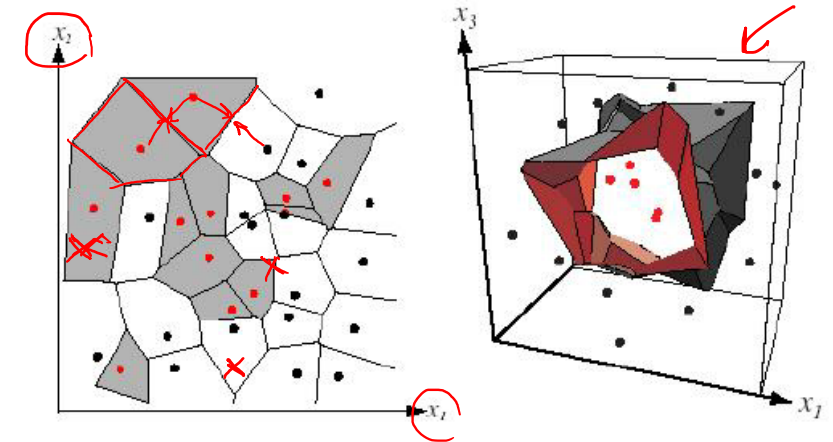


NEAREST NEIGHBOR

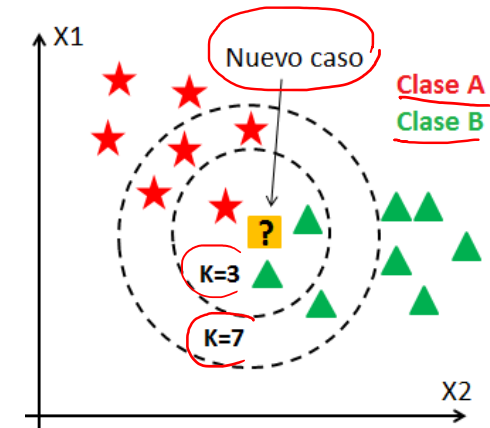
- La idea es asignar a un nuevo caso la etiqueta del punto de entrenamiento más cercano
- La partición del espacio de características es a través de celdas de Voronoi
 - Las fronteras se definen como:
$$R_k = \{x \in X | d(x, P_k) \leq d(x, P_j) \forall j \neq k\}$$
 - La forma de las celdas pueden depender de la definición de distancia considerada
- Es fácil de implementar
- Es muy demandante en cuanto a recursos de memoria y procesamiento en la fase de testeo
 - Se deben recordar todos los datos de entrenamiento
 - Al comprobar un nuevo caso se debe evaluar la distancia contra todos los datos de entrenamiento (o utilizar un KD tree)
- k-NN: Robustece el criterio de decisión
 - Para un nuevo punto encuentra los k puntos más cercanos
 - Realiza una votación entre esos k puntos para clasificar
 - Este método de consenso funciona muy bien, aunque sigue siendo intensivo en el uso de datos.

R → P
N → N

Pattern Classification
Peter Hart, Richard Duda, David Stork



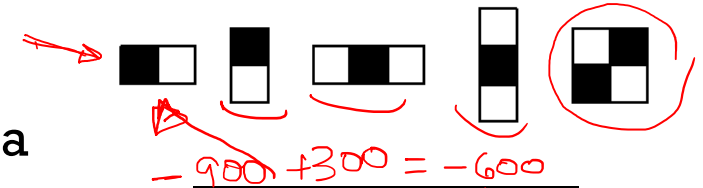
Negros: Casos positivos
Rojos: Casos negativos



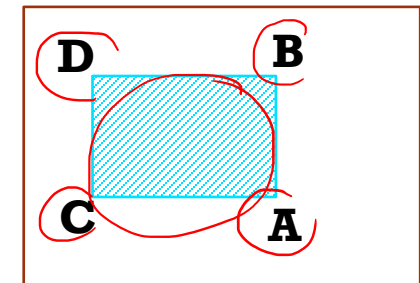
FILTROS DE HAAR

- Buscan detectar patrones de brillo como características de la imagen (diferencias de intensidad entre zonas adyacentes)
- Se pueden definir a múltiples escalas y posiciones
- El resultado del filtro es la diferencia en la suma de los valores de los píxeles entre zonas claras y oscuras
- Se pueden computar muy rápidamente utilizando imágenes integrales → solo tres operaciones de suma para encontrar el área de cualquier tamaño de rectángulo
- Considerando todos los posibles parámetros (posición y escala y tipo) de estos filtros se pueden armar en imágenes de 24x24 más de 180.000 características

Filtros de Haar Básicos



| | | | | | | | |
|-----|-----|-----|-----|-----|-----|-----|-----|
| 200 | 200 | 100 | 100 | 200 | 200 | 100 | 100 |
| 250 | 250 | 50 | 50 | 250 | 250 | 50 | 50 |
| 255 | 255 | 255 | 255 | 100 | 100 | 100 | 100 |
| 255 | 255 | 255 | 255 | 100 | 100 | 100 | 100 |
| 200 | 200 | 100 | 100 | 200 | 200 | 100 | 100 |
| 250 | 250 | 50 | 50 | 250 | 250 | 50 | 50 |
| 255 | 255 | 255 | 255 | 100 | 100 | 200 | 200 |
| 255 | 255 | 255 | 255 | 100 | 100 | 250 | 250 |

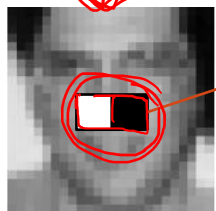


Suma = A - B - C + D



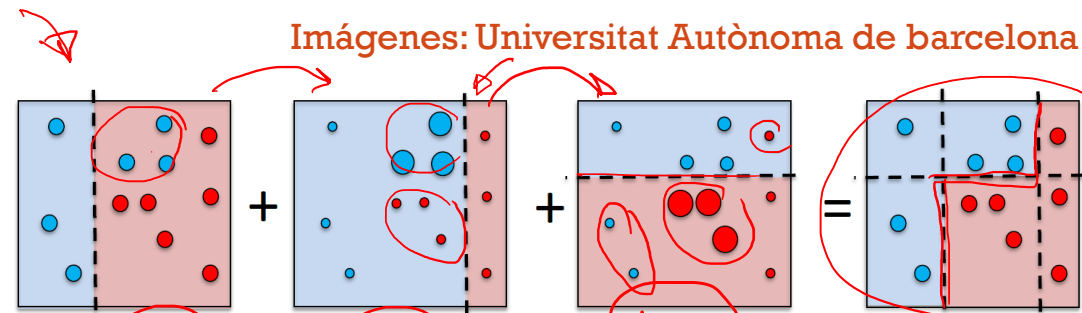
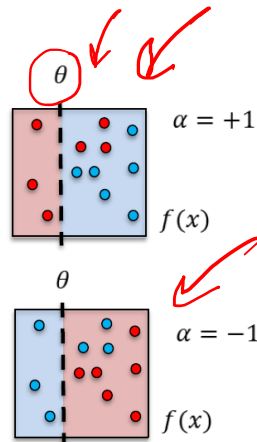
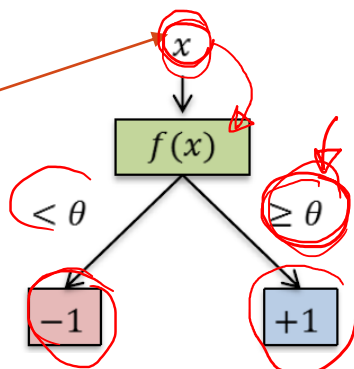
BOOSTING (ADABOOST)

- "Rapid Object Detection using a Boosted Cascade of Simple Features" – Paul Viola, Michael Jones, 2001
- Es un método iterativo que se basa en dos pasos
 1. Encontrar el clasificador "débil" que obtiene el menor error de entrenamiento
 2. Elevar los pesos de las muestras de entrenamiento mal clasificadas por ese clasificador
- En el caso de características de Haar cada clasificador débil se basa en el resultado de la aplicación de una característica de Haar sobre la imagen.

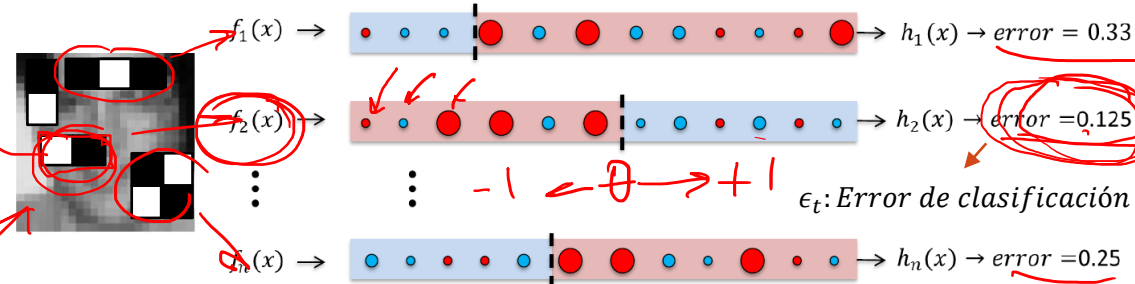
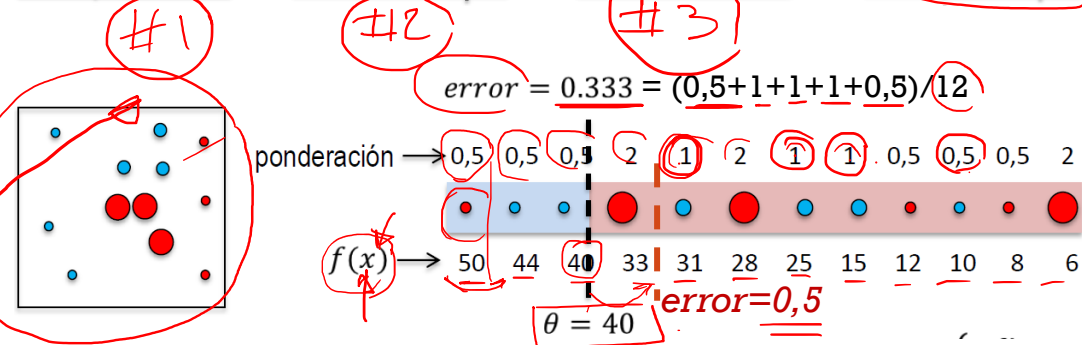


$$h(x) = \begin{cases} -\alpha & f(x) < \theta \\ \alpha & f(x) \geq \theta \end{cases}$$

$$\alpha \in \{-1, +1\}$$



Imágenes: Universitat Autònoma de barcelona



Actualización de pesos

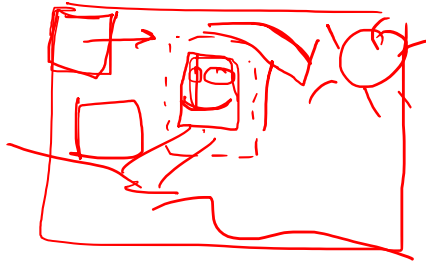
Clasificador fuerte

$$w_i(t+1) = \begin{cases} \frac{1}{2} \frac{w_i(t)}{\epsilon_t} & h(x_i) \neq y_i \text{ (casos mal clasificados)} \\ \frac{1}{2} \frac{w_i(t)}{1 - \epsilon_t} & h(x_i) = y_i \text{ (casos bien clasificados)} \end{cases}$$

$$H(x) = \text{signo} \left(\sum_{i=1}^T \alpha_i h_i(x) \right)$$

$$\alpha_t = \frac{1}{2} \log \frac{1 - \epsilon_t}{\epsilon_t}$$

ADABOOST



▪ Detector en cascada

1. La mayoría de las ventanas en una búsqueda de rostros (u objetos) no contiene rostros!
2. Detectar más rápido si NO hay rostros que si hay
3. Combinación secuencial de clasificadores
4. Una imagen es solo reconocida como rostro si todos los clasificadores la aceptan

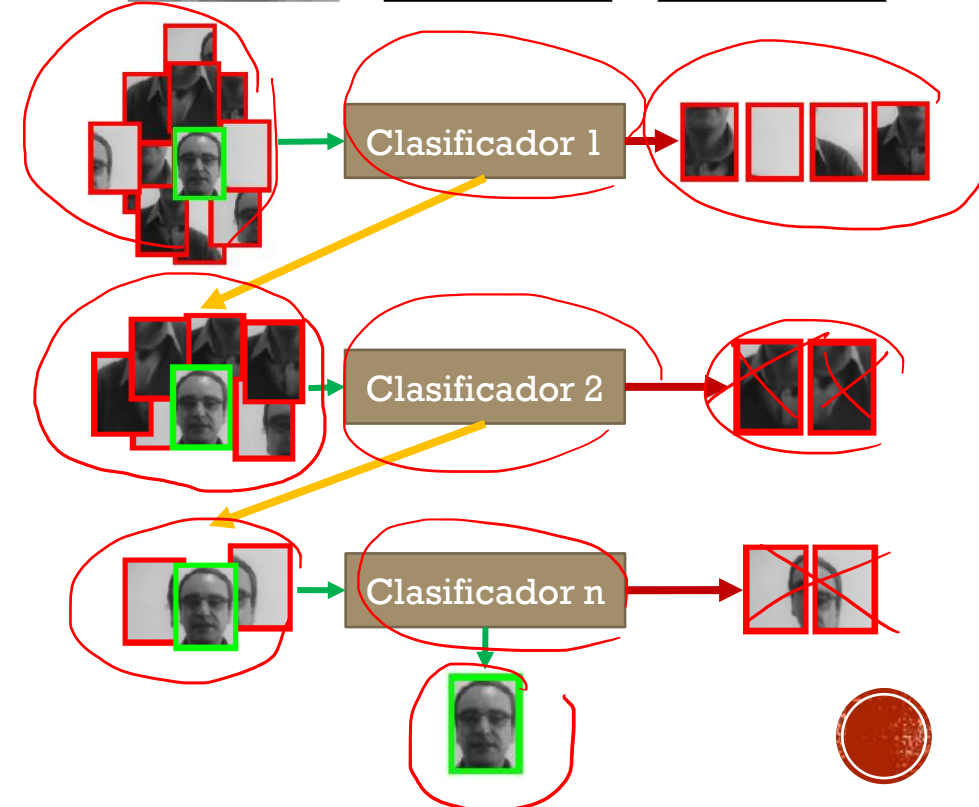
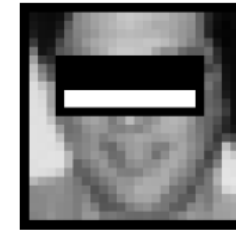
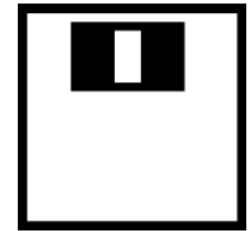
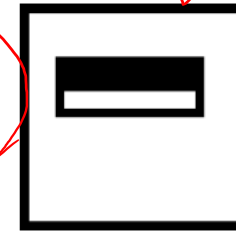
▪ Ventajas

1. Integra la clasificación con la selección de características
2. Flexibilidad en la selección de los clasificadores débiles (esquema de boosting)
3. La complejidad es lineal en el número de ejemplos de entrenamiento
4. La etapa de testeo es realmente rápida
5. Es fácil de implementar

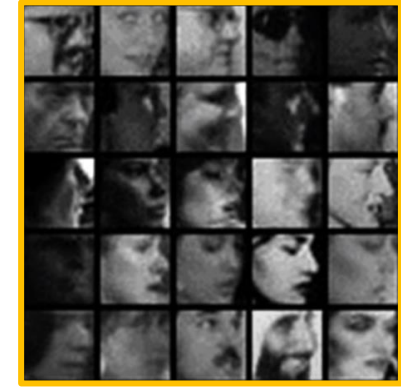
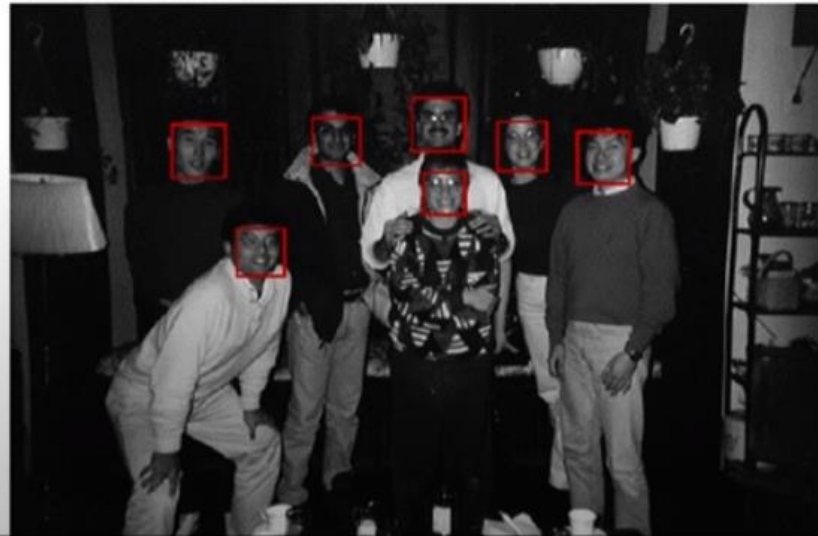
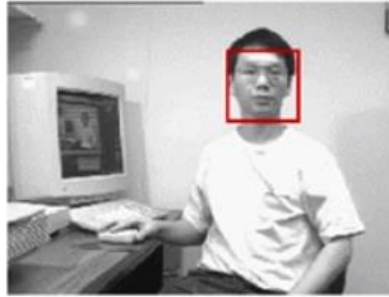
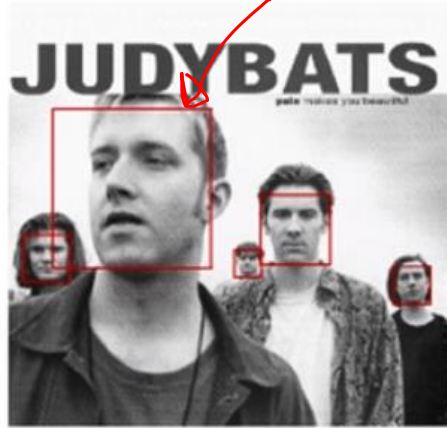
▪ Desventajas

1. Necesita muchos datos de entrenamiento
2. Principalmente en casos con muchas clases no anda tan bien como otros métodos (por ejemplo SVM)

Primeras características seleccionadas por el método anterior



ADABOOST



SVM (SUPPORT VECTOR MACHINES)

- Clasificadores lineales

- Podemos escribir la recta en coordenadas homogéneas como

$$px + qy + b = 0$$

$$\bar{w} = \begin{bmatrix} p \\ q \end{bmatrix} \quad \bar{x} = \begin{bmatrix} x \\ y \end{bmatrix} \Rightarrow \bar{w}^t \bar{x} + b = 0$$

- La distancia de punto (x_0, y_0) a esta recta será (en este caso \bar{w} puede no ser unitario)

$$D = \frac{|px_0 + qy_0 + b|}{\sqrt{p^2 + q^2}} = \frac{|\bar{w}^t \bar{x} + b|}{\|\bar{w}\|}$$

- ¿Ahora...cuál es la mejor recta que separa las poblaciones?

- Línea de separación óptima (en el caso 2D)
- Las líneas punteadas son la distancia más corta a puntos de las poblaciones
- La idea es maximizar el **margen** entre los casos positivos y negativos
- Puntos especiales \rightarrow puntos o “vectores” de soporte (modelo discriminativo)
- Si cambian las muestras del vector de soporte cambia la solución del SVM
- Si cambia cualquier otra muestra no cambia la solución del SVM

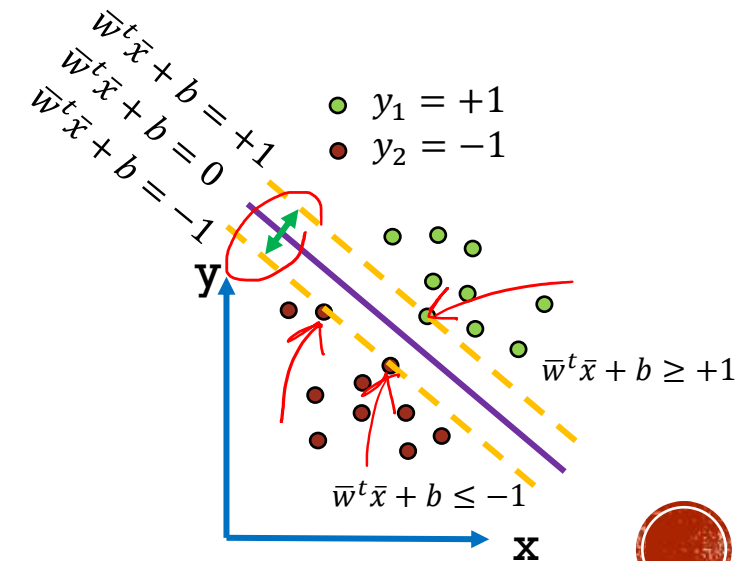
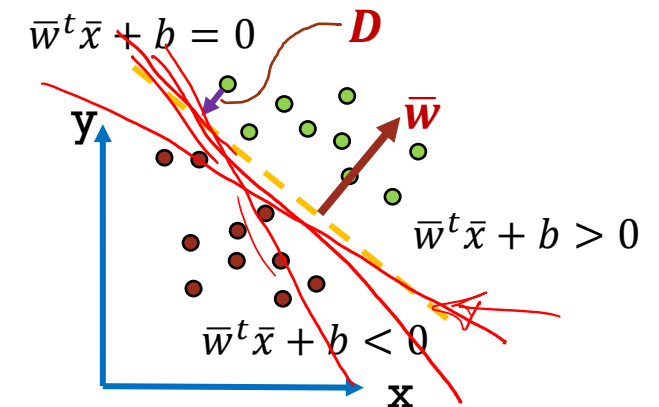
- La condición de clasificación podemos escribirla como

$$y_i(\bar{w}^t \bar{x} + b) \geq 1$$

- ¿Cómo medimos el margen?

$$D_{(+1,-1)} = \frac{|\bar{w}^t \bar{x} + b|}{\|\bar{w}\|} = \frac{|\pm 1|}{\|\bar{w}\|}$$

$$M = D_{+1} + D_{-1} = \frac{2}{\|\bar{w}\|}$$



SVM (SUPPORT VECTOR MACHINES)

- Entonces la idea es maximizar el margen...que sería lo mismo que minimizar

$$\frac{1}{2} \bar{w}^t \bar{w}, \text{ restringido a: } y_i (\bar{w}^t \bar{x}_i + b) \geq 1$$

- Y este es un problema de optimización cuadrática estándar!. Ya sabemos plantear este problema a través del Lagrangiano.

$$\mathcal{L}(\bar{w}, b, \lambda_i) = \frac{1}{2} \bar{w}^t \bar{w} - \sum_i \lambda_i [y_i (\bar{w}^t \bar{x}_i + b) - 1]$$

- La minimización del Lagrangiano (SVM Primal) respecto de \bar{w} y b implica

$$\frac{\partial \mathcal{L}}{\partial \bar{w}} = \bar{w} - \sum_i \lambda_i y_i \bar{x}_i = 0 \Rightarrow \bar{w} = \sum_i \lambda_i y_i \bar{x}_i$$

$$\frac{\partial \mathcal{L}}{\partial b} = - \sum_i \lambda_i y_i = 0 \Rightarrow \sum_i \lambda_i y_i = 0, \lambda_i \geq 0$$

- Reemplazando estas soluciones nuevamente en el Lagrangiano y operando podemos obtener una ecuación para la resolución de optimización (maximización) del SVM Dual

$$\Theta(\lambda) = -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \bar{x}_i^t \bar{x}_j + \sum_i \lambda_i$$

- A partir de la solución para \bar{w} y sabiendo que $b = y_i - \bar{w}^t \bar{x}_i$ para cada vector de soporte podemos plantear una función de clasificación

$$f(\bar{x}) = \text{signo}(\bar{w}^t \bar{x} + b)$$

$$f(\bar{x}) = \text{signo} \left(\sum_i \lambda_i y_i \bar{x}^t \bar{x}_i + b \right)$$

nueva muestra.
vectores soporte



SVM (SUPPORT VECTOR MACHINES)

- ¿Qué pasa si no estamos en un caso linealmente separable?

- Mapeamos los datos a una dimensión superior

$$\Phi: x \rightarrow \varphi(x)$$

- Entonces el producto escalar se transforma en

$$K(\bar{x}_i, \bar{x}_j) = \varphi(\bar{x}_i)^T \varphi(\bar{x}_j)$$

- Este K es el “kernel”. Una función kernel es una función de similitud que corresponde a un producto interno en un espacio expandido de características

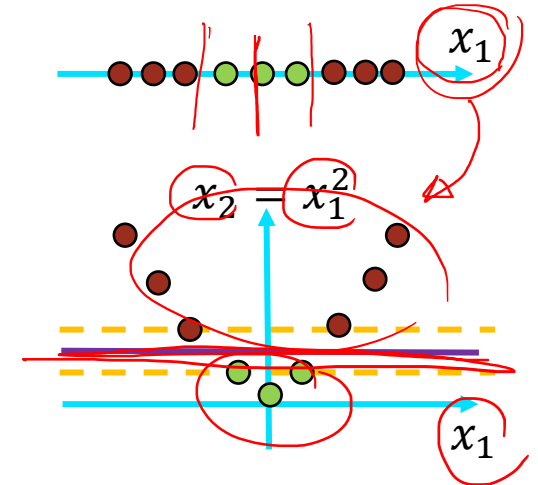
- Ejemplo: $\bar{x} = (x_1, x_2)$, $K(\bar{x}_i, \bar{x}_j) = (1 + \bar{x}_i^t \bar{x}_j)^2$

$$\begin{aligned} K(\bar{x}_i, \bar{x}_j) &= 1 + x_{i1}^2 x_{j1}^2 + 2x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2x_{i1} x_{j1} + 2x_{i2} x_{j2} \\ &= \begin{bmatrix} 1 & x_{i1}^2 & \sqrt{2}x_{i1}x_{i2} & x_{i2}^2 & \sqrt{2}x_{i1} & \sqrt{2}x_{i2} \end{bmatrix}^t \begin{bmatrix} 1 & x_{j1}^2 & \sqrt{2}x_{j1}x_{j2} & x_{j2}^2 & \sqrt{2}x_{j1} & \sqrt{2}x_{j2} \end{bmatrix} \\ &= \varphi(\bar{x}_i)^T \varphi(\bar{x}_j) \end{aligned}$$

$$\varphi(\bar{x}) = [1 \quad x_1^2 \quad \sqrt{2}x_1x_2 \quad x_2^2 \quad \sqrt{2}x_1 \quad \sqrt{2}x_2] \rightarrow \text{cumple}$$

- Entonces en lugar de calcular explícitamente la transformación $\varphi(\bar{x})$, definimos la función kernel K y la incluimos en la función de clasificación para tener “**una decisión de borde no lineal en el espacio original de características**”.

$$f(\bar{x}) = \text{signo} \left(\sum_i \lambda_i y_i \bar{x}^t \bar{x}_i + b \right) \rightarrow \text{signo} \left(\sum_i \lambda_i y_i K(\bar{x}, \bar{x}_i) + b \right)$$



Ejemplos de Kernel

Polinomial

$$K(\bar{x}, \bar{z}) = \langle \bar{x}, \bar{z} \rangle^d$$

Func. Base radial (Gaussian)

$$K(\bar{x}, \bar{z}) = e^{-\|\bar{x} - \bar{z}\|^2 / 2\sigma}$$

Sigmoide

$$K(\bar{x}, \bar{z}) = \tanh(k \langle \bar{x}, \bar{z} \rangle - \delta)$$

Multi-cuadrático inverso

$$K(\bar{x}, \bar{z}) = (\|\bar{x} - \bar{z}\|^{1/2} 2\sigma + c^2)^{-1}$$

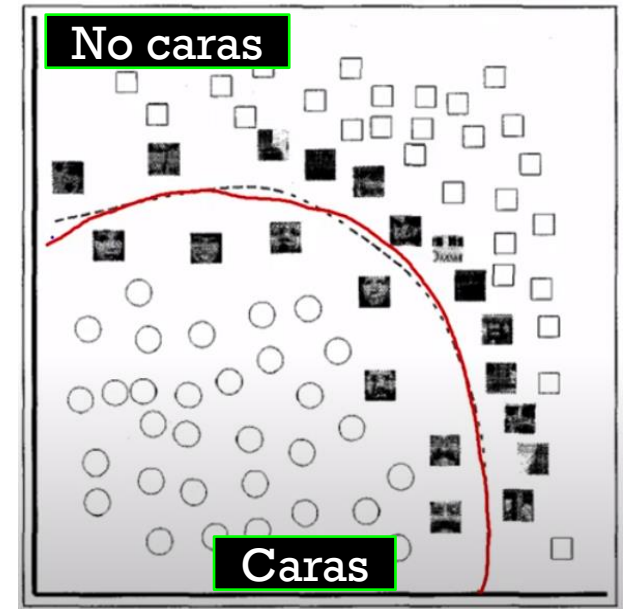
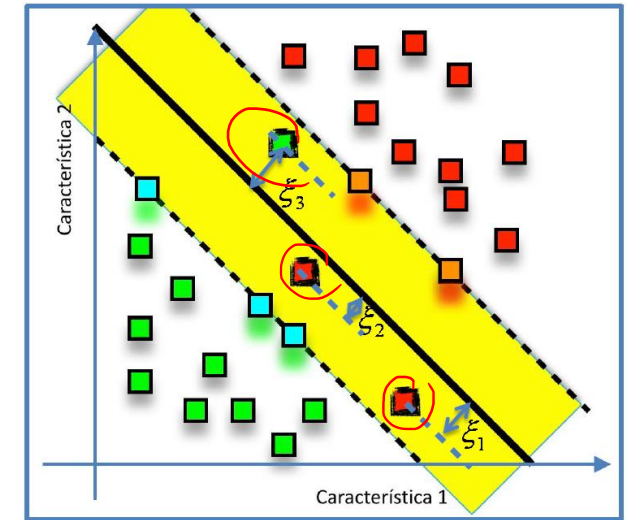
Kernel de intersección

$$K(\bar{x}, \bar{z}) = \sum_{i=1}^n \min(x_i, z_i)$$



SVM (SUPPORT VECTOR MACHINES)

- ¿Qué pasa si aún así no son linealmente separables?
 - Se puede permitir un cierto “margen suave” que permita pequeños errores de clasificación en la frontera
 - Se introduce una variable de holgura $\xi_i \geq 0$ y pedimos
 1. Minimizar $\Phi(\bar{w}) = \frac{1}{2} \bar{w}^t \bar{w} + C \sum_i \xi_i$ sujeto a: $y_i (\bar{w}^t \varphi(\bar{x}_i) + b) \geq 1 - \xi_i$
 2. Maximizar $\Theta(\lambda) = -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j K(\bar{x}, \bar{x}_i) + \sum_i \lambda_i$ sujeto a:
$$\begin{cases} 0 \leq \lambda_i \leq C & b = y_i(1 - \xi_i) - \sum_j y_j \lambda_j K(\bar{x}_j, \bar{x}_i) \\ \sum_i \lambda_i y_i = 0 & f(\bar{x}) = \text{signo} \left(\sum_i \lambda_i y_i K(\bar{x}, \bar{x}_i) + b \right) \end{cases}$$
- Entonces, ¿cuál es el procedimiento?
 - En etapa de entrenamiento
 1. Definir la representación (espacio de características)
 2. Elegir la función de kernel
 3. Calcular los pares del kernel sobre las muestras etiquetadas
 4. Usar esa “matriz kernel” para resolver para los vectores de soporte y los pesos.
 - En etapa de testeo
 1. Calcular los valores del kernel entre la nueva muestra y los vectores de soporte
 2. Aplicar los pesos
 3. Verificar la salida (signo)



SVM (SUPPORT VECTOR MACHINES)

- ¿Y si tenemos muchas clases distintas? → Multi-Class SVMs

1. Uno Vs. Todos

- **Entrenamiento:** Aprender un SVM para cada clase Vs. el resto
- **Testeo:** Aplicar cada SVM a la muestra ejemplo y asignarla a la clase del SVM que devuelva el valor de decisión más alto

2. Uno Vs. Uno

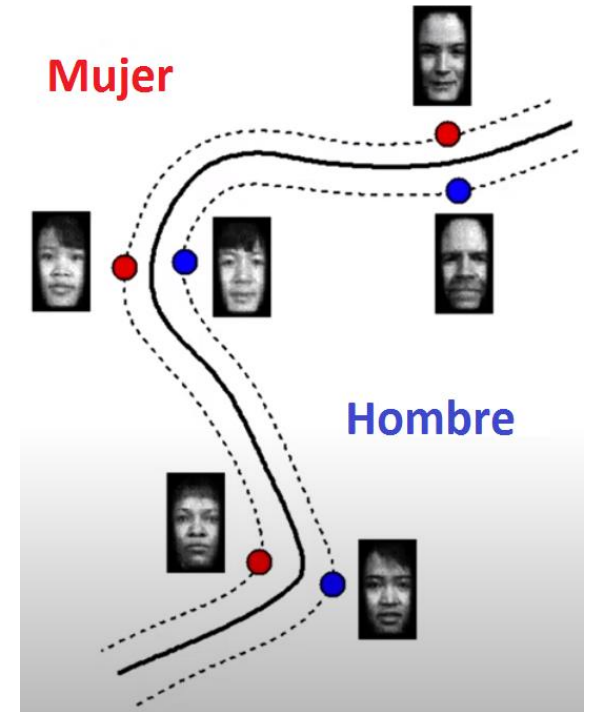
- **Entrenamiento:** Aprender un SVM para cada par de clases
- **Testeo:** Cada SVM “vota” por la clase a ser asignada a la muestra ejemplo

▪ Ventajas

- Hay muchas librerías públicas disponibles
 - <http://www.kernel-machines.org/software>
 - <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>
 - [OpenCV](#)
- La estructura basada en kernels es muy potente y flexible
- Suele devolver un pequeño número de vectores soporte (es rápido en testeo)
- Funciona muy bien en la práctica (inclusive mejor que boosting)

▪ Desventajas

- No hay manera (nativa) de entrenar un SVM multi-clase (se deben combinar)
- Puede ser demandante elegir el kernel correcto para un problema
- Difícil de programar (se suele usar a través de librerías). Se pierde intuición sobre los resultados
- Tiempo de cómputo y memoria
 - Durante el entrenamiento se debe calcular la matriz de kernel para cada par de muestras
 - El aprendizaje puede tomar mucho tiempo para problemas a gran escala



TP6

- Utilizando la webcam
 1. Implementar un algoritmo utilizando filtros de Haar en cascada que:
 - a) Detecte Rostros
 - b) Dentro de los rostros detecte ojos
 - c) Dentro de los rostros detecte sonrisas

