

Analyse et Reconnaissance d'Images

James L. Crowley

M2R IVR

Premier Semestre 2007/2008

Séance 8

16 décembre 2007

Discrimination Linéaire

Notations	2
Fonctions de Discrimination (rappel de Séance 6).....	3
Classification Linéaire Bayesienne.....	5
Le cas des variances blanches ("Matched Filter").....	5
Exemple : Intercorrelation de motifs (NCC).....	6
Quelques Faits sur les Hyperplans.....	7
Estimation par Moindres de Carrées.....	9
Discrimination Linéaire MultiClasse.....	10
Perceptrons	11
Méthodes à Noyaux (Kernel Methods).....	13
Boosting.....	14
Résumé de l'algorithme de Boosting quand $K = 2$	14
Détection des Visage par Ada-Boost.....	15
Integral Images :.....	15
Détection par cascade de classifications.....	17

Sources Bibliographiques :

Viola, P. and Jones, Rapid object detection using a boosted cascade of simple features. In IEEE Conference on Computer Vision and Pattern Recognition., IEEE CVPR 2001.

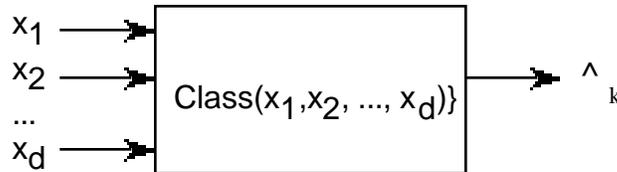
Notations

x	Une variable
X	Une valeur aléatoire (non-prévisible).
N	Le nombre de valeurs possible pour x ou X
x	Un vecteur de D variables
X	Un vecteur aléatoire (non-prévisible).
D	Nombre de dimensions de x ou X
T_k	La classe k
k	Indice d'une classe
K	Nombre de classes
(y_m, X_m)	Une exemple d'événement.
y_m	Un variable d'Indication pour l'exemple X_m $y_m = \{-1, +1\}$
Y_m	Un vecteur d'indication l'exemple X_m pour k classes
M	Nombre totale d'exemples de toutes les classes

Fonctions de Discrimination (rappel de Séance 6)

Soit les événements E décrits par un vecteur de caractéristiques X : (E, X).
 Soit K classes d'événements $\{k\} = \{1, 2, \dots, K\}$ avec $E \in \{k\}$ et i, j tel que $i \neq j$:

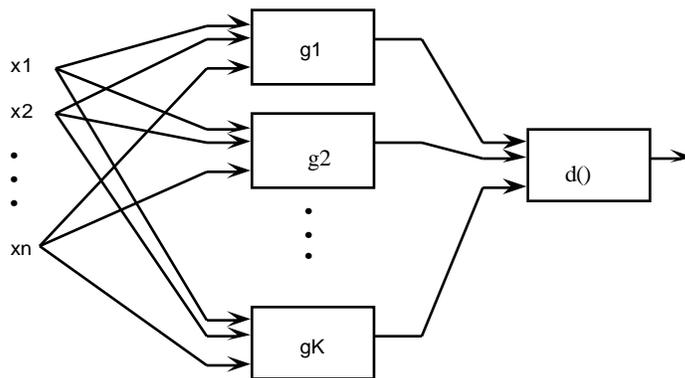
La classification est un processus d'estimation de l'appartenance d'un événement E à une des classes k fondée sur les caractéristiques de l'événement, X.



\hat{k} est la proposition que $(E = k) : \hat{k} = d\{g(X)\}$.
 La fonction de classification est composée de deux parties d() et $g_k()$:

$g(X)$: Une fonction de discrimination : $R^D \rightarrow R^K$
 $d()$: Une fonction de décision : $R^K \rightarrow \{K\}$

Dans cette forme, le classificateur est une machine qui calcule K fonctions $g_k(x)$ suivie d'une sélection de la décision.



Dans le cas générale $D > 1$ et les caractéristiques suivent une densité Normale :

$$p(X | k) = \mathcal{N}(X; \mu_k, C_k) \quad \text{et} \quad g_k(X) = \text{Log}\{ \mathcal{N}(X; \mu_k, C_k) \cdot p(k) \}$$

La fonction de discrimination devient :

$$g_k(x) = -\frac{1}{2} \text{Log}\{\det(C_k)\} - \frac{1}{2}(X - \mu_k)^T C_k^{-1} (X - \mu_k) + \text{Log}\{p(k)\}$$

Ceci peut être traduit dans une forme canonique :

$$g_k(X) = X^T B_k X + W_k^T X + b_k.$$

avec

- une terme 2^{ème} ordre : $B_k = \frac{1}{2} C_k^{-1}$
- une terme 1^{ère} ordre : $d_k = C_k^{-1} \mu_k$
- Constant : $b_k = -\frac{1}{2}(\mu_k^T C_k^{-1} \mu_k) - \frac{1}{2} \text{Log}\{\det(C_k)\} + \text{Log}\{p(\mu_k)\}$

Aujourd'hui nous allons regarder des méthodes linéaires pour concevoir les $g_k(X)$.

Il existe plusieurs techniques simples d'estimation de $g_k()$ comme fonction linéaire. Certains sont très simples et très efficaces. Il ne repose pas sur l'hypothèse de bruit Gaussienne. On peut les combiner une classification linéaire avec la méthode de noyau pour réaliser les systèmes de classification dit "discriminative".

Classification Linéaire Bayésienne.

Le cas des variances blanches ("Matched Filter").

Si $\sigma_{ij} = \sigma^2 \delta_{ij}$ dans ce cas. $\det(C) = (\sigma^2)^n$ et $C_k^{-1} = \frac{1}{\sigma^2} I$

Parce que les termes $\frac{n}{2} \text{Log}\{\sigma^2\}$, $\frac{1}{\sigma^2} I$ et $(\sigma^2)^n$ sont indépendants de i et j ,

$$g_k(x) = -\frac{\|(x - \mu_k)\|^2}{2\sigma^2} + \text{Log}\{p(x_k)\}$$

Ce cas arrive quand les observations sont corrompues par un bruit additif blanc indépendant des classes et d'une puissance égales pour toutes les caractéristiques. Ce cas se rencontre dans les systèmes de réception des signaux hertziens, ainsi que pour la numérisation des images et des sons. En électronique, il est connu comme le cas de la réception "optimal" ("matched Filter")

Ceci est la forme d'un détecteur optimal étudié en théorie de la communication. Pour chaque classe, T_k , le vecteur μ_k est utilisé comme "prototype" ou motif. Avec cette formule, C. Shannon a fait une révolution pour la communication hertzienne en 1946. Son résultat a résolu un problème posé depuis la naissance du télégraphe en 1840 : Combien de message peut-on placer sur un canal de communication ?

Mais parce que $B = \frac{1}{2} C^{-1}$ est indépendant de k , on peut l'éliminer.

Le terme linéaire s'exprime : $W_k = C^{-1} \mu_k^T$

et le constant est $b_k = -\frac{1}{2}(\mu_k^T C^{-1} \mu_k) + \text{Log}\{p(x_k)\}$

Mais, si tous les message ont le même probabilité :

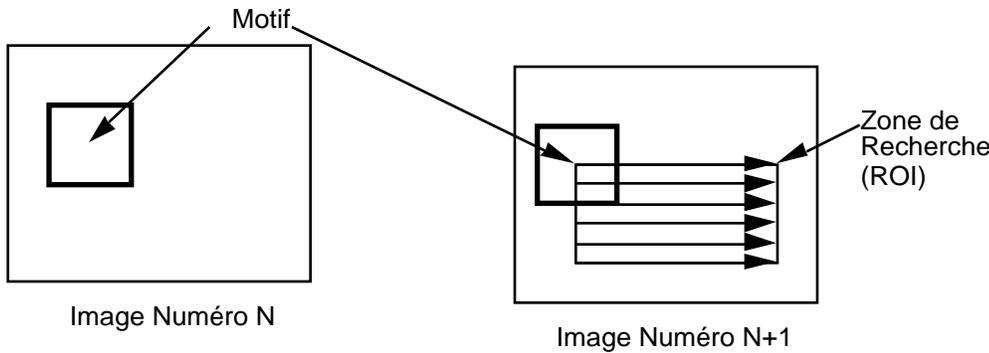
$$b_k = -\frac{1}{2}(\mu_k^T \mu_k) = -\frac{1}{2} \|\mu_k\|^2$$

Donc, on peut écrire $g_k(X) = W_k^T X + b_k = \mu_k^T X - \frac{1}{2} \|\mu_k\|^2$

Exemple : Intercorrelation de motifs (NCC).

(Normalised Cross Correlation). Il s'agit d'une technique d'analyse d'image utilisée pour suivi de cible dans une séquence d'images.

Problème : Soit deux image $S_t(i,j)$ et $S_{t+1}(i,j)$.
 Soit le voisinage (imagerie) $M(i, j)$ issu du $S_t(i,j)$ a position (i_0, j_0) .
 Retrouver sa position (i_1, j_1) dans l'image $S_{t+1}(i,j)$.



Les classes sont les imagerie de l'Image t+1. Ils sont égaux : $p(i) = p(j)$.
 les variances sont égales : $\sigma_i = \sigma_j$
 les variances des pixels sont indépendantes : $\sigma_{ij} = 0$ Donc $C_k = I$

Pour éviter les variations d'intensité, on normalise les imagerie.

Si les vecteurs M et S ont une longueur unitaire, le produit scalaire est un cosinus de l'angle entre les vecteurs.

$$M_u(m, n) = \frac{M}{\|M\|} = \frac{M(m, n)}{\sqrt{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} M(m, n)^2}}$$

$$S_u(m, n) = \frac{S}{\|S\|} = \frac{S(i_1+m, j_1+n)}{\sqrt{\sum_{m=0}^{M-1} \sum_{n=0}^{N-1} S(i_1+m, j_1+n)^2}}$$

On obtient un inter corrélation "normalisée" par l'énergie (NCC) :

$$NCC(i_1, j_1) = \langle M_u, S_u \rangle = \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} \frac{S(i_1+m, j_1+n)}{\|S\|} \frac{M(m, n)}{\|M\|}$$

Le NCC est le cosinus entre les vecteurs M et S . Sa valeur est entre -1 et 1 .

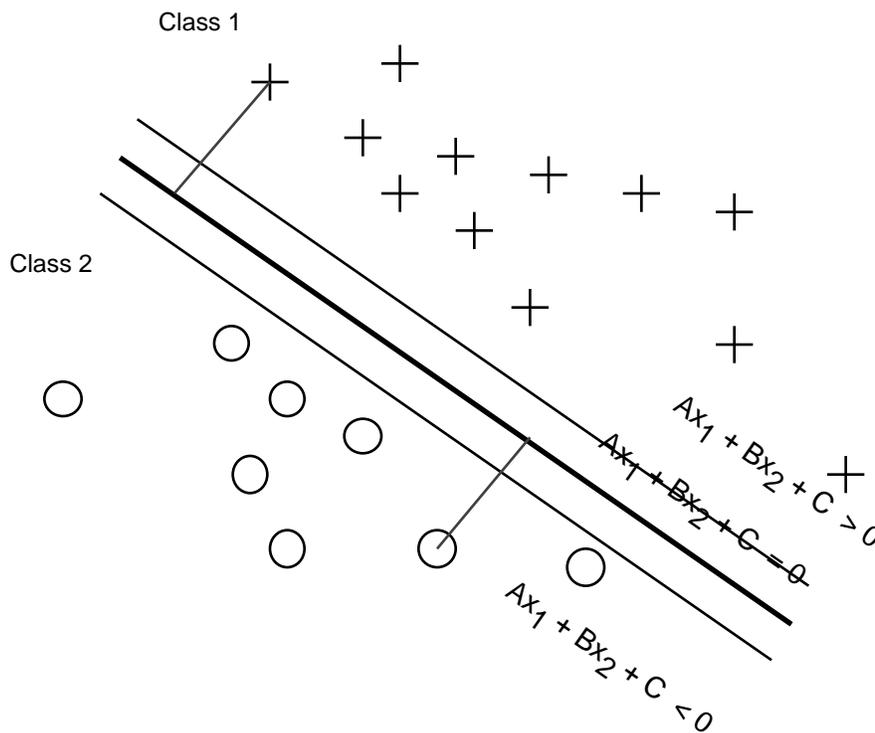
Quelques Faits sur les Hyperplans

Soit $K = 2$ (Deux classes)

Soit M exemples $\{y_m, X_m\}$ tel que $y=+1$ pour les événements de T_1 , et $y = -1$ pour les événements de T_2

Les y_m s'appel les variables d'indication.

Notre objective est d'estimer un plan (ou hyperplan si $D > 2$) qui sépare les deux classes. dans ce cas la fonction de décision devient $d() = \text{sgn}()$.



Un (hyper)plan est un ensemble de points tel quel

$$w_1x_1 + w_2x_2 + \dots + w_Dx_D + b = 0$$

En forme de vecteur : $W \cdot X + b = 0$

ou $W = \begin{pmatrix} w_1 \\ w_2 \\ \dots \\ w_D \end{pmatrix}$ est la norme du plan.

Si $\|W\| = 1$, alors pour tous les points hors du plan,

$b = -W \cdot X$ est la distance perpendiculaire à l'origine.

si $\|W\| = 1$ alors utilise $W' = \frac{W}{\|W\|}$, et $b' = \frac{b}{\|W\|}$

Dans ce cas, on peut dire que $y = W \cdot X + b$

projet les caractéristiques d'un événement sur une norme, W .

Le plus que y est grande, le plus que X est semblable à T_1 .

Pour certaines opérations, il nous faut les coordonnées homogènes.

$$\hat{X} = \begin{pmatrix} X \\ 1 \end{pmatrix} \quad \hat{W} = \begin{pmatrix} W \\ b \end{pmatrix}$$

Dans ce cas on peut écrire $y_m = \hat{W} \cdot \hat{X}_m = \hat{X}_m \cdot \hat{W} = w_1x_1 + w_2x_2 + \dots + w_Dx_D + b$

Estimation par Moindres de Carrées

Discrimination Binaire : $K=2$

Soit M exemples $\{y_m, X_m\}$ tel que $y=+1$ pour les événements de T_1 ,
et $y = -1$ pour les événements de T_2

On cherche $y = g(X) = W X + b$ (en coordonnées Homogènes).

qui minimise la fonction de "Loss" : $L(\hat{W}) = \sum_{m=1}^M (y_m - W X_m)^2$

Pour l'ensemble des M exemples $\{Y_m, X_m\}$. on compose le matrice \mathbf{X} et une vecteur \mathbf{Y}

$$\mathbf{X} = (X_1, X_2, \dots, X_M) = X_m^d \text{ (taille } D \text{ lignes, et } M \text{ colonnes)}$$

$$\mathbf{Y} = (Y_1, Y_2, \dots, Y_M) = Y_m \text{ (taille } M \text{ lignes).}$$

$$\text{on a } L(W, b) = L(\hat{W}) = (\mathbf{Y} - \hat{W} \mathbf{X}) (\mathbf{Y} - \hat{W} \mathbf{X})$$

Le Minimum est trouvé quand :

$$\frac{\partial L(\hat{W})}{\partial \hat{W}} = -2 \mathbf{X} \mathbf{Y} + 2 \hat{W} \mathbf{X} \mathbf{X} = 0$$

$$\text{Donc : } \hat{W} \mathbf{X} \mathbf{X} = \mathbf{X} \mathbf{Y} \text{ et donc } \hat{W} = (\mathbf{X} \mathbf{X})^{-1} \mathbf{X} \mathbf{Y}$$

Discrimination Linéaire MultiClasse

Pour le cas ou $K > 2$ Nous allons chercher k vecteurs W_k tel que

$$k = \arg\text{-max}_k \{ W_k \cdot X + b \}$$

Pour chaque X On définit un K dimensionnel vecteur d'indication, Y .

$$Y = Y^k = \begin{pmatrix} y_1 \\ y_2 \\ \dots \\ y_K \end{pmatrix}$$

Pour une exemple de la classe k , le k th coefficient vaut 1, les autres -1 .

$$y_k \hat{=} \begin{cases} 1 & \text{E} \quad k \\ -1 & \text{sinon} \end{cases}$$

Pour l'ensembles des M exemple $\{Y_m, X_m\}$. on compose des matrices X et Y

$$X = X_m^d = (X_1, X_2, \dots, X_M)$$

$$Y = Y_m^k = (Y_1, Y_2, \dots, Y_M)$$

X a D lignes et M colonnes. Y a K lignes et M colonnes.

Soit $K=2, D=3, M = 4$.

$$X_m^d = \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} \quad Y_m^k = \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix}$$

$$W = (X X^T)^{-1} X Y^T = \left(\begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} \right) \begin{pmatrix} \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \\ \bullet & \bullet & \bullet & \bullet \end{pmatrix} = \begin{pmatrix} \bullet & \bullet \\ \bullet & \bullet \\ \bullet & \bullet \end{pmatrix}$$

W est une matrice composée de D lignes et K Colonnes.

$$k = \arg\text{-max}_k \{ W_k \cdot X + b \}$$

Perceptrons

Un "perceptron" est une méthode incrémentale d'apprentissage inventé par Frank Rosenblatt en 1956. Il s'agit d'une méthode "en-ligne", dirigé par les erreurs. Une perception génère une ensemble d'Hyperplans pour séparer les exemples des classes. Si les exemples peuvent être parfaitement séparé, on dit que les classes sont "séparables". Sinon, ils sont "non-separable".

Le "marge", γ , est le plus petit séparation entre deux classes.

Si les exemples sont "séparables", l'algorithme d'apprentissage utilise les erreurs pour une mise a jour du plan jusqu'à l'il n'y a plus d'erreur. Si les exemples ne sont pas séparables, la méthode ne convergera pas, et il faut arrêter après un certain nombre de cycles.

A chaque cycle, on utilise les erreurs pour adapter le plan de séparation.

Note que pour tous les M exemples :

$$y_m(W \cdot X_m + b) = \begin{cases} 1 & \text{Si la classification est correcte} \\ -1 & \text{Si la classification est en erreur} \end{cases}$$

l'algorithme de perceptron utilise un "gain" positif η pour déterminer la vitesse d'apprentissage.

Algorithme :

```

W0 = 0; b0 = 0; i = 0;
R = max { || Xm || }
REPEAT
  FOR m = 1 TO M DO
    IF ym(Wi · Xm + bi) ≤ 0 THEN
      Wi+1 = Wi + η ym Xm;
      bi+1 = bi + η ym R2;
      i = i + 1;
    END IF
  END FOR
UNTIL no mistakes in FOR loop.

```

La marge pour chaque exemple est :

$$\gamma_m = y_m(W_i \cdot X_m + b_i)$$

Si les coefficients sont normalisés, le marge devient "géométrique" ou "Euclidienne".

$$W' = \frac{W}{\|W\|}, \quad b' = \frac{b}{\|W\|}$$

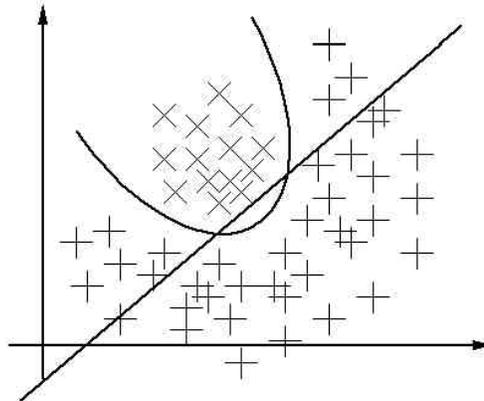
Le "qualité" d'un perceptron est donnée par la distribution des, par exemple, par un histogramme des marges géométrique.

La règle de décision est $d(g(X)) = \text{sgn}(g(X))$

Méthodes à Noyaux (Kernel Methods)

Parce que les fonctions de discrimination linéaire sont si simples à estimer pour $D \gg 1$, il est intéressant de voir si on peut les appliquer dans les cas où les données ne sont pas séparables par les plans.

Par exemple, si les covariances ne sont pas égales, une frontière quadratique donne une meilleure séparation entre les classes.



On peut transformer une discrimination linéaire en discrimination quadratique par substitution de variables.

Il s'agit de projeter un vecteur avec D dimensions dans un espace en $P > D$ dimensions avec un noyau, $K()$.

Par exemple :

$x = (x_1, x_2, \dots, x_D)$ peut être projeté sur un vecteur de $P = \frac{D(D+1)}{2}$ dimension

$w = (x_1, x_2, \dots, x_D, x_1^2, x_1x_2, x_1x_3, \dots, x_{D-1}x_D, x_D^2)$

Ainsi, une fonction quadratic en $D = 2$ dimensions est linéaire en $P = 5$ dimensions

$$x = (x_1, x_2) \quad K(x) = w = (x_1, x_2, x_1^2, x_1x_2, x_2^2)$$

$$g(K(x)) = g(w) = aw_1 + b w_2 + c w_3 + d w_4 + e w_5 \\ = ax_1 + b x_2 + c x_1^2 + d x_1x_2 + e x_2^2$$

Autres Noyaux populaires :

$$\text{Gaussien : } K(x) = \mathcal{N}(x; \mu, \sigma^2)$$

$$\text{Log : } K(x) = \ln(x)$$

Les noyaux sont utilisés avec les méthodes linéaires telles que la Régression ou les Perceptrons pour rendre le système efficace.

Boosting

Une des idées les plus innovantes en apprentissage des dernières années est le "boosting". Le principe est de composer un comité de classificateurs linéaires faibles.

Leur combinaison par vote donne un classificateur fort.

Avec l'algorithme ADA Boost, on peut déterminer une ensemble de classifieurs linéaire faible pour lequel le taux d'erreur est arbitrairement limité.

L'idée est d'appliquer un poids aux exemples, et de renforcer le poids des exemples mal classés, et ré-estimer une nouvelle classifieur.

Résumé de l'algorithme de Boosting quand $K = 2$

Soit deux classes T_1 et T_2 . Soit une ensemble des exemples $\{Y_m, X_m\}$

1) Initialiser un vecteur de M coefficients A avec $a_m = 1$. Initialiser $S = M$. $i = 1$ (a_m serait le poids de chaque exemple, $\{Y_m, X_m\}$. S est la somme des poids.

2) Pour l'ensemble des M exemples $\{Y_m, X_m\}$. on compose le matrice X_m^d et un vecteur Y_m pour calculer le i^{eme} fonction de discrimination avec le poids a_m .

$$\mathbf{X} = (\hat{X}_1, \hat{X}_2, \dots, \hat{X}_M) \quad (\text{taille } D \text{ lignes, et } M \text{ colonnes})$$

$$\mathbf{Y} = (a_1 Y_1, a_2 Y_2, \dots, a_m Y_M) \quad (\text{taille } M \text{ lignes}).$$

On trouve une classification, par exemple avec :

$$\mathbf{W}_i = (\mathbf{X} \mathbf{X}^T)^{-1} \mathbf{X} \mathbf{Y}^T$$

3) Pour chaque exemple en $\{X_m\}$, si la classification est en erreur, on accroître son poids a_m

si $d\{g_n(X)\} \quad k$ alors error :
 $a_m = a_m + 1, S = S + 1,$

4) Normalise les poids a_m : $a_m := a_m / \sum_{m=1}^M a_m$

5) $i = i+1$; Repeter étape 2.

La Classification final est fait par une Vote des classifieurs "boostés"

$$k = \arg\text{-max}_k \left\{ \sum_{i=1}^I W_{ik} \mathbf{X} + b_i \right\}$$

Détection des Visage par Ada-Boost.

Integral Images :

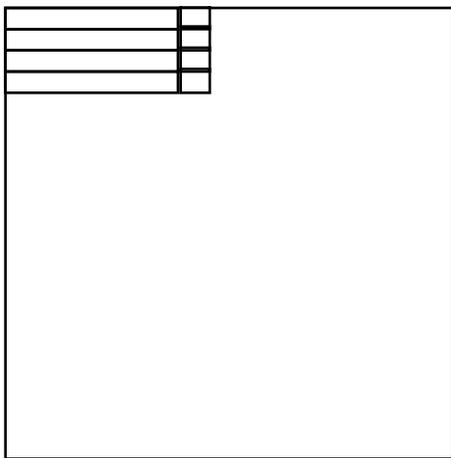
Soit une image $P(i, j)$.

L'integral emage mn est :

$$S_{mn}(i, j) = \sum_{k=1}^m \sum_{l=1}^n P(i+k, j+l)$$

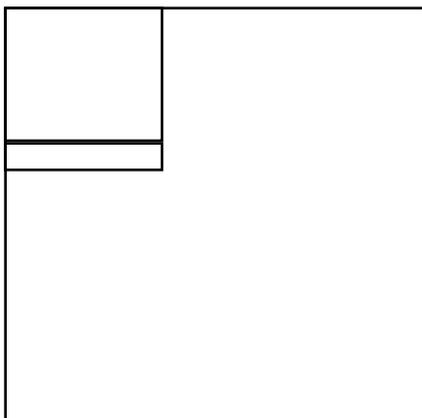
Il existe un algorithme rapide pour le calcul somme des voisinage de tous les tailles.

2) Pour chaque m, calculer $R_m(i, j) = \sum_{k=0}^m P(i+k, j) = R_{m-1}(i, j) + P(i+m, j)$



$$R_m(i, j) = R_{m-1}(i, j) + P(i+m, j)$$

3) pour chaque $S_{mn}(i, j) = \sum_{m=0}^M \sum_{n=0}^N P(i+m, j+n) = S_{n-1}(i, j) + R(i, j+n)$



Les differences des Integrals Images donnees les caracteristiques

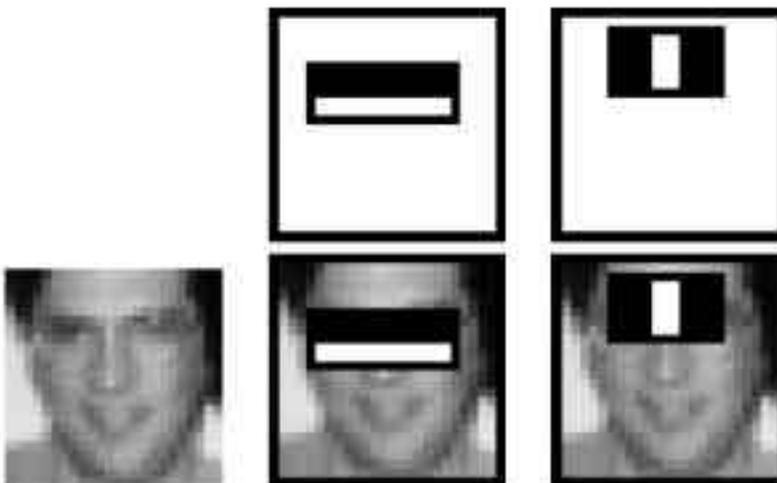
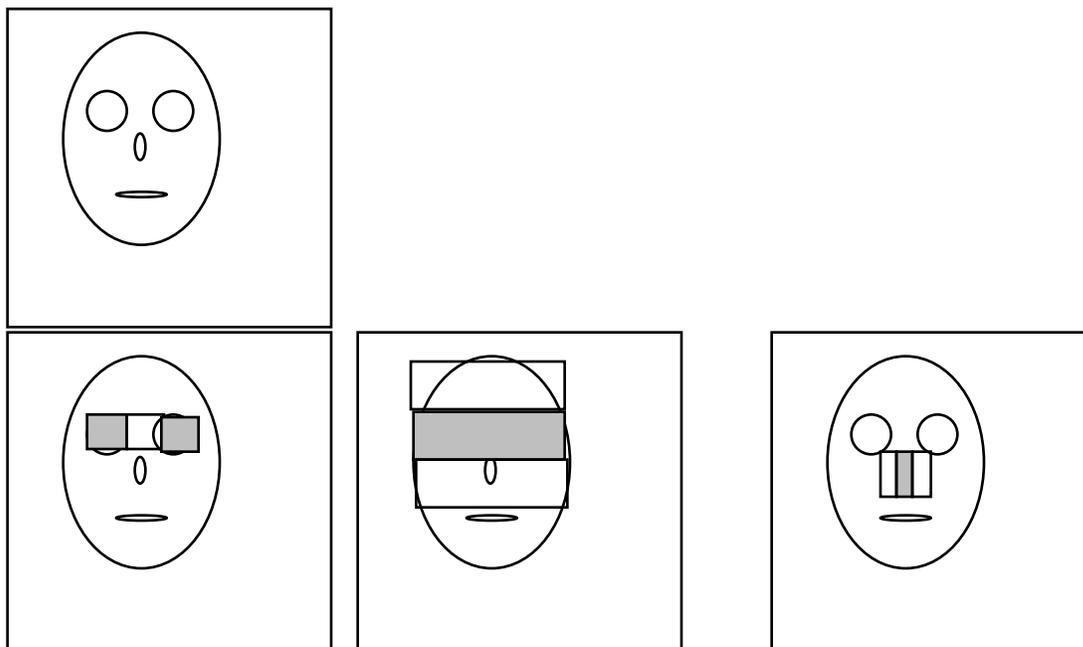
Ceci donne une grande quantité (D) de differences de fenetres uniformes : (les ondelettes de Haar).

Chaque ondelette est une hyperplan dans l'Espace a D dimensions.

On utilise Boosting pour determiner les plans les plus discriminant.

Exemples :

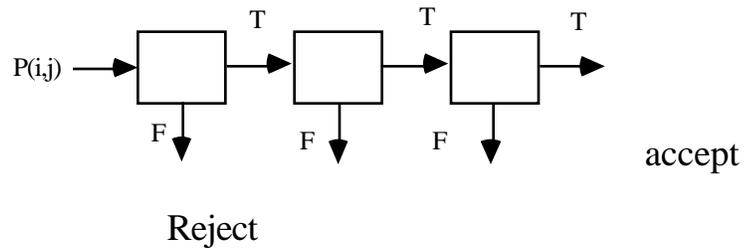
Exemples



Exemple extrait du Viola Jones 2001

Détection par cascade de classifications.

Une cascade est appris par boosting.



Exemple : 38 etages et 6000 caracteristiques

Le detecteur est appliqué à une gamme de positions, i, j , est à une ensemble de échelles, S .