

Systemes Intelligents : Raisonnement et Reconnaissance

James L. Crowley

Deuxième Année ENSIMAG

Deuxième Semestre 2005/2006

Séance 11

19 avril 2006

Discrimination Linéaire

Notations.....	2
Fonctions de Discrimination (rappel).....	3
Rappelle de quelques faits sur les plans et les Hyperplans	4
Estimation par Moindres de Carrées.....	6
Discrimination Linéaire MultiClasse.....	7
Perceptrons	8
Méthodes à Noyaux (Kernel Methods).....	9
Boosting	11
Résumé de l'algorithme de Boosting quand $K = 2$	11
La discriminante linéaire de Fisher	12
Le Discriminant de Fisher avec $K > 2$	15
Classification Linéaire Bayésienne.	18
Le cas des variances blanches ("Matched Filter").....	18
Exemple : Intercorrelation de motifs (NCC).....	20

Sources Bibliographique :

N. Cristianini, J. Shawe-Taylor, "Support Vector Machine and other Kernel based learning methods", Cambridge University Press, 2000.

Notations

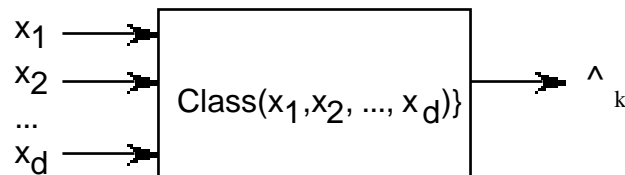
x	Une variable
X	Une valeur aléatoire (non-prévisible).
N	Le nombre de valeurs possible pour x ou X
x	Un vecteur de D variables
X	Un vecteur aléatoire (non-prévisible).
D	Nombre de dimensions de x ou X
T_k	La classe k
k	Indice d'une classe
K	Nombre de classes
(y_m, X_m)	Une exemple d'événement.
y_m	Un 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)

Soit les événements E décrivent par un vecteur de caractéristiques X : (E,X).

Soit K classes d'événements $\{k\} = \{1, 2, \dots, K\}$ avec une classe $k \in \{k\}$

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



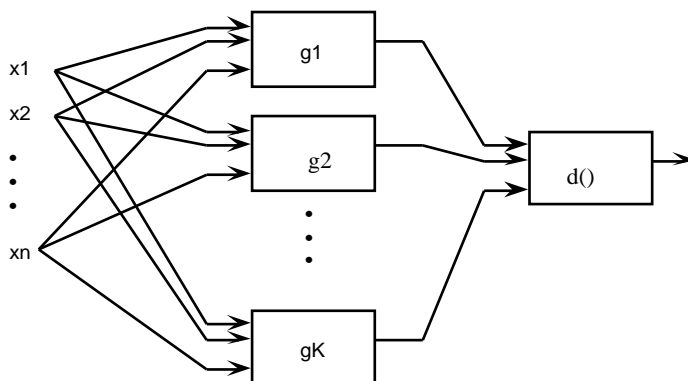
\hat{k} est la proposition que $(E \in k) : \hat{k} = d\{g(X)\}$.

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.



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".

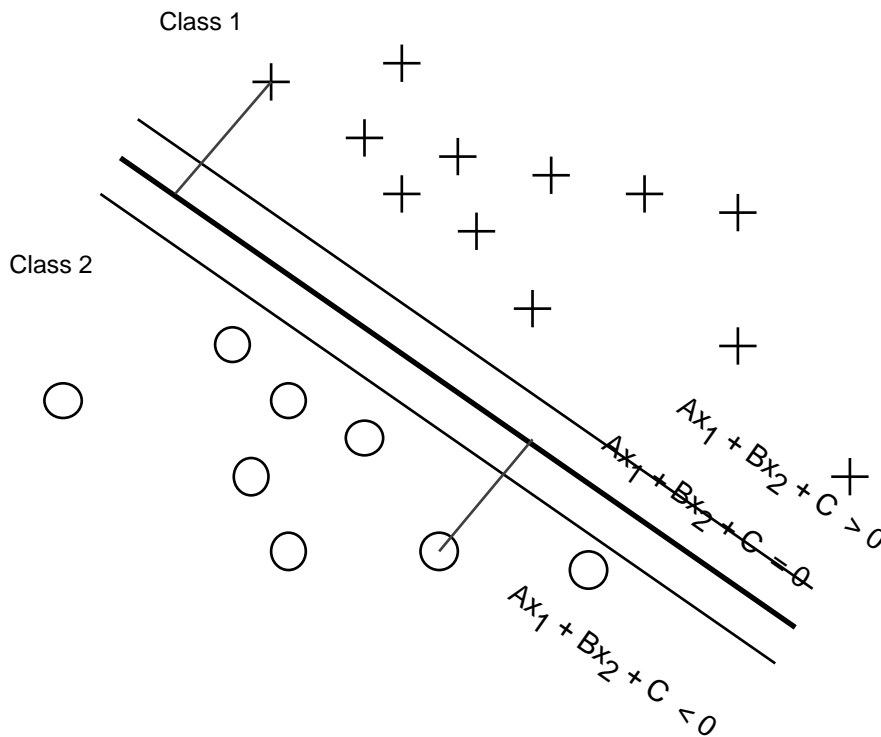
Rappelle de quelques faits sur les plans et les Hyperplans

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 $W = \begin{matrix} w_1 \\ w_2 \\ \dots \\ w_D \end{matrix}$ 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 Carrés

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 $y = g(X) = W X + b = \hat{X}_m \hat{W}$ (en coordonnées Homogènes).

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

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

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

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

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

Le Minimum est trouvé quand :

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

$$\text{Donc : } X Y = X X \hat{W} \text{ et donc } \hat{W} = \begin{matrix} W \\ b \end{matrix} = (X X)^{-1} X Y$$

Discrimination Linéaire MultiClasse

Pour le cas où $K > 2$ On définit un K dimensionnel vecteur d'indication, Y .

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

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

$$y_k \hat{=} \begin{cases} 1 & \text{E} & k \\ 0 & \text{sinon} \end{cases}$$

(Les autres classes ne contribuent pas au hyperplane).

Ensuite $Y = (Y_1, Y_2, \dots, Y_M)$

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

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

Pour un événement X , $\hat{X} = \begin{pmatrix} X \\ 1 \end{pmatrix}$ $\hat{W} = \begin{pmatrix} W \\ b \end{pmatrix}$

$Y = \hat{W}^T \hat{X} = w_1 x_1 + w_2 x_2 + \dots + w_D x_D + b$ est un vecteur de K valeurs de vraisemblance.

On cherche $k = \arg\text{-max}_{y_k} \{ Y = \hat{W}^T \hat{X} \}$

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-séparable".

Le "marge". , 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.

À 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 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 + ym Xm;
      bi+1 = bi + ym R2;
      i = i + 1;
    END IF
  END FOR
UNTIL no mistakes in FOR loop.

```


La marge pour chaque exemple est :

$$m = y_m(W_i X_m + b_i)$$

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

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

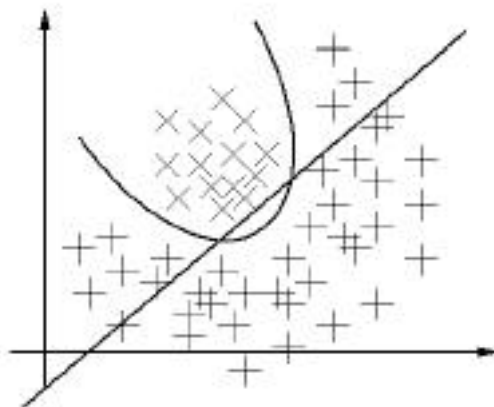
Le "qualité" d'un perceptron est donnée par la distribution des, par exemple, par un histogramme des marges géométrique.

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, 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

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

$$\begin{aligned} g(K(x)) &= g(w) = a w_1 + b w_2 + c w_3 + d w_4 + e w_5 \\ &= a x_1 + b x_2 + c x_1^2 + d x_1x_2 + e x_2^2 \end{aligned}$$

Autres Noyaux populaires :

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

$$K(x) = \ln(x)$$

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 un 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 un ensemble des exemples $\{Y_m, X_m\}$

1) Initialiser un vecteur de M coefficients $a_m = 1$. Initialiser $S = M$.

(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 la matrice \mathbf{X} et un vecteur \mathbf{Y} pour calculer la $N^{\text{ième}}$ fonction de discrimination.

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

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

On trouve une classification, par exemple avec :

$$\hat{\mathbf{W}} = \frac{\mathbf{W}}{b} = (\mathbf{X} \mathbf{X})^{-1} \mathbf{X} \mathbf{Y}$$

3) Pour chaque exemple en $\{X_m\}$, si $d\{g_n(X)\} > k$ alors error :

$$a_m = a_m + 1, S = S + 1$$

4) Répéter étape 2 pour classifier $N+1$.

La discriminante linéaire de Fisher

Dans beaucoup de domaines, il existe une multitude de caractéristiques utilisables pour la reconnaissance. Chaque caractéristique semble à apporter une contribution pour un cas ou dans une autre. Il semble souhaitable de les inclure dans le vecteur x . Mais ceci induit une croissance exponentielle dans les nombres d'exemples nécessaires, M . En Anglais, on appelle ce problème "The Curse of Dimensionality".

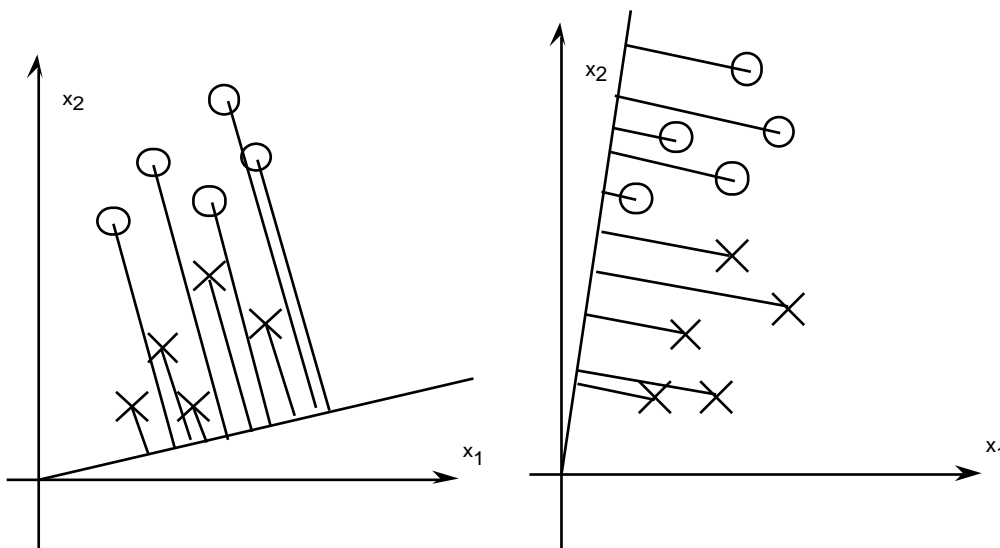
La technique de Fisher permet une réduction dans le nombre de dimensions, D avec une faible augmentation dans la probabilité d'erreur.

Le principe de Fisher est de projeter le vecteur de caractéristique, x de D_x dimensions vers un espace z de D_z par une transformation linéaire F choisit tel quel

- 1) $D_z \ll D_x$ et
- 2) Les exemples des classes A_k sont séparés.

$$z = F^T x$$

En général, s'il y a K classes, nous allons chercher $D_z = K-1$



La discriminabilité des classes dépend de la direction de F

Pour déterminer la meilleure projection, on appuie sur une mesure de la séparation entre classes.

Soit $K=2$ classes T_1 et T_2 représenté par les exemples X_{1m} et X_{2m}

Dans ce cas, $D_Z = 2-1 = 1$. Z est un scalaire

Pour chaque exemple :

$$Z_{km} = F^T X_{km}$$

Nota que F est une projection telle que $\|F\| = 1$

Les moyennes des exemples pour chaque classe sont

$$\mu_k = E\{X_{km}\} = \frac{1}{M_k} \sum_{m=1}^{M_k} X_{km}$$

Les moments sont les invariants affines. Donc, la moyenne (1ere moment) d'une projection est la projection de la moyenne.

$$\beta_k = E\{Z_{km}\} = \frac{1}{M_k} \sum_{m=1}^{M_k} Z_{km} = \frac{1}{M_k} \sum_{m=1}^{M_k} F^T X_{km} = F^T \mu_k$$

La distance entre les classes est $d_{12} = \|\beta_1 - \beta_2\| = \|F^T (\mu_1 - \mu_2)\|$

On veut rendre la distance entre classes aussi grandes que possible, sans disperser les classes.

La dispersion ("Scatter") pour une ensemble $\{X_{km}\}$ d'exemples et pour une classe k est une matrice

$$S_k = M_k C_k = \sum_{m=1}^{M_k} (X_{km} - \mu_k)(X_{km} - \mu_k)^T$$

La transformation F projet le vecteur X sur la scalaire Z . La dispersion ("Scatter") pour la projection des exemples de la classe k est

$$\tilde{S}_k = \sum_{m=1}^{M_k} (Z_{km} - \beta_k)^2$$

Le critère de Fisher est de maximiser le ratio de la séparation des deux classes par rapport à leurs dispersions.

$$J(F) = \frac{(\mu_1 - \mu_2)^2}{\tilde{S}_1 + \tilde{S}_2} = \frac{\|F^T(\mu_1 - \mu_2)\|^2}{\tilde{S}_1 + \tilde{S}_2}$$

Fisher cherche la transformation F^T tel quel

$$F = \arg\text{-max}_F \left\{ \frac{\|F^T(\mu_1 - \mu_2)\|^2}{\tilde{S}_1 + \tilde{S}_2} \right\}$$

Soit $M = \sum_{k=1}^K M_k$ exemples, X_{km} .

Pour $K=2$, $M = M_1 + M_2$

La moyenne de chaque classe est

$$\mu_k = \frac{1}{M_k} \sum_{m=1}^{M_k} X_{km}$$

La moyenne de TOUS les exemples est

$$\mu = \frac{1}{M} \sum_{k=1}^K M_k \mu_k = \frac{1}{M} (M_1 \mu_1 + M_2 \mu_2)$$

La matrice de dispersion inter-classes S_B (B voudrait dire "between") est la dispersion des moyennes des classes.

$$S_B = (\mu_1 - \mu)(\mu_1 - \mu)^T + (\mu_2 - \mu)(\mu_2 - \mu)^T$$

La dispersion intra-classe S_W (En Anglais W pour "within") est la covariance moyenne.

$$S_W = \sum_{k=1}^K S_k = \sum_{k=1}^K \frac{1}{M_k} \sum_{m=1}^{M_k} (X_{km} - \mu_k)(X_{km} - \mu_k)^T = S_1 + S_2$$

La meilleure transformation F est celle que

$$F = \operatorname{argmax}_F \left\{ \frac{\|F^T S_B F\|}{\|F^T S_w F\|} \right\}$$

Dans notre exemple avec $K=2$, $D_z = 1$ (donc $F^T S_B F$ est un scalaire)

On définit :

$$J(F) = \frac{F^T S_B F}{F^T S_w F}$$

en physique, ceci est connu comme le quotient de Rayleigh. Il est possible de montrer que

$$S_B F = S_w F.$$

donc

$$S_w^{-1} S_B F = F.$$

Le facteur d'échelle n'est pas important est-on peut déterminer directement

$$\text{Donc } F = S_w^{-1} S_B = S_w^{-1}(\mu_1 - \mu_2)$$

Ceci est le discriminant linéaire de Fisher pour deux classes.

Il maximise la dispersion entre les classes.

On rappelle que la surface de décision linéaire entre deux classes a la forme :

$$F^T X + b_0 = 0 \quad \text{où } F = C^{-1}(\mu_1 - \mu_2)$$

et b_0 est un constant

Le Discriminant de Fisher avec $K > 2$

Comment généraliser en $D_z = K-1$ dimensions?

Pour le cas de K classes, la généralisation naturelle est avec $K-1$ fonctions de Fisher.

Il est supposé que $D = K$.

$$S_w = \sum_{k=1}^K S_k$$

$$\text{où } S_k = \sum_{m=1}^{M_k} (X_{km} - \mu_k)(X_{km} - \mu_k)^T$$

$$\text{et } \mu_k = \frac{1}{M_k} \sum_{m=1}^{M_k} X_{km}$$

On peut définir une moyenne globale, μ et une matrice de dispersion totale S_T comme :

$$\mu = \frac{1}{M} \sum_{m=1}^M X_m = \frac{1}{M} \sum_{k=1}^K M_k \mu_k$$

$$\text{et } S_T = \sum_{k=1}^K \sum_{m=1}^{M_k} (X_{km} - \mu)(X_{km} - \mu)^T$$

S_T est la dispersion "totale".

On peut définir la dispersion entre classe ("between class") comme

$$S_B = \sum_{k=1}^K M_k (\mu_k - \mu)(\mu_k - \mu)^T$$

Il est possible de démontrer que

$$S_T = S_w + S_B$$

Pour chaque classe, k , on obtient une transformation F_k dans la forme d'un vecteur à D dimensions.

$$z_k = F_k^T X$$

Si on aligne les transformations dans une matrice de taille $D \times (K-1)$ on a

$$Z = F^T X$$

Ou Z est un vecteur de $K-1$ coefficients.

Par invariance des moments, on peut montrer que

$$\begin{aligned}\tilde{S}_W &= F^T S_W F \\ \tilde{S}_B &= F^T S_B F\end{aligned}$$

Le Critère de Fisher est de maximiser

$$J(W) = \frac{\det(\tilde{S}_B)}{\det(\tilde{S}_W)} = \frac{\det(F^T S_B F)}{\det(F^T S_W F)}$$

La solution est rendue par une analyse en composant principales de \tilde{S}_B

$$\tilde{S}_B F_k = \lambda_k \tilde{S}_W F_k \quad \text{Résoudre } F_k \text{ telle que } (\tilde{S}_B - \lambda_k \tilde{S}_W) F_k = 0$$

Les colons peuvent être calculé par une simple orthogonalisation par l'algorithme de Gram-Schmidt des vecteurs

$$(\mu_k - \mu) \quad \text{pour } k = 1, \dots, K-1.$$

On note que F n'est pas unique. Il existe une classe d'équivalence avec les rotations et multiplications par une constant.

Classification Linéaire Bayésienne.

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

Si $\sigma_{i,j} = \sigma^2 \delta_{ij}$ dans ce cas.

$$\det(C) = (\sigma^2)^n \text{ et } C_k^{-1} = \frac{1}{\sigma^2} I$$

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

$$g_k(x) = -\frac{\|x - \mu_k\|^2}{2\sigma^2} + \log\{p(w_k)\}$$

Ce cas arrive quand les observations sont corrompues par un bruit additif blanc indépendant des classes et d'une puissance égale 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 ?

Si les caractéristiques suivent une densité Normale :

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

La fonction de discrimination devient est une fonction quadratic

$$g_k(x) = -\frac{1}{2} \log\{\det(C)\} - \frac{1}{2}(x - \mu_k)^T C^{-1}(x - \mu_k) + \log\{p(k)\}$$

On peut réécrire $g_k(x)$ comme

$$g_k(X) = X^T \mathbf{B} X + b_k^T X + b_{k0}.$$

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

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

et le constant est $b_{k0} = -\frac{1}{2}(\mu_k^T \mathbf{C}^{-1} \mu_k) + \text{Log}\{p(k)\}$

Mais, si tous les messages ont la même probabilité :

$$b_{k0} = -\frac{1}{2}(\mu_k^T \mu_k) = -\frac{1}{2} \|\mu_k\|^2$$

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

La problématique de la communication est formalisée comme :

Comment chercher un message "M(t)" dans un signal S(t)?

Pour la communication hertzienne, les caractéristiques sont les canaux de bande passante, et le bruit est blanc. Il est indépendant du signal.

Pour chaque message T_k , on prend une moyenne d'échantillon comme un "prototype". On normalise l'énergie (L'énergie ne dépend pas du message) pour former un prototype.

$$b_k = b(t) = \frac{\mu_k}{\|\mu_k\|}$$

Ensuite, à chaque t_0 pour chaque k , on cherche $\int b(t) S(t-t_0) dt < \text{seuil}$. Si oui,

$$k = \arg\text{-max}_k \left\{ \int b(t) S(t-t_0) dt \right\}$$

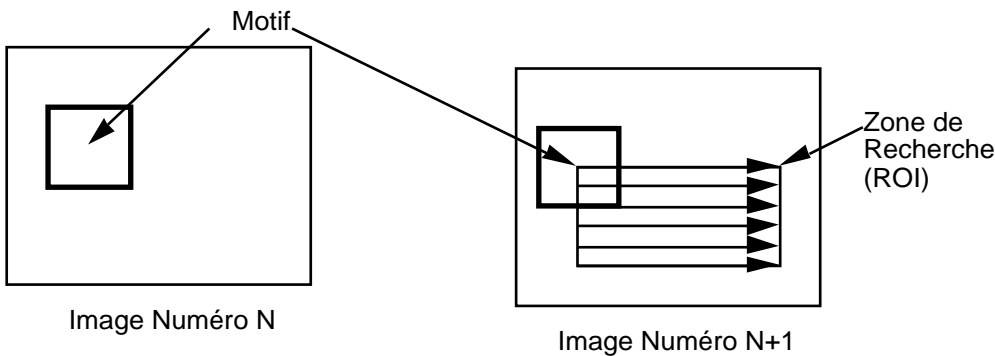
Exemple : Intercorrelation de motifs (NCC).

(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 : $d \quad dd^2 = 2$

les variances des pixels sont indépendantes : $i, j \quad i \quad j \quad ij^2 = 0$ Donc $C_k = 2 I$

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.

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

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

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.

