

# Análise de Séries Temporais

## VII. Métodos estatísticos no domínio da frequência

Fernando Lucambio

Departamento de Estatística  
Universidade Federal do Paraná

Julho, 2023

O problema da discriminação entre explosões de mineração e terremotos é um representante razoável para o problema de discriminação.

Esse problema é de fundamental importância para o monitoramento de um tratado abrangente de proibição de testes.

Os problemas de classificação de séries temporais não se restringem a aplicações geofísicas, mas ocorrem sob muitas e variadas circunstâncias em outros campos.

Tradicionalmente, a detecção de um sinal incorporado em uma série de ruídos é analisada na literatura de engenharia por técnicas estatísticas de reconhecimento de padrões.

As abordagens históricas para o problema da discriminação entre diferentes classes de séries temporais podem ser divididas em duas categorias distintas.

A abordagem de otimização, como encontrada na literatura de engenharia e estatística, faz suposições Gaussianas específicas sobre as funções de densidade dos grupos separados e, em seguida, desenvolve soluções que satisfazem critérios de erros mínimos bem definidos.

Uma segunda classe de técnicas, que pode ser descrita como uma abordagem de extração de recursos, prossegue-se de maneira mais heurística, observando quantidades que tendem a ser bons discriminadores visuais para populações bem separadas e que têm alguma base na teoria ou intuição física.

Para séries temporais multivariadas mais longas que possam ser consideradas estacionárias após a subtração da média comum, a abordagem do domínio da frequência será mais fácil computacionalmente porque o vetor  $n \times p$  dimensional será reduzido para separar os cálculos feitos nas transformadas discretas de Fourier (DFT's)  $p$ -dimensionais.

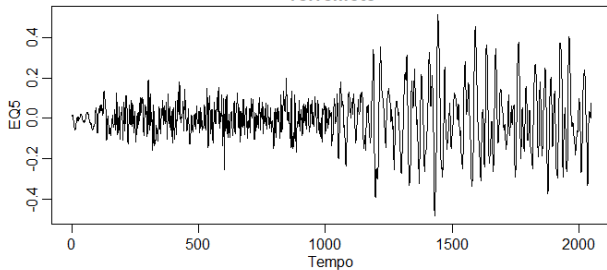
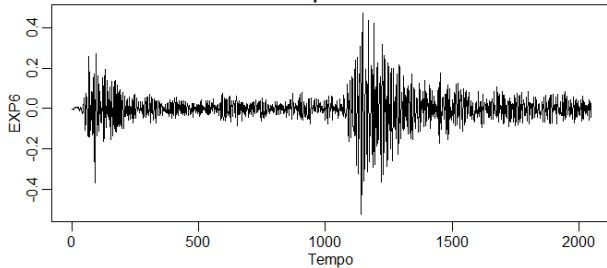
Isso acontece devido à independência aproximada das transformadas discretas de Fourier,  $X(\omega_k)$ ,  $0 \leq \omega \leq 1$ , uma propriedade que usamos frequentemente anteriormente.

Por fim, as propriedades de agrupamento de medidas como informações sobre discriminação e estatísticas baseadas na verossimilhança podem ser usadas para desenvolver medidas de disparidade para agrupar séries temporais multivariadas.

O problema geral de classificação de uma série temporal vetorial ocorre da seguinte maneira. Observamos uma série temporal  $X$  conhecida por pertencer a uma das  $g$  populações, denotadas por  $\Pi_1, \Pi_2, \dots, \Pi_g$ .

O problema geral é atribuir ou classificar essa observação em um dos  $g$  grupos da melhor maneira possível.

Um exemplo pode ser as  $g = 2$  populações de terremotos e explosões. Gostaríamos de classificar o evento desconhecido, mostrado no gráfico a seguir, como pertencente às populações de terremotos  $\Pi_1$  ou explosões  $\Pi_2$ . Exemplo I.7 Terremotos e explosões.

**Terremoto****Explosão**

Para resolver esse problema, precisamos de um critério de otimização que leve a uma estatística  $T(X)$  que possa ser usada para atribuir o evento às populações de terremotos ou explosões.

Para medir o sucesso da classificação, precisamos avaliar os erros que podem ser esperados no futuro em relação ao número de terremotos classificados como explosões (alarmes falsos) e o número de explosões classificadas como terremotos (sinais perdidos).

O problema pode ser formulado assumindo que a série observada  $X$  tem uma função de densidade  $p_i(x)$  quando a série observada é da população  $\Pi_i$  para  $i = 1, \dots, g$ . Em seguida, particione o espaço abrangido pelo processo  $n \times p$ -dimensional  $x$  em  $g$  regiões mutuamente exclusivas  $R_1, R_2, \dots, R_g$  tal que, se  $x$  cair em  $R_i$ , atribuímos  $x$  à população  $\Pi_i$ .

À probabilidade de classificação incorreta é definida como a probabilidade de classificar a observação na população  $\Pi_j$  quando ela pertence a  $\Pi_i$ , para  $j \neq i$  e seria dada pela expressão

$$P(j|i) = \int_{R_j} p_i(x) dx \cdot$$

A probabilidade do erro total depende também das probabilidades a priori, por exemplo,  $\pi_1, \pi_2, \dots, \pi_g$  de  $x$  pertencer a um dos  $g$  grupos. Por exemplo, a probabilidade de uma observação  $x$  se originar de  $\Pi_i$  e depois ser classificada em  $\Pi_j$  é obviamente  $\pi_i P(j|i)$  e a probabilidade total de erro se torna

$$P_\epsilon = \sum_{i=1}^g \pi_i \sum_{j \neq i} P(j|i) \cdot$$



Embora os custos não tenham sido incorporados na expressão acima, é possível fazê-lo multiplicando  $P(j | i)$  por  $C(j | i)$ , o custo de atribuir uma série da população  $\Pi_i$  a  $\Pi_j$ .

O erro geral  $P_\epsilon$  é minimizado classificando  $x$  em  $\Pi_i$  se

$$\frac{p_i(x)}{p_j(x)} > \frac{\pi_j}{\pi_i},$$

para todos os  $j \neq i$  (ver, por exemplo, Anderson, 1984).

Uma quantidade de interesse, da perspectiva Bayesiana, é a probabilidade a posteriori de uma observação pertencer à população  $\Pi_i$ , condicionada à observação de  $x$ , digamos

$$P(\Pi_i | x) = \frac{\pi_i p_i(x)}{\sum_j \pi_j p_j(x)}.$$

O procedimento que classifica  $x$  na população  $\Pi_j$  para a qual a probabilidade a posteriori é maior é equivalente ao critério implícito anterior. As probabilidades a posteriori dão uma idéia intuitiva plausível das chances relativas de pertencer a cada uma das populações.

Muitas situações ocorrem, como na classificação de terremotos e explosões, nas quais existem apenas  $g = 2$  populações de interesse. Para duas populações, o Lema de Neyman – Pearson implica, na ausência de probabilidades a priori, classificar uma observação em  $\Pi_1$  quando

$$\frac{p_1(x)}{p_2(x)} > K,$$

minimiza cada uma das probabilidades de erro para um valor fixo do outro.

A regra é idêntica à regra de Bayes do erro geral  $P_\epsilon$ , o qual é minimizado classificando  $x$  em  $\Pi_i$  quando  $p_i(x)/p_j(x) > \pi_j/\pi_i$ , para todos os  $j \neq i$ . A regra acima é idêntica à regra de Bayes quando  $K = \pi_2/\pi_1$ .

A teoria apresentada acima assume uma forma simples quando o vetor  $x$  tem uma distribuição normal  $p$ -variada com vetores de média  $\mu_j$  e matrizes de covariância  $\Sigma_j$  sob  $\Pi_j$  para  $j = 1, 2, \dots, g$ .

Nesse caso, basta usar

$$p_j(x) = (2\pi)^{-p/2} |\Sigma_j|^{-1/2} \exp \left\{ -\frac{1}{2} (x - \mu_j)^\top \Sigma_j^{-1} (x - \mu_j) \right\}.$$

As funções de classificação são expressas por quantidades proporcionais aos logaritmos das densidades, por exemplo,

$$g_j(x) = -\frac{1}{2} \ln(|\Sigma_j|) - \frac{1}{2} x^\top \Sigma_j^{-1} x + \mu_j^\top \Sigma_j^{-1} x - \frac{1}{2} \mu_j^\top \Sigma_j^{-1} \mu_j + \ln(\pi_j).$$

Nas expressões que envolvem a log-verossimilhança, em geral ignoraremos os termos que envolvem a constante  $-2 \ln(2\pi)$ . Podemos então atribuir uma observação  $x$  à população  $\Pi_i$  sempre que

$$g_i(x) > g_j(x)$$

para  $i \neq j$ ,  $i = 1, 2, \dots, g$  e a probabilidade a posteriori tem a forma

$$p(\Pi_i | x) = \frac{\exp(g_i(x))}{\sum_j \exp(g_j(x))}.$$

Uma situação comum que ocorre nas aplicações é a classificação sob a premissa de normalidade multivariada e matrizes de covariância iguais; ou seja,  $\Sigma_1 = \Sigma_2 = \Sigma$ .

Então, o critério de atribuição de uma observação  $x$  à população  $\Pi_i$  pode ser expresso em termos da função discriminante linear como

$$d_i(x) = g_1(x) - g_2(x) = (\mu_1 - \mu_2)^\top \Sigma^{-1} x - \frac{1}{2} (\mu_1 - \mu_2)^\top \Sigma^{-1} (\mu_1 + \mu_2) + \ln(\pi_1/\pi_2),$$

onde classificamos em  $\Pi_1$  ou  $\Pi_2$  de acordo com  $d_i(x) \geq 0$  ou  $d_i(x) < 0$ .

A função discriminante linear é claramente uma combinação de variáveis normais e, para o caso  $\pi_1 = \pi_2 = 0.5$ , terá média  $D^2/2$  sob  $\Pi_1$  e média  $-D^2/2$  sob  $\Pi_2$ , com variâncias dadas por  $D^2$  nas duas hipóteses, em que

$$D^2 = (\mu_1 - \mu_2)^\top \Sigma^{-1} (\mu_1 - \mu_2)$$

é a distância de Mahalanobis entre os vetores médios  $\mu_1$  e  $\mu_2$ . Nesse caso, as duas probabilidades de classificação incorreta são

$$P(1|2) = P(2|1) = \Phi(-D/2),$$

e o desempenho está diretamente relacionado à distância de Mahalanobis.

Para o caso de matrizes de covariâncias diferentes, a função discriminante assume uma forma diferente, com a diferença  $g_1(x) - g_2(x)$  na forma

$$d_q(x) = -\frac{1}{2} \ln \left( \frac{|\Sigma_1|}{|\Sigma_2|} \right) - \frac{1}{2} x^\top (\Sigma_1^{-1} - \Sigma_2^{-1}) x + (\mu_1^\top \Sigma_1^{-1} - \mu_2^\top \Sigma_2^{-1}) x + \ln \left( \frac{\pi_1}{\pi_2} \right),$$

para  $g = 2$  grupos.

Essa função discriminante difere do caso de covariâncias iguais no termo linear e em um termo quadrático não linear envolvendo as diferentes matrizes de covariância.

A distribuição teórica não é tratável para o caso quadrático, portanto, nenhuma expressão conveniente está disponível para as probabilidades de erro da função discriminante quadrática.

Uma dificuldade em aplicar a teoria acima a dados reais é que os vetores médios  $\mu_j$  e as matrizes de covariância  $\Sigma_j$  raramente são conhecidos.

Alguns problemas de engenharia, como a detecção de um sinal no ruído branco, assumem que os vetores médios e os parâmetros de covariância são conhecidos exatamente e isso pode levar a uma solução ótima.



Na situação multivariada clássica, é possível coletar uma amostra de  $N_i$  vetores de treinamento do grupo  $\Pi_i$ , por exemplo,  $x_{i,j}$ , para  $j = 1, \dots, N_i$  e usá-los para estimar os vetores médios e as matrizes de covariância para cada um dos grupos  $i = 1, 2, \dots, g$ ; ou seja, basta escolher  $\bar{x}_i$  e

$$S_i = \frac{1}{N_i - 1} \sum_{j=1}^{N_i} (x_{i,j} - \bar{x}_i)(x_{i,j} - \bar{x}_i)^\top$$

como estimadores para  $\mu_i$  e  $\Sigma_i$ , respectivamente. Caso as matrizes de covariância sejam consideradas iguais, basta usar o estimador agrupado

$$S = \frac{1}{\sum_i N_i - g} \sum_i (N_i - 1) S_i.$$

Para o caso de uma função discriminante linear, podemos usar

$$\hat{g}_i(x) = \bar{x}_i^\top S^{-1}x - \frac{1}{2}\bar{x}_i^\top S^{-1}\bar{x}_i + \log(\pi_i),$$

como um estimador simples para  $g_i(x)$ .

Para amostras grandes,  $\bar{x}_i$  e  $S$  convergem para  $\mu_i$  e  $\Sigma$  em probabilidade, de modo que  $\hat{g}_i(x)$  converge em distribuição para  $g_i(x)$  nesse caso.

O procedimento funciona razoavelmente bem para o caso em que cada  $N_i$ ,  $i = 1, \dots, g$  são grandes, em relação ao comprimento da série  $n$ , um caso relativamente raro na análise de séries temporais. Por esse motivo, recorreremos ao uso de aproximações espectrais para o caso em que os dados são fornecidos como séries temporais longas.

O desempenho das funções discriminantes amostrais podem ser avaliadas de várias maneiras diferentes. Se os parâmetros populacionais forem conhecidos, a distância de Mahalanobis e as probabilidades de classificação incorreta poderão ser avaliadas diretamente. Se os parâmetros forem estimados, a distância estimada de Mahalanobis  $\hat{D}^2$  pode ser substituída pelo valor teórico em amostras muito grandes.

Outra abordagem é calcular as taxas de erro aparentes usando o resultado da aplicação do procedimento de classificação às amostras de treinamento. Se  $n_{i,j}$  denota o número de observações da população  $\Pi_j$  classificadas em  $\Pi_i$ , as taxas de erro amostral podem ser estimadas pela razão

$$\hat{P}(i|j) = \frac{n_{i,j}}{\sum_i n_{i,j}}, \quad i \neq j.$$

Se as amostras de treinamento não forem grandes, esse procedimento pode ser tendencioso e uma opção de reamostragem, como validação cruzada ou o bootstrap, pode ser empregada.

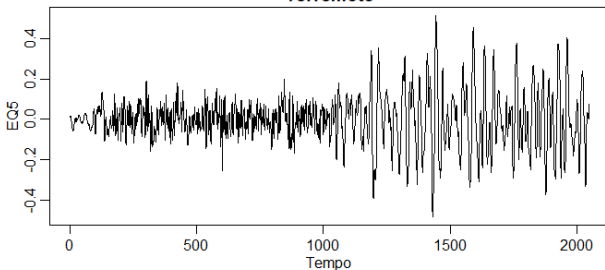
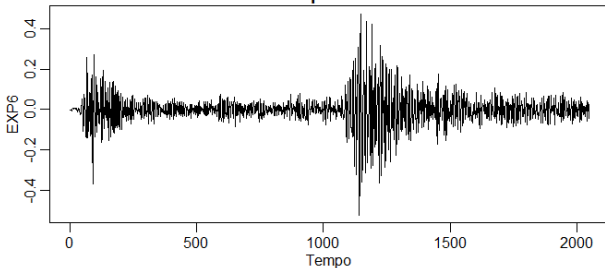
Uma versão simples da validação cruzada é o jackknife proposto por Lachenbruch and Mickey (1968), que sustenta a observação a ser classificada, derivando a função de classificação das demais observações.

A repetição deste procedimento para cada um dos membros da amostra de treinamento e o cálculo da taxa de erro amostral acima para as amostras de validação levam a melhores estimadores das taxas de erro.

## Exemplo VII.10. Análise discriminante usando amplitudes.

Os logaritmos (base 10) das amplitudes máximas pico a pico dos componentes  $P$  e  $S$ , denotados por  $\log_{10}(P)$  e  $\log_{10}(S)$ , podem ser considerados vetores de característica bidimensionais, por exemplo,  $x = (x_1, x_2)^T = (\log_{10}(P), \log_{10}(S))^T$ , de uma população normal bivariada com médias e covariâncias diferentes.

Os dados originais, de Kakizawa et al. (1998), são mostrados na figura anterior. Os dados incluem o evento Novaya Zemlya (NZ) de origem desconhecida. A tendência dos terremotos de terem valores mais altos para  $\log_{10}(S)$ , em relação ao  $\log_{10}(P)$ , foi notada por muitos e o uso do logaritmo da razão, ou seja,  $\log_{10}(P) - \log_{10}(S)$  em algumas referências, por exemplo, Lay (1997) é um indicador tácito de que uma função linear dos dois parâmetros será um discriminante útil.

**Terremoto****Explosão**

As médias amostrais foram

$$\bar{x}_1 = (0.3477384, 1.0244672)$$

e

$$\bar{x}_2 = (0.9222803, 0.9930151)$$

e as matrizes de covariância amostrais foram

$$S_1 = \begin{pmatrix} 0.025907196 & -0.007064702 \\ -0.007064702 & 0.010202513 \end{pmatrix}$$

e

$$S_2 = \begin{pmatrix} 0.025197761 & -0.000760573 \\ -0.000760573 & 0.010342222 \end{pmatrix}$$

com a matriz de covariância agrupada dada por

$$S = \begin{pmatrix} 0.025552479 & -0.003912638 \\ -0.003912638 & 0.010272368 \end{pmatrix}.$$

Embora as matrizes de covariância não sejam iguais, tentamos a função discriminante linear que produz, com probabilidades a priori iguais  $\pi_1 = \pi_2 = 0.5$ , as funções discriminantes amostrais

$$\hat{g}_1(x) = 29.68321x_1 + 106.1547x_2 - 59.53698$$

e

$$\hat{g}_2(x) = 51.95275x_1 + 112.8352x_2 - 80.5039$$

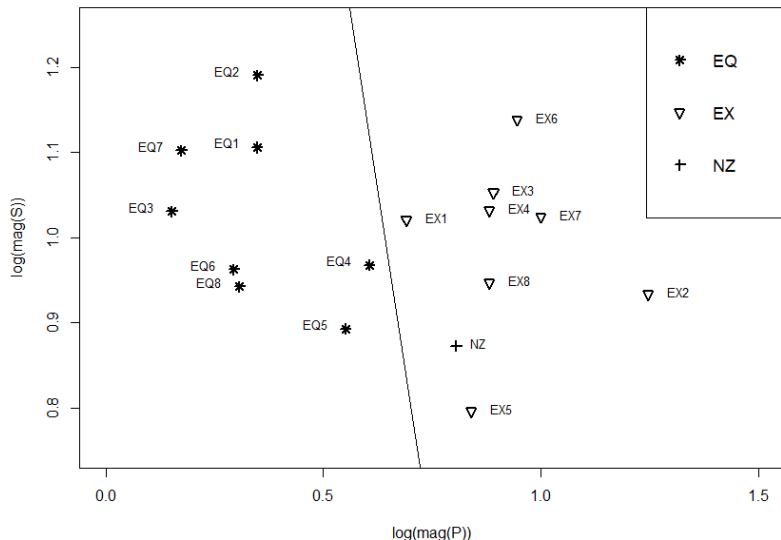
com a função discriminante linear estimada

$$\hat{g}_2(x) = -23.3795x_1 - 5.843191x_2 + 20.74047.$$

As probabilidades a posteriori jackknife de ser um terremoto para o grupo de terremotos variavam de 0.621 a 1.000, enquanto as probabilidades de explosão para o grupo de explosão variavam de 0.717 a 1.000. O evento desconhecido, NZ, foi classificado como explosão, com probabilidade a posteriori 0.960.



### Classificação baseada na magnitude



A abordagem de extração de recursos geralmente funciona bem para discriminar entre classes de séries univariadas ou multivariadas quando existe um vetor de baixa dimensão simples que parece capturar a essência das diferenças entre as classes.

Ainda parece sensato, no entanto, desenvolver métodos ótimos de classificação que explorem as diferenças entre as médias multivariadas e as matrizes de covariância no caso de séries temporais.

Tais métodos podem basear-se na aproximação de Whittle à log-verossimilhança.

Nesse caso, assume-se que os DFTs vetoriais,  $X(\omega_k)$ , sejam aproximadamente normais, com médias  $M_j(\omega_k)$  e matrizes espectrais  $f_j(\omega_k)$  para a população  $\Pi_j$  nas frequências  $\omega_k = k/n$ , para  $k = 0, 1, \dots, [n/2]$  e são aproximadamente não correlacionados em diferentes frequências, digamos,  $\omega_k$  e  $\omega_l$ , para  $k \neq l$ .

Em seguida, escrevendo as densidades normais complexas, leva ao critério,

$$g_j(X) = \ln(\pi_j) - \sum_{0 < \omega_k < 1/2} \left( \ln(|f_j(\omega_k)|) + X^*(\omega_k) f_j^{-1}(\omega_k) X(\omega_k) - 2M_j^*(\omega_k) f_j^{-1}(\omega_k) X(\omega_k) + M_j^*(\omega_k) f_j^{-1}(\omega_k) M_j(\omega_k) \right),$$

onde a soma ultrapassa as frequências nas quais  $|f_j(\omega_k)| \neq 0$ .

A regra de classificação é como anteriormente mencionado.

No caso de séries temporais, é mais provável que a análise discriminante envolva assumir que as matrizes de covariâncias são diferentes e as médias são iguais. Por exemplo, os testes, mostrados na figura do exemplo anterior, implicam que, para os terremotos e explosões, as diferenças primárias estão nas matrizes espectrais bivariadas e as médias são essencialmente as mesmas. Para este caso, será conveniente escrever a aproximação de Whittle para o logaritmo da verossimilhança na forma

$$\ln(p_j(\mathbf{X})) = \sum_{0 < \omega_k < 1/2} \left( -\ln(|f_j(\omega_k)|) - \mathbf{X}^*(\omega_k) f_j^{-1}(\omega_k) \mathbf{X}(\omega_k) \right),$$

onde omitimos as probabilidades a priori da equação.

O detector quadrático, neste caso, pode ser escrito na forma

$$\ln(p_j(X)) = \sum_{0 < \omega_k < 1/2} \left( -\ln(|f_j(\omega_k)|) - \text{tr}(I(\omega_k)f_j^{-1}(\omega_k)) \right),$$

onde

$$I(\omega_k) = X(\omega_k)X^*(\omega_k)$$

denota a matriz do periodograma. Para probabilidades a priori iguais, podemos atribuir uma observação  $x$  na população  $\Pi_j$  sempre que

$$\ln(p_j(x)) > \ln(p_i(x)),$$

para  $j \neq i, j = 1, 2, \dots, g$ .

Numerosos autores consideraram várias versões de análise discriminante no domínio da frequência. Por exemplo, Shumway and Unger (1974), Alagón (1989), Dargahi-Noubary and Laycock (1981), Taniguchi et al. (1994) e Shumway (1982).

## Medidas de discrepância

Antes de prosseguir com os exemplos de análise discriminante e de agrupamento será útil considerar a relação com a informação de discriminação de Kullback-Leibler (K-L).

Usando a aproximação espectral e observando que a matriz do periodograma tem esperança aproximada

$$E_j(I(\omega_k)) = f_j(\omega_k),$$

podemos escrever as informações aproximadas de discriminação como

$$I(f_1; f_2) = \frac{1}{n} E_1 \left( \ln \left( \frac{p_1(X)}{p_2(X)} \right) \right).$$

A classificação, neste caso, prossegue simplesmente escolhendo a população  $\Pi_j$  que minimiza  $I(\hat{f}; f_j)$ , ou seja, atribuindo  $x$  à população  $\Pi_i$  sempre que

$$I(\hat{f}; f_i) < I(\hat{f}; f_j),$$

para  $j \neq i, j = 1, 2, \dots, g$ . Kakizawa et al. (1998) propuseram o uso da medida de informação de Chernoff (CH) (Chernoff, 1952, Renyi, 1961), definida como

$$B_\alpha(1; 2) = -\ln \left( E_2 \left( \frac{p_2(x)}{p_1(x)} \right)^\alpha \right),$$

onde a medida é indexada por um parâmetro de regularização  $\alpha$ , para  $0 < \alpha < 1$ . Quando  $\alpha = 0.5$ , a medida de Chernoff é a divergência simétrica proposta por Bhattacharya (1943). Para o caso multivariado,

$$B_\alpha(1; 2 : x) = -\ln \left( E_2 \left( \frac{p_2(x)}{p_1(x)} \right)^\alpha \right).$$

## Exemplo VII.11. Análise discriminante em dados sísmicos.

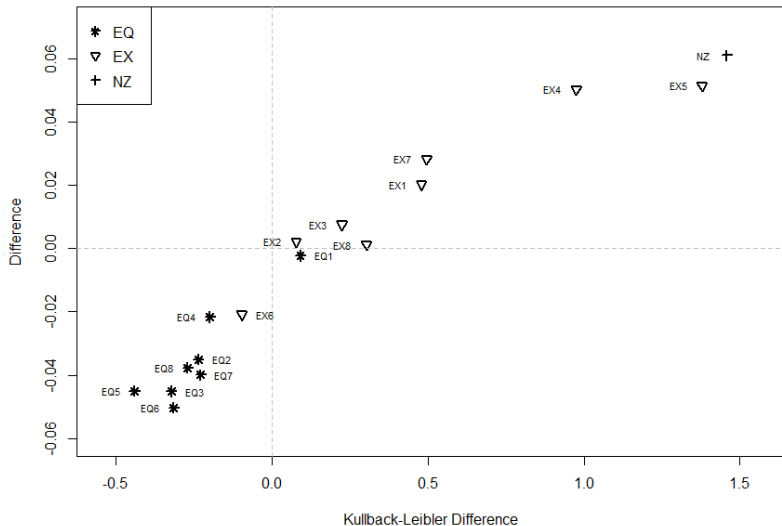
Um esforço considerável foi despendido no uso de várias razões espectrais envolvendo as fases bivariadas P e S como características de discriminação. Kakizawa et al. (1998) mencionam uma série de medidas que têm sido usadas na literatura sísmológica como recursos. Esses recursos incluem relações de potência para as duas fases e relações de componentes de potência em bandas de alta e baixa frequência.

O uso de tais características do espectro sugere que um procedimento ótimo baseado na discriminação entre as matrizes espectrais de dois processos estacionários seria razoável. O fato de que a hipótese de que as matrizes espectrais eram iguais, também foi amplamente rejeitada sugere o uso de uma função discriminante baseada em diferenças espectrais.





### Classification Based on Chernoff and K-L Distances



Para fins de agrupamento, pode ser mais útil considerar uma medida de disparidade simétrica e apresentamos a medida J-Divergência

$$J(f_1; f_2) = I(f_1; f_2) + I(f_2; f_1)$$

e o número de Chernoff simétrico

$$JB_\alpha(f_1; f_2) = B_\alpha(f_1; f_2) + B_\alpha(f_2; f_1)$$

para esse propósito.

Neste caso, definimos a disparidade entre a matriz espectral amostral de um único vetor  $x$ , e a população  $\Pi_j$  como

$$J(\hat{f}; f_j) = I(\hat{f}; f_j) + I(f_j; \hat{f})$$

e

$$JB_\alpha(\hat{f}; f_j) = B_\alpha(\hat{f}; f_j) + B_\alpha(f_j; \hat{f}),$$

respectivamente e use-as como quase distâncias entre o vetor e a população  $\Pi_j$ .

As medidas de disparidade podem ser usadas para agrupar séries temporais multivariadas. As medidas simétricas de disparidade, conforme definidas acima, garantem que a disparidade entre  $f_i$  e  $f_j$  seja igual à disparidade entre  $f_j$  e  $f_i$ .

Portanto, consideraremos as formas simétricas mostradas acima como quase distâncias com o propósito de definir uma matriz de distância para entrada em um dos procedimentos de agrupamento padrão (ver Johnson and Wichern, 1992).

Em geral, podemos considerar métodos de agrupamento hierárquico ou particionado usando a matriz de quase distância como entrada.

Para fins de ilustração, podemos usar a divergência simétrica  $J(\hat{f}; f_j)$ , que implica que a quase-distância entre as séries amostrais com matrizes espectrais estimadas  $\hat{f}_i$  e  $\hat{f}_j$  seria  $J(\hat{f}_i; \hat{f}_j)$ ; ou seja,

$$J(\hat{f}_i; \hat{f}_j) = \frac{1}{n} \sum_{0 < \omega_k < 1/2} \left( \text{tr} \left( \hat{f}_i(\omega_k) \hat{f}_j^{-1}(\omega_k) \right) + \text{tr} \left( \hat{f}_j(\omega_k) \hat{f}_i^{-1}(\omega_k) \right) - 2p \right)$$

para  $i \neq j$ .

Também podemos usar a forma comparável para a divergência de Chernoff, mas podemos não querer fazer uma suposição para o parâmetro de regularização  $\alpha$ .

No agrupamento hierárquico, começamos agrupando os dois elementos da população que minimizam a medida de disparidade  $J(\hat{f}_i; \hat{f}_j)$ . Então, esses dois itens formam um cluster e podemos calcular distâncias entre itens não agrupados como antes.

A distância entre elementos não agrupados e um agrupamento atual é definida aqui como a média das distâncias aos elementos no agrupamento. Novamente, combinamos objetos que estão mais próximos. Também podemos calcular a distância entre os itens não agrupados e os itens agrupados como a distância mais próxima, em vez da média. Quando uma série está em um cluster, ela permanece lá. Em cada estágio, temos um número fixo de clusters, dependendo do estágio de fusão.

Alternativamente, podemos pensar no agrupamento como uma partição da amostra em um número pré-especificado de grupos. MacQueen (1967) propôs isso usando agrupamento de k-médias, usando a distância de Mahalanobis entre uma observação e os vetores de média do grupo.

Em cada estágio, uma redesignação de uma observação em seu grupo de afinidade mais próximo é possível. Para ver como este procedimento se aplica ao contexto atual, considere uma partição preliminar em um número fixo de grupos e defina a disparidade entre a matriz espectral da observação, digamos,  $\hat{f}$  e a matriz espectral média do grupo, digamos,  $\hat{f}_i$ , como  $J(\hat{f}; \hat{f}_i)$ , como  $J(\hat{f}; \hat{f}_i)$ , onde a matriz espectral do grupo pode ser estimada por

$$\hat{f}_i(\omega_k) = \frac{1}{N_i} \sum_{j=1}^{N_i} \hat{f}_{ij}(\omega_k).$$

Em qualquer passagem, uma única série é reatribuída ao grupo para o qual sua disparidade é minimizada.

O procedimento de reatribuição é repetido até que todas as observações fiquem em seus grupos atuais. Obviamente, o número de grupos deve ser especificado para cada repetição do algoritmo de particionamento e uma partição inicial deve ser escolhida. Essa atribuição pode ser aleatória ou escolhida a partir de um agrupamento hierárquico preliminar, conforme descrito acima.

## Exemplo VII.12. Análise de cluster em dados sísmicos.

É instrutivo tentar um procedimento de agrupamento na população de terremotos e explosões conhecidos. A figura abaixo mostra os resultados da aplicação do algoritmo de agrupamento Partitioning Around Medoids (PAM), que é essencialmente uma robustez do procedimento k-means, sob a suposição de que dois grupos são apropriados.

A partição de dois grupos tende a produzir uma partição final que concorda intimamente com a configuração conhecida com o terremoto 1 (EQ1) e a explosão 8 (EX8) sendo classificados incorretamente; como nos exemplos anteriores, o evento NZ é classificado como uma explosão.



## Clustering Results for Explosions and Earthquakes

