

O metodach identyfikacji składowych okresowych w wielowymiarowych szeregach czasowych

Maciej Kawecki

Instytut Matematyki i Informatyki
Politechnika Wrocławska

„Statystyka Matematyczna Wiśła 2009”

7-11 grudnia 2009

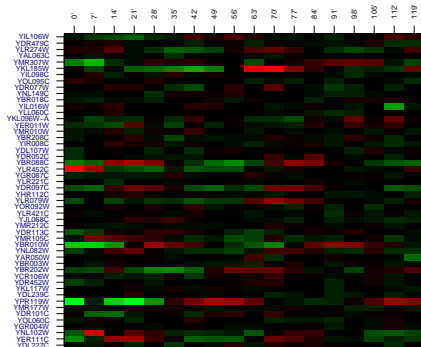
Plan prezentacji

- 1 Wprowadzenie - geneza i opis problemu
- 2 Przegląd stosowanych metod
- 3 Testy klasycznej analizy spektralnej wraz z modyfikacjami
- 4 Zastosowanie modelu RRRVAR
- 5 Porównanie efektywności testów
- 6 Analiza danych rzeczywistych



Charakterystyka analizowanych szeregów czasowych

- Długość szeregu czasowego $T \leq 50$
- Liczba składowych $M > 400$
- Składowe silnie skorelowane
- Dwie grupy składowych: okresowe i nieokresowe



Rysunek: Przykładowa *heatmapa* ekspresji genów komórki drożdży (Spellman, 1998)

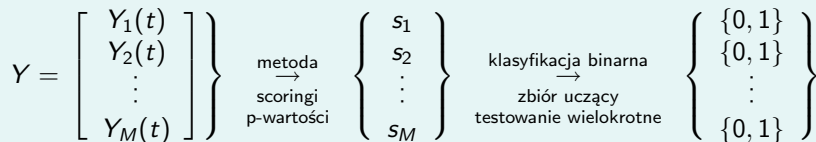
Cel: klasyfikacja

Każdej składowej chcemy przyporządkować wartość:

0 – nieokresowy komponent

1 – okresowy komponent

Przeprowadzamy klasyfikację scoringową



Jaką wiedzę posiadamy?

Ogólny model dla i -tej składowej

$$Y_i(t) = a_{i,0} + a_{i,1} \cdot f_i(t + \tau_i) + \varepsilon_i(t)$$

Dodatkowe założenia

funkcja okresowa

$$f_i(t) = f(t)$$

okres

$$s_i = s$$

wariancja błędów

$$\sigma_i^2 = \sigma^2$$

wyraz wolny

$$a_{i,0} = 0$$

Dodatkowe informacje a priori

- postać funkcji okresowej, jej częstotliwość
- wariancja błędów, struktura korelacji składowych
- zbiór uczący
 - podzbiór składowych, o których wiadomo czy są okresowe

- **jednowymiarowe metody niewykorzystujące wiedzy a-priori**
 - a) klasyczna analiza spektralna (Brockwell i in., Wichert i in., 2004)
- **wielowymiarowe metody niewykorzystujące wiedzy a-priori**
 - a) dopasowanie modelu RRRVAR (Kustra, Zagdański, 2006)
 - b) połączony test regularności i okresowości (Lichtenberg, 2005)
- **metody wykorzystujące wiedzę a-priori**
 - a) klasyczna analiza spektralna
 - b) podejście bayesowskie (Andersson i in., 2006)
 - c) test korelacyjny (Spellman i in., 1998)
 - d) interpolacja współdzielonej funkcji okresowej za pomocą B-splajnów (Luan, 2004)

TESTY KLASYCZNEJ ANALIZY SPEKTRALNEJ

Idea

Model:

$$Y_i(t) = a_{i,0} + a_{i,1} \cdot f(t + \tau_i) + \varepsilon_i(t)$$

Weryfikowane hipotezy:

$$H_0 : a_{i,1} = 0 \quad (\text{składowa nieokresowa})$$

$$H_1 : a_{i,1} \neq 0 \quad (\text{składowa okresowa})$$

Sposób obliczania scoringów

Scoringi to wartości statystyk testowych

Używane testy statystyczne

Brockwell, *Time Series: Theory and Methods*

G1 test Fishera ukrytej okresowości o nieokreślonej częstotliwości (G-test Fishera)

- $f(t) = \sin(\omega t + \tau_i)$
- $G_1 = \frac{\max_k I(\omega_k)}{\frac{1}{qT} \sum_{k=1}^q I(\omega_k)}$

G2 test ukrytej okresowości o określonej częstotliwości ω_0

- $f(t) = \sin(\omega_0 t + \tau_i)$
- $G_2 = \frac{(n-3)I(\omega_0)}{\sum_{t=1}^T Y_i^2(t) - I(0) - 2I(\omega_0)}$

G3 test weryfikujący obecność niesinusoidalnego komponentu okresowego o znanym okresie p

- $f_i(t) = \sum_{k=1}^{\lfloor \frac{p-1}{2} \rfloor} \left[A_{i,k} \cos\left(\frac{2\pi kt}{p}\right) + B_{i,k} \sin\left(\frac{2\pi kt}{p}\right) \right] + A_{\frac{p}{2}} (-1)^t$
- $G_3 = \frac{(T-p)}{(p-1)} \frac{2 \sum_{1 \leq j \leq p/2} I(\omega_{rj}) + \delta_p I(\pi)}{\sum_{t=1}^T Y_i^2(t) - I(0) - \delta_p I(\pi) - 2 \sum_{1 \leq j \leq p/2} I(\omega_{rj})}$

Stosowane wersje estymatorów gęstości spektralnej

- standardowy periodogram

$$I(\omega_k) = \sum_{|h| < T} \hat{\gamma}(h) \exp^{-i\omega_k h}$$

- odporny estymator gęstości spektralnej (Ahdesmaki i in., 2005)

$$I_{*1}(\omega) = \sum_{k=-T-2}^{t=T-2} \hat{\rho}(k) \exp(-i\omega k)$$

- wygładzony dyskretnie periodogram

$$I_{*2}(\omega_j) = (2\pi)^{-1} \sum_{|k| \leq m} W(k) I(\omega_j + k)$$

- wygładzony jądrowo periodogram

$$I_{*3}(\omega_j) = \hat{f}(\omega_j, h) = \frac{1}{Th} \sum_{k=-q_T}^{q_T} K\left(\frac{\omega_j - \omega_k}{h}\right) I(\omega_k)$$

Otrzymujemy 12 wersji testów klasycznej analizy spektralnej.

ZASTOSOWANIE MODELU RRRVAR

Definition (Model RRVAR(p, r))

Szereg czasowy $\{Y_t\}$ jest szeregiem czasowym z modelu RRVAR(p, r) jeżeli:

$$Y_t = ABX_t + \epsilon_t \quad (t = 1, \dots, T)$$

gdzie:

- $X_t = [Y'_{t-1}, \dots, Y'_{t-p}]'$ — wektor kolumnowy wymiaru: $Kp \times 1$
- A — macierz współczynników wymiaru $K \times r$
- $B = [B_1, B_2, \dots, B_p]$ — macierz wymiaru $r \times Kp$
- B_i — macierz współczynników wymiaru $r \times K$
- ϵ_t — wektor błędu o nieosobliwej macierzy kowariancji

Zauważmy, że macierze współczynników z wyjściowego procesu VAR(p) wyrażają się przy pomocy współczynników RRVAR(p, r) w następujący sposób $\Phi_i = AB_i$.

RRRVAR – szczegóły techniczne

Dysponując parametrami p oraz r estymujemy współczynniki macierzy, kanoniczne wagi oraz szeregi czasowe.

Kanoniczne szeregi czasowe

Dla $j = 1, \dots, r$ definiujemy j -ty kanoniczny szereg czasowy jako:

$$w_t^j = l_j' Y_t$$

Kanoniczne wagi

Definiujemy:

$$l_j = \Sigma_\epsilon^{-\frac{1}{2}} V_j',$$

gdzie V_j jest j -tym wektorem własnym macierzy:

$$W = \Sigma_{YY}^{-1/2} \Sigma'_{XY} \Sigma_{XX}^{-1} \Sigma_{XY} \Sigma_{YY}^{-1/2},$$

odpowiadającym j -tej największej wartości własnej.

Do wyznaczenia estymatorów Σ stosuje się estymatory autokorelacji

$$\tilde{\Gamma}(0) = \hat{\Gamma}(0) + \lambda I_{M \times M}.$$

Parametr λ nazywany jest współczynnikiem kary grzbietowej.

Szergi kanoniczne mogą być interpretowane jako najczęściej występujące wzorce w wielowymiarowym szeregu czasowym (również okresowe), więc w celu znalezienia komponentów okresowych możemy przeprowadzić procedurę:

Procedura wyznaczenia scoringów

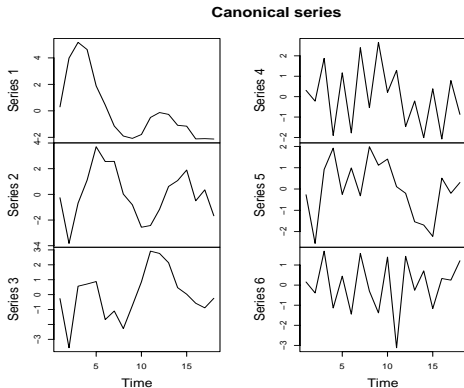
- 1 dopasujemy odpowiednie parametry p oraz r modelu
- 2 znajdujemy współczynniki macierzy, szeregi oraz wagi kanoniczne
- 3 testujemy, które szeregi kanoniczne są okresowe (otrzymujemy r_1 szeregów)
- 4 wybieramy wagi kanoniczne dla tych szeregów (otrzymujemy macierz $M \times r_1$ ($r_1 \leq r$))
- 5 dla każdego wiersza obliczamy maksimum z wartości bezwzględnych wybranych wag
w ten sposób uzyskujemy „scoring okresowości” dla każdego szeregu

Cykl podziału komórki drożdży - eksperyment alpha

- mamy 6178 genów oraz 18 obserwacji dla każdego szeregu.
- znamy 104 nazwy genów okresowych.
- Spellman zidentyfikował 800 okresowych genów.

Dopasowaliśmy model RRRVAR(2,6) używając kryteriów wyboru modelu (CV)

Otrzymujemy następujące szeregi czasowe.



SYMULACYJNE PORÓWNIANIE METOD

Jak porównać efektywność metod scoringowych?

uwagi ogólne

Czy dysponujemy zbiorem uczącym?

1 tak

- pole pod krzywymi ROC (AUC)
- błąd klasyfikacji
obliczony metodą k-fold cross validation

2 nie

- metody porównujące rangi
- klasyfikacja '*TOP N*' – diagramy Venna
porównanie pierwszych *N* najwyższej sklasyfikowanych składowych

Jak porównać efektywność metod scoringowych?

scenariusz symulacyjny

Ogólny model symulacyjny

$$Y(t) = [\underbrace{Y_{1,1}(t), \dots, Y_{1,M_1}(t)}_{\text{składowy 1}}, \dots, \underbrace{Y_{i,1}(t), \dots, Y_{i,M_i}(t)}_{\text{składowy } i}, \dots, \underbrace{Y_{k,1}(t), \dots, Y_{k,M_k}(t)}_{\text{składowy } k}]$$

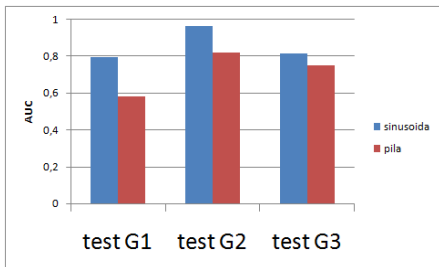
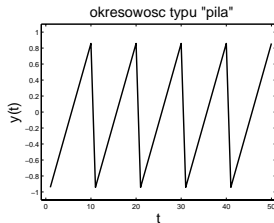
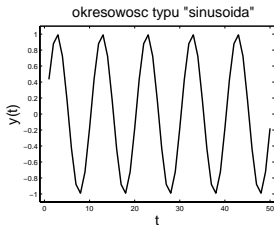
$$Y_{k,i}(t) = \mu_{k,i} + \alpha_k f_k(t - \tau_{k,i}) + \varepsilon_{k,i}(t)$$

Parametry modelu symulacyjnego

- długość szeregu czasowego i liczba składowych
- macierz kowariancji białego szumu
- postać funkcji okresowej i częstotliwość
- liczność frakcji składowych okresowych

Porównanie testów - wpływ typu funkcji okresowej

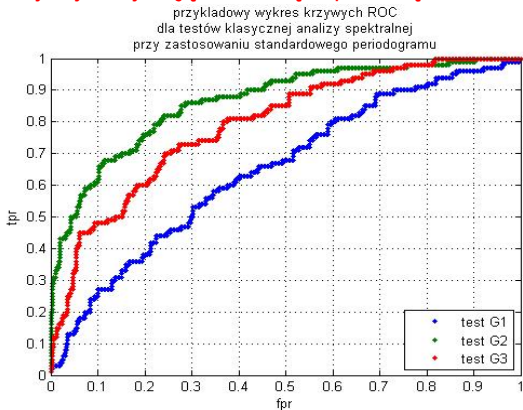
wyniki testów dla różnych funkcji okresowych
rodzaj funkcji okresowej ma wpływ na efektywność metody



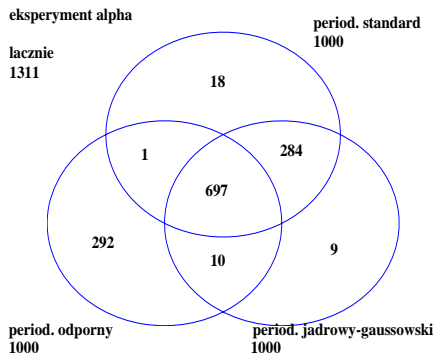
Porównanie testów - krzywe ROC

metody wykorzystujące wiedzę a-priori vs metody niewykorzystujące wiedzy a-priori

metody wykorzystujące wiedzę a-priori są skuteczniejsze

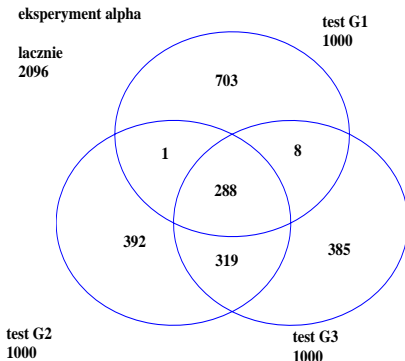


Stabilność testu w zależności od typu periodogramu



Rysunek: Diagram Venna – liczba genów z klasyfikacji 'TOP 1000' dla eksperymentu 'alpha' dla testu G-Fishera w zależności od zastosowanego periodogramu

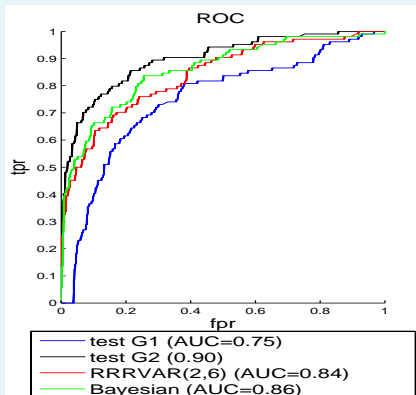
Rozbieżność wyników w zależności od zastosowanego testu



Rysunek: Diagram Venna – liczba genów z klasyfikacji 'TOP 1000' dla eksperymentu 'alpha' przy zastosowaniu standardowego periodogramu

Konstrukcja porównania krzywych ROC

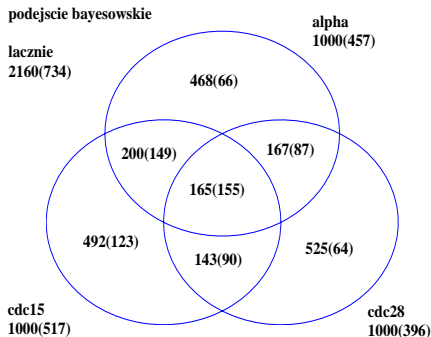
- zbiór testowy – 104 znane okresowe geny oraz 512 genów nie-Spellmana jako nieokresowe
- stosujemy różne metody w celu otrzymania scoringów
- dla każdej metody wyznaczamy ROC



- 1 nową metodę identyfikacji składowych okresowych można sprawdzać używając opisanego schematu
- 2 metody wykorzystujące wiedzę a-priori mają mniejszy błąd klasyfikacji niż metody niewykorzystujące tej wiedzy
- 3 metoda RRRVAR jest porównywalna z metodami jednowymiarowymi
- 4 metoda RRRVAR jest lepsza od standardowego testu Fishera (obie nie wykorzystują żadnej wiedzy a-priori)
- 5 znacznie lepiej używać metod właściwych dla określonego typu funkcji okresowej
– tę informację uzyskujemy używając np. RRRVAR

Jeden test – różne eksperymenty – brak zgodności wyników

- aktywność genów w czasie w cyklu podziału komórki drożdży (Spellman 1998)
- przeprowadzono 3 eksperymenty „cdc15”, „cdc28”, „alpha”
- dla każdego zbioru danych wyznaczmy „TOP 1000” genów w scoringu uzyskanym podejściem bayesowskim
- porównujemy listy otrzymanych genów okresowych
- w nawiasach liczba genów ze zbioru 800-Spellmana



Rysunek: Diagram Venna – liczba genów z klasyfikacji 'TOP 1000' dla różnych eksperymentów mierzących okresową aktywność genów cyklu podziału komórki drożdży

- 1 nową metodę identyfikacji składowych okresowych można sprawdzać używając opisanego schematu
- 2 metody wykorzystujące wiedzę a-priori mają mniejszy błąd klasyfikacji niż metody niewykorzystujące tej wiedzy
- 3 metoda RRRVAR jest porównywalna z metodami jednowymiarowymi
- 4 metoda RRRVAR jest lepsza od standardowego testu Fishera (obie nie wykorzystują żadnej wiedzy a-priori)
- 5 znacznie lepiej używać metod właściwych dla określonego typu funkcji okresowej
– tę informację uzyskujemy używając np. RRRVAR
- 6 należy zastosować kilka metod scoringowych do różnych zbiorów, a następnie „połączyć” rezultaty

Dziękuję za uwagę



Obliczanie scoringów

Wartości statystyk testowych traktujemy jako scoringi.

Obliczanie p-wartości dla różnych wersji estymatora gęstości spektralnej

- 1 dokładny rozkład statystyk testowych
 - standardowy periodogram
- 2 testowanie permutacyjne
 - odporny estymator gęstości spektralnej
 - wygładzony dyskretnie periodogram
- 3 metoda bootstrap w dziedzinie spektralnej
 - wygładzony jądrowo periodogram

Używane testy statystyczne (Brockwell)

G1 test Fishera ukrytej okresowości o nieokreślonej częstotliwości (G-test Fishera)

- $f(t) = \sin(\omega t + \tau_i)$
- $G_1 = \frac{\max_k I(\omega_k)}{\frac{1}{qT} \sum_{k=1}^q I(\omega_k)}$

G2 test ukrytej okresowości o określonej częstotliwości ω_0

- $f(t) = \sin(\omega_0 t + \tau_i)$
- $G_2 = \frac{(n-3)I(\omega_0)}{\sum_{t=1}^T Y_i^2(t) - I(0) - 2I(\omega_0)}$

G3 test weryfikujący obecność niesinusoidalnego komponentu okresowego o znanym okresie p

- $f_i(t) = \sum_{k=1}^{\lfloor \frac{p-1}{2} \rfloor} \left[A_{i,k} \cos\left(\frac{2\pi kt}{p}\right) + B_{i,k} \sin\left(\frac{2\pi kt}{p}\right) \right] + A_{\frac{p}{2}} (-1)^t$
- $G_3 = \frac{(T-p)}{(p-1)} \frac{2 \sum_{1 \leq j \leq p/2} I(\omega_{rj}) + \delta_p I(\pi)}{\sum_{t=1}^T Y_i^2(t) - I(0) - \delta_p I(\pi) - 2 \sum_{1 \leq j \leq p/2} I(\omega_{rj})}$

Połączony test regularności i okresowości

Przeprowadzamy dwa testy – regularności i okresowości.
Następnie łączymy otrzymane rezultaty.

Regularność - podejście wielowymiarowe

Składowe okresowe wielowymiarowego szeregu czasowego mają zazwyczaj większą wariancję wartości niż pozostałe składowe.

p-wartość testu regulacji wyznaczamy jako:

$$p_{i,reg} = \frac{1}{B} \sum_{j=1}^B 1 \left\{ std(Y_i) \leq std(X_j) \right\},$$

gdzie $X_j(t)$ dla $j = 1, \dots, B$ – B szeregów bazowych wyznaczonych metodą bootstrap poprzez replikowanie ze zbioru $\{Y_1(t), \dots, Y_M(t)\}$

Okresowość - podejście jednowymiarowe

Dla każdej składowej obliczamy p-wartość $p_{i,period}$ dowolnego testu weryfikującego występowanie okresowości.

Określenie scoringu

$$s_i = (p_{i,reg} \cdot p_{i,period}) \left[1 + \left(\frac{p_{i,reg}}{a} \right)^2 \right] \left[1 + \left(\frac{p_{i,period}}{a} \right)^2 \right],$$

gdzie $a \leq 1$ – współczynnik kary – powoduje, że składowe, które są nieistotne ani regularnie ani okresowo, otrzymują wyższy scoring
Końcowy scoring otrzymujemy poprzez przekształcenie:

$$s_{i,final} = \exp(-s_i)$$

Idea

Zakładamy model

$$Y_i(t) = a_{i,0} + a_{i,1} \cdot f(t - \tau_i) + Z_i(t)$$

Konstruujemy kompozytowy (zagregowany) szereg czasowy

$$\begin{aligned} X(t) &= \beta^* U(t) \\ f_X(\omega, \beta) &= \beta^* f_U(\omega) \beta \end{aligned}$$

gdzie

$$\begin{aligned} U(t) &= V^{-1/2} Y(t) \\ V &= \text{diag}\{\sigma_{Y_1}^2, \dots, \sigma_{Y_M}^2\} \end{aligned}$$

Definicja obwiedni spektralnej

$$\lambda(\omega) = \sup_{\beta \neq 0} \{f_X(\omega, \beta)\}$$

Po przekształceniach

$$\lambda(\omega) = \sup_{b \neq 0} \left\{ \frac{b^* f_Y(\omega) b}{b^* V b} \right\},$$

gdzie $b = V^{-1/2} \beta$

Rozwiązanie

$\lambda(\omega)$ jest największą wartością własną macierzy

$$f_U(\omega) = V^{-1/2} f_Y(\omega) V^{-1/2}, \quad (1)$$

natomiast $\beta(\omega)$ jest odpowiadającym wektorem własnym.

Zaproponowana interpretacja

- 1 $\beta(\omega)$ – optymalne M -wymiarowe skalowanie gęstości spektralnej f_Y w punkcie ω
- 2 $|\beta(\omega_0)|$ – wektor scoringów
- 3 $\lambda(\omega)$ – narzędzie wstępnej analizy szeregu czasowego

Stosowane wersje estymatorów gęstości spektralnej

- standardowy periodogram

$$I(\omega_k) = \sum_{|h| < T} \hat{\gamma}(h) \exp^{-i\omega_k h}$$

- odporny estymator gęstości spektralnej

$$I_{*1}(\omega) = \sum_{k=-T-2}^{t=T-2} \hat{\rho}(k) \exp(-i\omega k)$$

- wygładzony dyskretnie periodogram

$$I_{*2}(\omega_j) = (2\pi)^{-1} \sum_{|k| \leq m} W(k) I(\omega_{j+k})$$

- wygładzony jądrowo periodogram

$$I_{*3}(\omega_j) = \hat{\gamma}(\omega_j, h) = \frac{1}{Th} \sum_{k=-qT}^{qT} K\left(\frac{\omega_j - \omega_k}{h}\right) I(\omega_k)$$

Podejście bayesowskie - idea

Model:

$$H_0 : Y_i(t) = a_{i,0,0} + Z_i(t) \quad (\text{składowa nieokresowa})$$

$$H_1 : Y_i(t) = a_{i,1,0} + a_{i,1,1} \cdot \cos(\omega_0 t + \tau_i) + Z_i(t) \quad (\text{składowa okresowa})$$

Obliczamy prawdopodobieństwa a posteriori dla każdej składowej.
Stosując aproksymację bayesowskiego kryterium informacyjnego (BIC)
obliczmy scoringi.

Prawdopodobieństwa a posteriori

$$P(H_k|D) = \frac{P(D|H_k)P(H_k)}{\sum_{j=\{1,2\}} P(D|H_j)P(H_j)}$$

$$P(D|H_k) = \int_{\phi_{i,k} \in \Phi_{i,k}} p(D|\phi_{i,k}, H_k)p(\phi_{i,k}|H_k) d\phi_{i,k}$$

Założenia a priori

- parametry niezależne
- znana średnia μ_{ω_0} , wariancja σ_{ω_0}
- rozkład normalny parametru ω_0
- rozkład jednostajny parametrów $a_{i,k,j}$
- rozkład normalny zaburzeń Z_i

Wyznaczenie scoringów

Aproksymacja bayesowskiego kryterium informacyjnego

$$BIC(H_k) = \log p(D|\phi_{i,k}^{MAP}, H_k) - \frac{d}{2} \log T \sim \log P(D|H_k)$$

$$g(D) = BIC(H_1) - BIC(H_0)$$

$$s_i = \frac{\exp^{g(D)}}{1 + \exp^{g(D)}}$$

Funkcje B-sklejane

P - wymiar bazy

$1 = t_0 \leq t_1 \leq \dots \leq t_P = T$ - węzły

$\{b_{p,k}(t), p = 0, 1, \dots, P-1\}$ - baza funkcji B-sklejanych rzędu k

Wzory rekurencyjne Coxa-de Boora

$$b_{p,0}(t) = \begin{cases} 1 & , \quad t_p \leq t \leq t_{p+1} \\ 0 & , \quad t \notin (t_p, t_{p+1}) \end{cases} , \quad \text{dla } p = 0, 1, \dots, P-1$$

$$b_{p,j}(t) = 0, \quad \text{gdyn } p \geq P \text{ i } j \geq 0,$$

$$b_{p,k}(t) = \frac{t - t_p}{t_{p+k} - t_p} b_{p,k-1}(t) + \frac{t_{p+k+1} - t}{t_{p+k+1} - t_{p+1}} b_{p+1,k-1}(t)$$

dla $k > 0$ i $p < P$.

Baza sześciennych funkcji B-sklejanych - aproksymacja funkcji

$$B_p(t) = b_{p-1,3}(t), \quad \text{dla } p = 1, \dots, P$$

$$f(t) = \sum_{p=1}^P \gamma_p B_p(t)$$

Model szeregu czasowego

$$Y_i(t) = a_{i,0} + a_{i,1}f(t - \tau_i) + Z_i(t), \quad \text{dla } i = 1, \dots, M \text{ oraz } t = 1, \dots, T$$

Zbiór uczący

szeregi czasowe współdzielące wspólną funkcję okresową $f(t)$

$$U_l(t) = a_{l,0} + a_{l,1}f(t - \tau_l) + Z_l(t), \quad \text{dla } l = 1, \dots, L$$

KROK 1 szukamy $a_{l,0}, a_{l,1}, \tau_l$,
przyjmujemy $f(t) = \sin(t)$

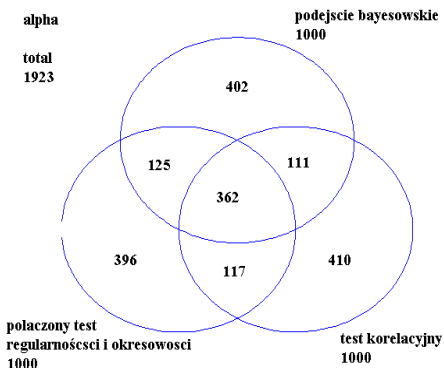
$$\min \left\{ \frac{1}{T} \sum_{t=1}^T [U_l(t) - a_{l,0} - a_{l,1} f(t - \tau_l)]^2 \right\}$$

KROK 2 szukamy γ_p ,
przyjmujemy $a_{l,0}, a_{l,1}, \tau_l$ obliczone w kroku 1

$$\min \left\{ \sum_{l=1}^L \frac{1}{T} \sum_{t=1}^T [u_l(t) - a_{l,0} - a_{l,1} \sum_{p=1}^P \gamma_k B_p(t - \tau_l)]^2 \right\}$$

KROK 3 szukamy $a_{i,0}, a_{i,1}, \tau_i$,
przyjmujemy γ_p obliczone w kroku 2

$$\min \left\{ \frac{1}{T} \sum_{t=1}^T [y_i(t) - a_{i,0} - a_{i,1} f(t - \tau_i)]^2 \right\}$$



Rysunek: Przykładowy diagram Venna – liczba genów z klasyfikacji 'TOP 1000' dla eksperymentu 'alpha'