

**Politechnika Gdańska**  
Wydział Fizyki Technicznej i Matematyki Stosowanej

**Anna Wieżel**

Nr albumu: 132540

# **Funkcjonalne Modele Liniowe**

**Praca magisterska**  
**na kierunku MATEMATYKA**  
**w zakresie MATEMATYKA FINANSOWA**

Praca wykonana pod kierunkiem  
**dra hab. Karola Dziedziula**  
Katedra Analizy Matematycznej i Numerycznej

Wrzesień 2015

## **Oświadczenie kierującego prac**

Potwierdzam, że niniejsza praca została przygotowana pod moim kierunkiem i kwalifikuje się do przedstawienia jej w postępowaniu o nadanie tytułu zawodowego.

Data

Podpis kierującego prac

## **Oświadczenie autora (autorów) pracy**

Świadom odpowiedzialności prawnej oświadczam, że niniejsza praca dyplomowa została napisana przeze mnie samodzielnie i nie zawiera treści uzyskanych w sposób niezgodny z obowiązującymi przepisami.

Oświadczam również, że przedstawiona praca nie była wcześniej przedmiotem procedur związanych z uzyskaniem tytułu zawodowego w wyższej uczelni.

Oświadczam ponadto, że niniejsza wersja pracy jest identyczna z załączoną wersją elektroniczną.

Data

Podpis autora (autorów) pracy

## **Streszczenie**

The paper's motivation is to contribute to popularization of mathematical statistics on infinite dimensional function Hilbert spaces. The author presents the fully functional linear model in form  $Y = \beta X + \varepsilon$  and its significance test proposed by Kokoszka et al. The test detects nullity of operator  $\beta$  which implies the lack of linear dependence between  $X$  and  $Y$ . Using the principal component decomposition it is concluded with test statistic convergent by distribution to chi-squared.

The test is further used for magnetic field data collected in some stations in different latitudes. The results show linear dependence between horizontal intensities of the magnetic field in mid- and low-latitude stations with high-latitude station data with a day or two delay but they contradict the linear dependence between data with more than a two-day lag.

## **Słowa kluczowe**

analiza danych funkcyjnych, dane funkcyjne, funkcyjny model liniowy, test istotności

## **Dziedzina pracy (kody wg programu Socrates-Erasmus)**

11.1 Matematyka

11.2 Statystyka

## **Klasyfikacja tematyczna**

62 Statistics

62-07 Data analysis

62J12 Generalized linear models

## **Tytuł pracy w języku angielskim**

Functional Linear Models



# Spis treści

<b>Wstęp</b>	5
<b>1. Preliminaria</b>	7
1.1. Klasyfikacja operatorów liniowych	7
1.2. Przestrzeń $L^2$	9
1.3. Zmienne funkcyjne w $L^2$ . Pojęcie średniej i operatora kowariancji	11
1.4. Estymacja średniej, funkcji kowariancji i operatora kowariancji	12
1.5. Estymacja wartości własnych i funkcji własnych operatora kowariancji	13
1.6. Funkcyjny model liniowy	13
<b>2. Test istotności w funkcyjnym modelu liniowym</b>	17
2.1. Procedura testowa	17
2.2. Formalne podstawy [nazwa?]	19
<b>3. Przykład zastosowania</b>	25
3.1. Ameryka Północna (Kanada)	26
3.2. Europa (Polska)	26
<b>A. Kod w R</b>	27
<b>Bibliografia</b>	29



# Wstęp

[już tu: Przykłady danych funkcjonalnych?]

[już tu: próba = punkty - ostatecznie: funkcja gładka?]

[Odpowiednik testu istotności dla prostego modelu regresji = F-test (+ t-test) [patrz: artykuł]]

[pakiet w R: fda]

Praca opiera się głównie na artykule [Kokoszka et al. (2008)], który to został użyty i rozwinięty w książce [Horváth, Kokoszka].

Ze względu na to, że analiza danych funkcjonalnych (*ang.* Functional Data Analysis, FDA) jest stosunkowo nowym działem statystyki i jest wciąż mało popularna w polskiej literaturze, wiele pojęć czy określeń zawartych w pracy nie posiada jeszcze ogólnie przyjętych polskich odpowiedników. Dlatego zostały one przetłumaczone przez autora według własnego uznania, przytaczając oczywiście oryginalne (angielskie) nazwy.

ACKNOWLEDGEMENTS/podziękowania?

The results presented in this paper rely on data collected at magnetic observatories. We thank the national institutes that support them and INTERMAGNET for promoting high standards of magnetic observatory practice ([www.intermagnet.org](http://www.intermagnet.org)).





# Rozdział 1

## Preliminaria

Przestrzenią funkcyjną  $E$  nazywać będziemy przestrzeń liniową funkcji z dowolnego zbioru  $A$  do zbioru  $B$ .

**Definicja 1.0.1** [Ferraty, Vieu]

Zmienną losową  $X$  nazywamy **zmienną funkcjonalną** (ang. *functional variable*) wtedy i tylko wtedy, gdy przyjmuje wartości w nieskończenie wymiarowej przestrzeni (przestrzeni funkcyjnej). Obserwację  $\chi$  zmiennej  $X$  nazywamy **daną funkcjonalną** (ang. *functional data*).

Jeśli zmienna funkcjonalna  $X$  (odpowiednio obserwacja  $\chi$ ) jest krzywą, to możemy przedstawić  $X$  w następującej postaci  $X = \{X(t), t \in T\}$  (odp.  $\chi = \{\chi(t), t \in T\}$ ), gdzie zbiór indeksów  $T \subset \mathbb{R}$ . Taką zmienną funkcjonalną możemy zatem utożsamiać z procesem stochastycznym z nieskończenie wymiarową przestrzenią stanów. W szczególności, zmienna funkcjonalna może być powierzchnią, czyli dwuwymiarowym wektorem krzywych - wtedy, analogicznie,  $T$  będzie dwuwymiarowym zbiorem indeksów tj.  $T \subset \mathbb{R}^2$  - lub dowolnie wymiarowym wektorem krzywych.

W niniejszej pracy skupimy się na zmiennych funkcjonalnych przyjmujących postać krzywych.

[przykłady? czy tylko we wstępie?]

[tu: próba = punkty - ostatecznie: funkcja gładka?]

Aby zbudować pojęcie operatora kowariancji dla zmiennych funkcjonalnych wprowadzimy niezbędne pojęcia z dziedziny operatorów liniowych.

### 1.1. Klasyfikacja operatorów liniowych

Niech  $(\Omega, \mathcal{F}, P)$  będzie przestrzenią probabilistyczną,  $\Omega$  jest zatem zbiorem scenariuszy  $\omega$ ,  $\mathcal{F}$  jest  $\sigma$ -algebrą podzbiorów  $\Omega$ , a  $P$  miarą prawdopodobieństwa nad  $\mathcal{F}$ . Dla uproszczenia zakładamy zupełność zadanej przestrzeni probabilistycznej. Rozważmy proces stochastyczny z czasem ciągłym  $X = \{X_t, t \in T\}$ , gdzie  $T$  jest przedziałem w  $\mathbb{R}$ , zdefiniowany na przestrzeni probabilistycznej  $(\Omega, \mathcal{F}, P)$ , taki, że  $X_t(\omega)$  należy do przestrzeni funkcyjnej  $E$  dla wszystkich  $\omega \in \Omega$ .

W pracy rozważać będziemy zmienne funkcjonalne przyjmujące wartości w przestrzeni Hilberta.

Rozważmy ośrodkową nieskończenie wymiarową rzeczywistą przestrzeń Hilberta  $H$  z iloczynem skalarnym  $\langle \cdot, \cdot \rangle$  zadającym normę  $\|\cdot\|$  i oznaczmy przez  $\mathcal{L}$  przestrzeń ciągłych (ograni-

czonych) operatorów liniowych w  $H$  z normą

$$\|\Psi\|_{\mathcal{L}} := \sup\{\|\Psi(x)\| : \|x\| \leq 1\}.$$

**Definicja 1.1.1** [Horváth, Kokoszka]

Operator  $\Psi \in \mathcal{L}$  nazywamy **operatorem zwartym**, jeśli istnieją dwie ortonormalne bazy w  $H$   $\{v_j\}_{j=1}^{\infty}$  i  $\{f_j\}_{j=1}^{\infty}$ , oraz ciąg liczb rzeczywistych  $\{\lambda_j\}_{j=1}^{\infty}$  zbieżny do zera, takie że

$$\Psi(x) = \sum_{j=1}^{\infty} \lambda_j \langle x, v_j \rangle f_j, \quad x \in H. \quad (1.1)$$

Bez straty ogólności możemy założyć, że w przedstawionej reprezentacji  $\lambda_j$  są wartościami dodatnimi, w razie konieczności wystarczy  $f_j$  zamienić na  $-f_j$ .

Równoważną definicją operatora zwartego jest spełnienie przez  $\Psi$  następującego warunku: zbieżność  $\langle y, x_n \rangle \rightarrow \langle y, x \rangle$  dla każdego  $y \in H$  implikuje  $\|\Psi(x_n) - \Psi(x)\| \rightarrow 0$ .

Inną klasą operatorów są operatory Hilberta-Schmidta, którą oznaczać będziemy przez  $\mathcal{S}$ .

**Definicja 1.1.2** [Bosq]

**Operatorem Hilberta-Schmidta** nazywamy taki operator zwarty  $\Psi \in \mathcal{L}$ , dla którego ciąg  $\{\lambda_j\}_{j=1}^{\infty}$  w reprezentacji (1.1) spełnia  $\sum_{j=1}^{\infty} \lambda_j^2 < \infty$ .

**Uwaga 1.1.1** [Bosq], [Horváth, Kokoszka]

Klasa  $\mathcal{S}$  jest przestrzenią Hilberta z iloczynem skalarnym

$$\langle \Psi_1, \Psi_2 \rangle_{\mathcal{S}} := \sum_{j=1}^{\infty} \langle \Psi_1(e_j), \Psi_2(e_j) \rangle, \quad (1.2)$$

gdzie  $\{e_j\}_{j=1}^{\infty}$  jest dowolną bazą ortonormalną w  $H$ .

Powyższy iloczyn skalarny zadaje normę

$$\|\Psi\|_{\mathcal{S}} := \left( \sum_{j=1}^{\infty} \lambda_j^2 \right)^{1/2}. \quad (1.3)$$

Dowód równości (1.3).

$$\begin{aligned} \|\Psi\|_{\mathcal{S}}^2 &= \langle \Psi, \Psi \rangle_{\mathcal{S}} = \sum_{n=1}^{\infty} \left\langle \sum_{j=1}^{\infty} \lambda_j \langle e_n, v_j \rangle f_j, \sum_{k=1}^{\infty} \lambda_k \langle e_n, v_k \rangle f_k \right\rangle \\ &= \sum_{n=1}^{\infty} \sum_{j=1}^{\infty} \sum_{k=1}^{\infty} \lambda_j \lambda_k \langle e_n, v_j \rangle \langle e_n, v_k \rangle \langle f_j, f_k \rangle = \sum_{n=1}^{\infty} \sum_{j=1}^{\infty} \lambda_j^2 \langle e_n, v_j \rangle^2 \\ &= \sum_{j=1}^{\infty} \lambda_j^2 \sum_{n=1}^{\infty} \langle e_n, v_j \rangle^2 \stackrel{\text{tożsamość Parsevala}}{=} \sum_{j=1}^{\infty} \lambda_j^2 \|v_j\|^2 = \sum_{j=1}^{\infty} \lambda_j^2. \end{aligned}$$

□

**Definicja 1.1.3** [Bosq]

Zwarty operator liniowy nazywamy **operatorem śladowym** (ang. nuclear operator), jeśli równość (1.1) spełniona jest dla ciągu  $\{\lambda_j\}_{j=1}^{\infty}$  takiego, że  $\sum_{j=1}^{\infty} |\lambda_j| < \infty$ .

**Uwaga 1.1.2** [Bosq]

Klasa operatorów śladowych  $\mathcal{N}$  z normą  $\|\Psi\|_{\mathcal{N}} := \sum_{j=1}^{\infty} |\lambda_j|$  jest przestrzenią Banacha.

**Definicja 1.1.4** [Horváth, Kokoszka]

Operator  $\Psi \in \mathcal{L}$  nazywamy **symetrycznym**, jeśli

$$\langle \Psi(x), y \rangle = \langle x, \Psi(y) \rangle, \quad x, y \in H,$$

oraz **nieujemnie określonym** (lub połowicznie pozytywnie określonym, ang. positive semi-definite), jeśli

$$\langle \Psi(x), x \rangle \geq 0, \quad x \in H.$$

**Uwaga 1.1.3** [Horváth, Kokoszka]

Symetryczny nieujemnie określony operator Hilberta-Schmidta  $\Psi$  możemy przedstawić w reprezentacji

$$\Psi(x) = \sum_{j=1}^{\infty} \lambda_j \langle x, v_j \rangle v_j, \quad x \in H, \quad (1.4)$$

gdzie ortonormalne  $v_j$  są **funkcjami własnymi**  $\Psi$ , tj.  $\Psi(v_j) = \lambda_j v_j$ . Funkcje  $v_j$  mogą być rozszerzone do bazy, przez dopełnienie ortogonalne podprzestrzeni rozpiętej przez oryginalne  $v_j$ . Możemy zatem założyć, że funkcje  $v_j$  w (1.4) tworzą bazę, a pewne wartości  $\lambda_j$  mogą być równe zero.

**1.2. Przestrzeń  $L^2$** 

Przestrzeń  $L^2 = L^2(T) = L^2(T, \mathcal{B}, \lambda)$  nad pewnym przedziałem  $T \subset \mathbb{R}$  jest zbiorem klas mierzalnych funkcji rzeczywistych całkowalnych z kwadratem określonych na  $T$ , tj.

$$x \in L^2(T) \iff x : T \rightarrow \mathbb{R} \wedge \int_T x^2(t) dt < \infty.$$

Przestrzeń  $L^2$  jest ośrodkową przestrzenią Hilberta z iloczynem skalarnym

$$\langle x, y \rangle := \int_T x(t)y(t)dt, \quad x, y \in L^2.$$

Normę zaś wyznacza wzór

$$\|x\|^2 = \langle x, x \rangle = \int x^2(t)dt, \quad x \in L^2.$$

Tak jak zwyczajowo zapisujemy  $L^2$  zamiast  $L^2(T)$ , tak w przypadku symbolu całki bez wskazania obszaru całkowania będziemy mieć na myśli całkowanie po całym przedziale  $T$ . Jeśli  $x, y \in L^2$ , równość  $x = y$  zawsze oznaczać będzie  $\int [x(t) - y(t)]^2 dt = 0$ .

Ważną klasę operatorów liniowych na przestrzeni  $L^2$  stanowią operatory całkowite.

**Definicja 1.2.1** [Pytlik]

**Operatorem całkowym** nazywamy operator liniowy  $\Psi$  dający się przedstawić w formie

$$\Psi(x)(t) = \int \psi(t, s)x(s)ds, \quad x \in L^2, \quad t \in T,$$

gdzie  $\psi$  jest mierzalną **[ciągłą = całkowalną?]** funkcją dwóch zmiennych nazywaną **jądrem całkowym** operatora  $\Psi$ .

Operator całkowy  $\Psi$  jest dobrze określony, jeśli spełnia pewnego rodzaju własność ograniczoności.

**Uwaga 1.2.1** [Wojtaszczyk]

Niech  $(T, \mu)$  będzie przestrzenią z miarą i niech  $\psi(t, s)$  będzie mierzalną funkcją na  $T \times T$ . Zdefiniujmy

$$\Psi x(s) = \int_T \psi(t, s)x(t)d\mu(t).$$

Jeśli  $1 < p < \infty$  oraz istnieje mierzalna dodatnia funkcja  $y$  na  $T$  oraz stałe  $a, b$  takie że dla  $\frac{1}{p} + \frac{1}{q} = 1$  mamy

$$\int_T |\psi(t, s)|y(t)^q d\mu(t) \leq [ay(s)]^q, \quad \mu - pr.w. \quad (1.5)$$

oraz

$$\int_T |\psi(t, s)|y(t)^p d\mu(t) \leq [by(s)]^p, \quad \mu - pr.w., \quad (1.6)$$

wtedy  $T : L^p(T, \mu) \rightarrow L^p(T, \mu)$ .

[całe twierdzenie? jakie dać oznaczenia?]

Dowód. Niech  $x \in L^p$  i niech  $y$  spełnia założenia twierdzenia. Mamy

$$\begin{aligned} |\Psi x(s)| &= \left| \int_T \psi(t, s)x(t)d\mu(t) \right| \leq \int_T |\psi(t, s)||x(t)|d\mu(t) \\ &= \int_T [|\psi(t, s)|^{1/q}y(t)] \cdot [|\psi(t, s)|^{1/p}|x(t)|y(t)^{-1}]d\mu(t) \\ &= \int_T [|\psi(t, s)|y^q(t)]^{1/q} \cdot [|\psi(t, s)|(|x|/y)^p(t)^{-1}]^{1/p}d\mu(t) \\ &\stackrel{\text{nierówn. Höldera}}{\leq} \int_T ay(s) \cdot \left[ \int_T |\psi(t, s)|(|x|/y)^p(t)d\mu(t) \right]^{1/p} d\mu(s). \end{aligned} \quad (1.7)$$

Stąd, korzystając z twierdzenia Fubiniego, otrzymujemy

$$\begin{aligned} \|\Psi x\|_p &\leq \left\| \int_T \left| \int_T \psi(t, s)x(t)d\mu(t) \right|^p d\mu(s) \right\|^{1/p} \\ &\stackrel{(1.7)}{\leq} a \left[ \int_T y^p(s) \int_T |\psi(t, s)|(|x|/y)^p(t)d\mu(t)d\mu(s) \right]^{1/p} \\ &\stackrel{\text{tw.F.}}{=} a \left[ \int_T (|x|/y)^p(t) \int_T y^p(s)|\psi(t, s)|d\mu(s)d\mu(t) \right]^{1/p} \\ &\stackrel{(1.6)}{\leq} ab \left[ \int_T (|x|/y)^p(t)y^p(t)d\mu(t) \right]^{1/p} = ab \left[ \int_T |x(t)|^p d\mu(t) \right]^{1/p} = ab\|x\|_p < \infty. \end{aligned}$$

Pokazaliśmy, że  $\Psi x \in L^p$ , co kończy dowód. □

**Uwaga 1.2.2** [Horváth, Kokoszka]

Operatory całkowe są operatorami Hilberta-Schmidta wtedy i tylko wtedy, gdy

$$\iint \psi^2(t, s)dtds < \infty. \quad (1.8)$$

Ponadto zachodzi

$$\|\Psi\|_S^2 = \iint \psi^2(t, s)dtds.$$

**Uwaga 1.2.3** (Twierdzenie Mercera) [Horváth, Kokoszka]

Niech operator  $\Psi$  będzie operatorem całkowym spełniającym (1.8). Jeśli ponadto jego jądro całkowe  $\psi$  spełnia  $\psi(s, t) = \psi(t, s)$  oraz  $\iint \psi(t, s)x(t)x(s)dtds \geq 0$ , to operator całkowy  $\Psi$  jest symetryczny i nieujemnie określony, zatem z Uwagi 1.1.3 mamy

$$\psi(t, s) = \sum_{j=1}^{\infty} \lambda_j v_j(t) v_j(s) \quad \text{w } L^2(T \times T),$$

gdzie  $\lambda_j, v_j$  są odpowiednio wartościami własnymi i funkcjami własnymi operatora  $\Psi$ . Jeżeli funkcja  $\psi$  jest ciągle, powyższe rozwinięcie jest prawdziwe dla wszystkich  $t, s \in T$  i szereg jest zbieżny jednostajnie.

### 1.3. Zmienne funkcjonalne w $L^2$ . Pojęcie średniej i operatora kowariancji

Rozważmy zmienną funkcjonalną  $X = \{X(t), t \in T\}$  będącą krzywą ( $T \subset \mathbb{R}$ ) jako element losowy z przestrzeni  $L^2(T)$  zaopatrzonej w  $\sigma$ -algebrę borelowskich podzbiorów  $T$ .

Mówimy, że zmienna  $X$  jest **całkowalna**, jeśli  $\mathbb{E}\|X\| = \mathbb{E}[\int X^2(t)dt]^{1/2} < \infty$ . Jeśli  $X$  jest całkowalna, to istnieje jedyna funkcja  $\mu \in L^2$  taka, że  $\mathbb{E}\langle y, X \rangle = \langle y, \mu \rangle$  dla dowolnej funkcji  $y \in L^2$  (zauważmy, że wartość oczekiwana jest funkcjonałem liniowym, możemy zatem skorzystać z twierdzenia Riesz). Zachodzi  $\mu(t) = \mathbb{E}[X(t)]$  dla prawie wszystkich  $t \in T$ , tak określoną funkcję  $\mu$  nazywać będziemy **funkcją średniej**. Ponadto, wartość oczekiwana jest przemieniana z operatorami ograniczonymi, tj. jeśli  $X$  jest całkowalna oraz  $\Psi \in \mathcal{L}$ , to mamy  $\mathbb{E}\Psi(X) = \Psi(\mathbb{E}X)$ .

**Definicja 1.3.1** [Bosq]

**Operator kowariancji** całkowalnej zmiennej funkcjonalnej  $X$  o funkcji średniej  $\mu_X$  przyjmującej wartości w przestrzeni funkcyjnej  $L^2$  spełniającej  $\mathbb{E}\|X\|^2 < \infty$  definiujemy jako ograniczony operator liniowy według wzoru

$$C_X(x) := \mathbb{E}[\langle X - \mu_X, x \rangle (X - \mu_X)], \quad x \in L^2.$$

Jeśli  $Y$  jest zmienną funkcjonalną o funkcji średniej  $\mu_Y$  spełniającą powyższe warunki, wtedy operator kowariancji między zmiennymi  $X$  i  $Y$  (ang. cross-covariance operator) przedstawiamy jako

$$C_{X,Y}(x) := \mathbb{E}[\langle (X - \mu_X), x \rangle (Y - \mu_Y)], \quad x \in L^2$$

oraz

$$C_{Y,X}(x) := \mathbb{E}[\langle (Y - \mu_Y), x \rangle (X - \mu_X)], \quad x \in L^2.$$

Operator kowariancji jest operatorem całkowym, czyli

$$C_X(x)(t) = \int c(t, s)x(s)ds,$$

gdzie jądro całkowe  $c(t, s)$  zdefiniowane następująco

$$c(t, s) = \mathbb{E}[(X(t) - \mu(t))(X(s) - \mu(s))]$$

nazywać będziemy **funkcją kowariancji**. Oczywiście jest, że  $c(t, s) = c(s, t)$  i mamy

$$\begin{aligned}\iint c(t, s)x(t)x(s)dtds &= \iint \mathbb{E}[(X(t) - \mu(t))(X(s) - \mu(s))]x(t)x(s)dtds \\ &= \mathbb{E} \left[ \left( \int X(t)x(t)dt \right)^2 \right] \geq 0.\end{aligned}$$

Zatem operator kowariancji  $C_X$  jest symetryczny oraz nieujemnie określony. Wartości własne  $\lambda_j$  operatora  $C_X$  są dodatnie i spełniony jest warunek  $\sum_{j=1}^{\infty} \lambda_j = \mathbb{E} \|X\|^2 < \infty$ .  $C_X$  jest operatorem Hilberta-Schmidta (a nawet operatorem śladowym) i posiada on następującą reprezentację

$$C_X(x) = \sum_{j=1}^{\infty} \lambda_j \langle x, v_j \rangle v_j, \quad x \in L^2.$$

[więcej w Bosq (2000), rozdz. 1]

## 1.4. Estymacja średniej, funkcji kowariancji i operatora kowariancji

Naturalnym problemem pojawiającym się przy danych funkcjonalnych jest wnioskowanie o obiektach nieskończenie wymiarowych na podstawie skończonej próbki danych.

Obserwujemy zatem  $N$  krzywych  $X_1, \dots, X_N$ , które możemy traktować jako realizacje losowej funkcji  $X$  lub obserwacje zmiennej funkcjonalnej  $X$  z przestrzeni  $L^2$ .

**Założenie 1.4.1** [Horváth, Kokoszka]

*Zakładamy, że  $X_1, \dots, X_N$  są niezależnymi zmiennymi losowymi w  $L^2$  o jednakowym rozkładzie jak zmienna  $X \in L^2$ .*

[rozkład?]

Poszukiwanymi parametrami są funkcja średniej, funkcja kowariancji oraz operator kowariancji, określone następująco

$$\begin{aligned}\text{funkcja średniej:} & \quad \mu(t) = \mathbb{E}[X(t)]; \\ \text{funkcja kowariancji:} & \quad c(t, s) = \mathbb{E}[(X(t) - \mu(t))(X(s) - \mu(s))]; \\ \text{operator kowariancji:} & \quad C = \mathbb{E}[\langle (X - \mu), \cdot \rangle (X - \mu)].\end{aligned}$$

Funkcję średniej  $\mu$  estymujemy średnią z funkcji z próby

$$\hat{\mu}(t) = \frac{1}{N} \sum_{n=1}^N X_n(t), \quad t \in T,$$

funkcję kowariancji ze wzoru

$$\hat{c}(t, s) = \frac{1}{N} \sum_{n=1}^N (X_n(t) - \hat{\mu}(t))(X_n(s) - \hat{\mu}(s)), \quad t, s \in T,$$

zaś operator kowariancji estymujemy

$$\hat{C}(x) = \frac{1}{N} \sum_{n=1}^N \langle X_n - \hat{\mu}, x \rangle (X_n - \hat{\mu}), \quad x \in L^2. \quad (1.9)$$

[więcej w [Horváth, Kokoszka], rozdz. 2]

## 1.5. Estymacja wartości własnych i funkcji własnych operatora kowariancji

W dalszej części pracy istotne będzie dla nas oszacowanie również wartości i funkcji własnych operatora kowariancji  $C$ . W szczególności interesować nas będzie  $p$  największych wartości własnych spełniających

$$\lambda_1 > \lambda_2 > \dots > \lambda_p > \lambda_{p+1}$$

oraz aby  $p$  pierwszych wartości własnych było zerowych.

Funkcje własne zdefiniowane są przez równanie  $Cv_j = \lambda_j v_j$ . Zauważmy, że (z definicji operatora liniowego), jeśli  $v_j$  jest funkcją własną, to również  $av_j$  jest funkcją własną, gdzie  $a \neq 0$  jest skalar.

[...]

Wartości i funkcje własne estymujemy według wzoru

$$\int \hat{c}(t, s) \hat{v}_j(s) ds = \hat{\lambda}_j \hat{v}_j(t), \quad j = 1, 2, \dots, N.$$

[...]

[EFPC: Interference... rozdz. 3]

## 1.6. Funkcjonalny model liniowy

Standardowy model liniowy dla par zmiennych skalarnych  $Y_n$  i wektorów  $\mathbf{X}_n$  (tworzonych przez  $p$  skalarnych zmiennych  $X_{ni}$ ,  $i = 1, \dots, p$ ), przy założeniu  $\mathbb{E}Y_n = 0$ ,  $\mathbb{E}\mathbf{X}_n = \mathbf{0}^1$  (gdzie  $n = 1, \dots, N$ ), przyjmuje postać

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}, \quad (1.10)$$

gdzie

- $\mathbf{Y}$  jest wektorem zmiennych objaśnianych długości  $N$ ,
- $\mathbf{X}$  jest macierzą zmiennych objaśniających wymiaru  $N \times p$ ,
- $\boldsymbol{\beta}$  jest wektorem parametrów długości  $p$ ,
- $\boldsymbol{\varepsilon}$  jest wektorem błędów losowych długości  $N$ .

[ Mając dane realizacje zmiennych  $\mathbf{Y}$  oraz  $\mathbf{X}$  poszukiwany wektor współczynników modelu  $\boldsymbol{\beta}$  znajdujemy metodą najmniejszych kwadratów. ]

Poza narzuconym już założeniem o scentrowanych zmiennych losowych  $\mathbf{Y}$  i  $\mathbf{X}$  (tu: jedynie aby uniknąć uwzględniania wyrazu wolnego<sup>2</sup>) najważniejszymi założeniami powyższego modelu liniowego są wymagania, aby zmienna losowa  $\boldsymbol{\varepsilon}$  opisująca błąd modelu również spełniała  $\mathbb{E}[\boldsymbol{\varepsilon}] = 0$  oraz aby nie była skorelowana ze zmiennymi  $X_n$ .

Rozważać będziemy odpowiednik modelu liniowego dla zmiennych funkcyjnych. Dla uproszczenia (podobnie jak wyżej) zakładamy, że zmienne objaśniane i objaśniające mają średnie równe zero. **Pełen model funkcyjny** (ang. *fully functional model*) przyjmuje postać

$$Y_n = \Psi X_n + \varepsilon_n, \quad n = 1, 2, \dots, N, \quad (1.11)$$

gdzie krzywe  $Y_n$ ,  $X_n$  oraz nieobserwowalny błąd  $\varepsilon_n$  należą do przestrzeni Hilberta  $L^2(T)$ . Operator  $\Psi : L^2 \rightarrow L^2$  jest ograniczonym operatorem liniowym, który jest operatorem całkowym. Jądro całkowe  $\psi(t, s)$  operatora  $\Psi$  jest funkcją całkowalną z kwadratem na  $T \times T$ .

<sup>1</sup>przenieść tę uwagę/wytłumaczenie do przypisu?

<sup>2</sup>przenieść tę uwagę/wytłumaczenie do przypisu?

Zauważmy ponadto, że, na mocy Uwagi 1.2.2, operator  $\Psi$  jest operatorem Hilberta-Schmidta. Równość (1.11) rozumiemy zatem następująco

$$Y_n(t) = \int \psi(t, s)X_n(s)ds + \varepsilon_n(t), \quad n = 1, 2, \dots, N. \quad (1.12)$$

Jak i w przypadku standardowego modelu liniowego, funkcjonalny model liniowy wymusza pewne założenia. Podobnie jak poprzednio, wymagamy, aby zmienna losowa  $\varepsilon_n$  opisująca błąd modelu spełniała  $\mathbb{E}[\varepsilon_n] = 0$  oraz aby nie była skorelowana ze zmiennymi  $X_n$ .

[”nieskorelowane zmienne” = ( operator kowariancji = 0 )?]

[inne założenia modelu? konsekwencje?]

[przykład - nawet jeśli nie zapisywać, to mieć w głowie]

Nazwa powyższego modelu wynika z faktu, że zarówno zmienne objaśniane  $Y_n$  jak i zmienne objaśniające  $X_n$  są zmiennymi funkcjonalnymi. Niewielkim uproszczeniem są pozostałe typy funkcjonalnych modeli liniowych, tj.

- model z odpowiedzią skalarną (ang. *scalar response model*)

$$Y_n = \int \psi(s)X_n(s)ds + \varepsilon_n, \quad n = 1, 2, \dots, N,$$

w którym tylko zmienne objaśniające  $X_n$  są zmiennymi funkcjonalnymi,

[przykład - nawet jeśli nie zapisywać, to mieć w głowie]

- model z odpowiedzią funkcyjną (ang. *functional response model*)

$$Y_n(t) = \psi(t)x_n + \varepsilon_n(t), \quad n = 1, 2, \dots, N,$$

w którym zmienne objaśniające  $x_n$  są deterministycznymi skalarami.

[przykład - nawet jeśli nie zapisywać, to mieć w głowie]

Naturalnym problemem pojawiającym się przy funkcjonalnym modelu liniowym jest estymacja operatora  $\Psi$  należącego do nieskończonej wymiarowej przestrzeni na podstawie skończonej próbki danych. Możliwym jest znalezienie operatora, który daje idealne dopasowanie do danych (dla którego wszystkie różnice od próbki są równe zero), nie narzucając dodatkowych założeń, ale przypomina on biały szum i jego interpretacja jest często problemowa i nie funkcjonalna. Jednym ze sposobów na rozwiązanie tego problemu jest poszukiwanie operatora należącego do podprzestrzeni generowanej przez funkcje własne operatora kowariancji danych z próby, nazywane **empirycznymi funkcjonalnymi głównymi składowymi** (ang. *empirical functional principal components, EFPC's*), które zostały opisane w podrozdziale 1.5. Główne składowe odpowiadają istotnym czynnikom zmienności zmiennych, dobrze służą zatem do przybliżania ich wartości.

...

[sposób znalezienia  $\Psi$ ]

Wykorzystany w dalszej części pracy pakiet *fda*, do programu *R-project*, do znalezienia operatora  $\Psi$  stosuje metodę najmniejszych kwadratów. Dlatego właśnie tę metodę przedstawiamy poniżej.

Niech  $\{\eta_k\}_{k=1}^\infty$  i  $\{\theta_l\}_{l=1}^\infty$  będą pewnymi ustalonymi bazami, niekoniecznie ortonormalnymi, np. bazami Fouriera lub splajnowymi. Ponadto, niech funkcje  $\eta_k$  dobrze przybliżają funkcje  $X_n$ , a  $\theta_l$  dobrze przybliżają  $Y_n$ . Nieznane jądro  $\psi$  estymujemy według postaci

$$\hat{\psi}(t, s) = \sum_{k=1}^K \sum_{l=1}^L p_{kl} \eta_k(s) \theta_l(t),$$



gdzie  $K$  i  $L$  są odpowiednio małymi liczbami wybranymi do wygładzenia przybliżenia  $X_n$  i  $Y_n$ . Podobnie jak w przypadku standardowego modelu linowego możemy znaleźć parametry  $p_{kl}$  metodą najmniejszych kwadratów przez minimalizację sumy kwadratów reszt

$$\sum_{n=1}^N \left\| Y_n - \int X_n(s) \hat{\psi}(s, \cdot) \right\|^2.$$

[jak w pakiecie w R] (1.12)



## Rozdział 2

# Test istotności w funkcjonalnym modelu liniowym

### 2.1. Procedura testowa

Jednym z podstawowych testów na efektywność modelu jest test istotności zmiennych objaśniających. Jak w przypadku modelu liniowego dla zmiennych skalarnych (postaci (1.10)) testuje się hipotezę o zerowaniu się wektora  $\beta$ , tak w przypadku funkcjonalnego modelu liniowego badamy zerowanie się operatora  $\Psi$ , tj. hipotezy

$$H_0 : \quad \Psi = 0 \quad \text{przeciw} \quad H_A : \quad \Psi \neq 0.$$

Zauważmy, że przyjęcie  $H_0$  nie oznacza braku związku między zmienną objaśnianą a objaśniającą. Prowadzi jedynie do stwierdzenia braku zależności liniowej.

Obserwujemy ciąg krzywych długości  $N$ . Zakładamy, że zmienna objaśniana  $Y_n$ , zmienne objaśniające  $X_n$  i błędy  $\varepsilon_n$  są scentrowanymi zmiennymi losowymi przyjmującymi wartości w przestrzeni Hilberta  $L^2$ . Oznaczając przez  $X$  (analogicznie  $Y$ ) zmienną funkcjonalną o tym samym rozkładzie co  $X_n$  ( $Y_n$ ) wprowadzamy operatory kowariancji

$$C(x) = \mathbb{E}[\langle X, x \rangle X], \quad \Gamma(x) = \mathbb{E}[\langle Y, x \rangle Y], \quad \Delta(x) = \mathbb{E}[\langle X, x \rangle Y], \quad x \in L^2. \quad (2.1)$$

Przez  $\hat{C}$ ,  $\hat{\Gamma}$ ,  $\hat{\Delta}$  oznaczamy ich estymatory (zgodnie z (1.9)), tj.

$$\hat{C}(x) = \frac{1}{N} \sum_{n=1}^N \langle X_n, x \rangle X_n, \quad \hat{\Gamma}(x) = \frac{1}{N} \sum_{n=1}^N \langle Y_n, x \rangle Y_n, \quad \hat{\Delta}(x) = \frac{1}{N} \sum_{n=1}^N \langle X_n, x \rangle Y_n, \quad x \in L^2.$$

Definiujemy również wartości i wektory własne  $C$  i  $\Gamma$

$$C(v_k) = \lambda_k v_k, \quad \Gamma(u_j) = \gamma_j u_j, \quad (2.2)$$

których estymatory będziemy oznaczać  $(\hat{\lambda}_k, \hat{v}_k)$ ,  $(\hat{\gamma}_j, \hat{u}_j)$ .

Test obejmuje obcięcie powyższych operatorów na podprzestrzeń skończenie wymiarowe. Podprzestrzeń  $\mathcal{V}_p = \text{span}\{v_1, \dots, v_p\}$  zawiera najlepsze przybliżenia  $X_n$ , które są liniowymi kombinacjami pierwszych  $p$  głównych składowych (ang. *Functional Principal Components, FPC*). Metodą głównych składowych wyznaczamy  $p$  największych wartości własnych operatora  $\hat{C}$  tak, że  $\hat{\mathcal{V}}_p = \text{span}\{\hat{v}_1, \dots, \hat{v}_p\}$  zawiera najlepsze przybliżenie  $X_n$ . Analogicznie  $\mathcal{U}_q = \text{span}\{u_1, \dots, u_q\}$  zawiera przybliżenia  $\text{span}\{Y_1, \dots, Y_N\}$ .

Z ogólnej postaci funkcjonalnego modelu liniowego

$$Y = \Psi X + \varepsilon$$

możemy wyprowadzić kolejne równości

$$\begin{aligned}\langle X, x \rangle Y &= \langle X, x \rangle \Psi X + \langle X, x \rangle \varepsilon \\ \mathbb{E}[\langle X, x \rangle Y] &= \mathbb{E}[\langle X, x \rangle \Psi X] + \mathbb{E}[\langle X, x \rangle \varepsilon].\end{aligned}$$

Korzystając z definicji operatorów  $C$  oraz  $\Delta$  (2.1), założenia, że  $\Psi$  jest operatorem ograniczonym oraz z założenia o braku korelacji między  $X$  a  $\varepsilon$  zachodzi

$$\Delta = \Psi C.$$

W szczególności, prawdziwa jest równość

$$\Delta(v_k) = \Psi C(v_k).$$

Na mocy definicji funkcji własnych (2.2), dla  $k \leq p$ , mamy

$$\Psi(v_k) = \lambda_k^{-1} \Delta(v_k).$$

Stąd,  $\psi$  zeruje się na  $\text{span}\{v_1, \dots, v_p\}$  wtedy i tylko wtedy, gdy  $\Delta(v_k) = 0$  dla każdego  $k = 1, \dots, p$ . Zauważmy, że

$$\Delta(v_k) \approx \hat{\Delta}(v_k) = \frac{1}{N} \sum_{n=1}^N \langle X_n, v_k \rangle Y_n.$$

Skoro zatem  $\text{span}\{Y_1, \dots, Y_N\}$  są dobrze aproksymowane przez  $\mathcal{U}_q$ , to możemy ograniczyć się do sprawdzania czy

$$\langle \hat{\Delta}(v_k), u_j \rangle = 0, \quad k = 1, \dots, p, \quad j = 1, \dots, q. \quad (2.3)$$

Jeśli  $H_0$  jest prawdziwa, to dla każdego  $x \in \mathcal{V}_p$ ,  $\psi(x)$  nie należy do  $\mathcal{U}_q$ . Co znaczy, że żadna funkcja  $Y_n$  nie może być opisana jako liniowa kombinacja  $X_n$ ,  $n = 1, \dots, N$ . Statystyka testowa powinna zatem sumować kwadraty iloczynów skalarnych (2.3). Twierdzenie 2.2.1 stanowi, że statystyka

$$\hat{T}_N(p, q) = N \sum_{k=1}^p \sum_{j=1}^q \hat{\lambda}_k^{-1} \hat{\gamma}_j^{-1} \langle \hat{\Delta}(\hat{v}_k), \hat{u}_j \rangle^2, \quad (2.4)$$

zbiega według rozkładu do rozkładu  $\chi^2$  z  $pq$  stopniami swobody.

Przy czym

$$\langle \hat{\Delta}(\hat{v}_k), \hat{u}_j \rangle = \left\langle \frac{1}{N} \sum_{n=1}^N \langle X_n, \hat{v}_k \rangle Y_n, \hat{u}_j \right\rangle = \frac{1}{N} \sum_{n=1}^N \langle X_n, \hat{v}_k \rangle \langle Y_n, \hat{u}_j \rangle$$

oraz  $\lambda_k = \mathbb{E} \langle X, v_k \rangle^2$  i  $\gamma_j = \mathbb{E} \langle Y, u_j \rangle^2$ .

**Uwaga 2.1.1** Oczywistym jest, że jeśli odrzucamy  $H_0$ , to  $\psi(v_k) \neq 0$  dla pewnego  $k \geq 1$ . Jednak ograniczając się do  $p$  największych wartości własnych, test jest skuteczny tylko jeśli  $\psi$  nie zanika na którymś wektorze  $v_k$ ,  $k = 1, \dots, p$ . Aczkolwiek takie ograniczenie jest intuicyjnie niegroźne, ponieważ test ma za zadanie sprawdzić czy główne źródła zmienności  $Y$  mogą być opisane przez główne źródła zmienności zmiennych  $X$ .

## Schemat przebiegu testu

1. Sprawdzamy założenie o liniowości metodą *FPC score predictor-response plots*.
2. Wybieramy liczbę głównych składowych  $p$  i  $q$  metodami *scree test* oraz *CPV*.
3. Wyliczamy wartość statystyki  $\hat{T}_N(p, q)$  (2.4).
4. Jeśli  $\hat{T}_N(p, q) > \chi_{pq}^2(1 - \alpha)$ , to odrzucamy hipotezę zerową o braku liniowej zależności. W przeciwnym razie nie mamy podstaw do odrzucenia  $H_0$ .

[rozwinąć i dopracować powyższe punkty]

## 2.2. Formalne podstawy [nazwa?]

**Założenie 2.2.1** [Kokoszka et al. (2008)], [Horváth, Kokoszka]

Trójka  $(Y_n, X_n, \varepsilon_n)$  tworzy ciąg niezależnych zmiennych funkcjonalnych o jednakowym rozkładzie, takich że  $\varepsilon_n$  jest niezależne od  $X_n$  oraz

$$\mathbb{E}X_n = 0, \quad \mathbb{E}\varepsilon_n = 0,$$

$$\mathbb{E}\|X_n\|^4 < \infty \quad i \quad \mathbb{E}\|\varepsilon_n\|^4 < \infty.$$

**Założenie 2.2.2** [Kokoszka et al. (2008)], [Horváth, Kokoszka]

Wartości własne operatorów  $C$  oraz  $\Gamma$  spełniają, dla pewnych  $p > 0$  i  $q > 0$

$$\lambda_1 > \lambda_2 > \dots > \lambda_p > \lambda_{p+1}, \quad \gamma_1 > \gamma_2 > \dots > \gamma_q > \gamma_{q+1}.$$

**Twierdzenie 2.2.1** [Kokoszka et al. (2008)], [Horváth, Kokoszka]

Jeśli spełnione są powyższe Założenia 2.2.1, 2.2.2 oraz  $H_0$ , to  $\hat{T}_N(p, q) \xrightarrow{d} \chi_{pq}^2$  przy  $N \rightarrow \infty$ .

**Twierdzenie 2.2.2** [Kokoszka et al. (2008)], [Horváth, Kokoszka]

Przy Założeniach 2.2.1, 2.2.2 oraz jeśli  $\langle \psi(v_k), u_j \rangle \neq 0$  dla  $k \leq p$  oraz  $j \leq q$ , to  $\hat{T}_N(p, q) \xrightarrow{P} \infty$  przy  $N \rightarrow \infty$ .

Dowody powyższych twierdzeń rozbijemy w krokach na kolejne lematy i wnioski. ...

Zauważmy, że konsekwencją prawdziwości  $H_0$  i przyjęcia modelu postaci  $Y_n = \Psi X_n + \varepsilon_n$  jest równość  $Y_n = \varepsilon_n$ . ?

**Lemat 2.2.1** [Kokoszka et al. (2008)], [Bosq]

Według oznaczeń podrozdziału 1.5, przy Założeniach 2.2.1, 2.2.2 spełnione są nierówności

$$\limsup_{N \rightarrow \infty} N \mathbb{E} \|v_k - \hat{v}_k\|^2 < \infty, \quad \limsup_{N \rightarrow \infty} N \mathbb{E} \|u_j - \hat{u}_j\|^2 < \infty,$$

$$\limsup_{N \rightarrow \infty} N \mathbb{E} \left[ |\gamma_k - \hat{\gamma}_k|^2 \right] < \infty, \quad \limsup_{N \rightarrow \infty} N \mathbb{E} \left[ |\lambda_j - \hat{\lambda}_j|^2 \right] < \infty,$$

dla  $k \leq p$  oraz  $j \leq q$ .

[

**Twierdzenie 2.2.3** *Centralne Twierdzenie Graniczne [Horváth, Kokoszka], [Bosq]*

Niech  $\{X_n\}_{n \geq 1}$  będzie ciągiem zmiennych funkcyjnych o jednakowym rozkładzie przyjmujących wartości w ośrodkowej przestrzeni Hilberta. Jeśli  $\mathbb{E}\|X_1\|^2 < \infty$ ,  $\mathbb{E}X_1 = \mu$  i  $C_{X_1} = C$ , wtedy

$$N^{-1/2} \sum_{n=1}^N X_n \xrightarrow{d} \mathcal{N},$$

gdzie  $\mathcal{N} \sim \mathcal{N}(0, C)$ .

]

**Lemat 2.2.2** [Kokoszka et al. (2008)], [Horváth, Kokoszka]

Jeśli spełnione są Założenia 2.2.1, 2.2.2 i  $H_0$ , to dla  $k \leq p$ ,  $j \leq q$

$$\sqrt{N} \langle \hat{\Delta} v_k, u_j \rangle \xrightarrow{d} \eta_{kj} \sqrt{\gamma_k \lambda_j}, \quad (2.5)$$

gdzie  $\eta_{kj} \sim N(0, 1)$ . Przy czym  $\eta_{kj}$  oraz  $\eta_{k'j'}$  są niezależne dla  $(k, j) \neq (k', j')$ .

Dowód. Przy  $H_0$

$$\sqrt{N} \langle \hat{\Delta} v_k, u_j \rangle = N^{-1/2} \sum_{n=1}^N \langle X_n, v_k \rangle \langle \varepsilon_n, u_j \rangle,$$

gdzie elementy pod sumą po prawej stronie powyższej równości mają średnie 0 i wariancje równe  $\lambda_k \gamma_j$ , co na mocy CTG (Twierdzenie 2.2.3) kończy dowód (2.5). [skalarne CTG?]

Aby udowodnić niezależność między  $\eta_{kj}$  i  $\eta_{k'j'}$  dla  $(k, j) \neq (k', j')$ , wystarczy pokazać, że  $\sqrt{N} \langle \hat{\Delta}(v_k), u_j \rangle$  i  $\sqrt{N} \langle \hat{\Delta}(v_{k'}), u_{j'} \rangle$  są nieskorelowane. Mamy

$$\begin{aligned} & \mathbb{E} \left[ \sqrt{N} \langle \hat{\Delta}(v_k), u_j \rangle, \sqrt{N} \langle \hat{\Delta}(v_{k'}), u_{j'} \rangle \right] \\ &= N \mathbb{E} \left[ \left\langle \frac{1}{N} \sum_{n=1}^N \langle X_n, v_k \rangle Y_n, u_j \right\rangle, \left\langle \frac{1}{N} \sum_{n'=1}^N \langle X_{n'}, v_{k'} \rangle Y_{n'}, u_{j'} \right\rangle \right] \\ &= N \mathbb{E} \left[ \left\langle \frac{1}{N} \sum_{n=1}^N \langle X_n, v_k \rangle (\Psi X_n + \varepsilon_n), u_j \right\rangle, \left\langle \frac{1}{N} \sum_{n'=1}^N \langle X_{n'}, v_{k'} \rangle (\Psi X_{n'} + \varepsilon_{n'}), u_{j'} \right\rangle \right] \\ &\stackrel{H_0}{=} \frac{1}{N} \mathbb{E} \left[ \sum_{n=1}^N \langle X_n, v_k \rangle \langle \varepsilon_n, u_j \rangle \sum_{n'=1}^N \langle X_{n'}, v_{k'} \rangle \langle \varepsilon_{n'}, u_{j'} \rangle \right] \\ &= \frac{1}{N} \sum_{n, n'=1}^N \mathbb{E} [\langle X_n, v_k \rangle \langle X_{n'}, v_{k'} \rangle] \mathbb{E} [\langle \varepsilon_n, u_j \rangle \langle \varepsilon_{n'}, u_{j'} \rangle] \\ &= \frac{1}{N} \sum_{n=1}^N \mathbb{E} [\langle X_n, v_k \rangle \langle X_n, v_{k'} \rangle] \mathbb{E} [\langle \varepsilon_n, u_j \rangle \langle \varepsilon_n, u_{j'} \rangle] \\ &= \langle C(v_k), v_{k'} \rangle \langle \Gamma u_j, u_{j'} \rangle = \gamma_k \delta_{kk'} \gamma_j \delta_{jj'}. \end{aligned}$$

[zastanowić się nad tym/dopracować]

□

Przypomnijmy, że norma Hilberta-Schmidta operatora Hilberta-Schmidta  $S$  zdefiniowana jest wzorem  $\|S\|_S^2 = \sum_{j=1}^{\infty} \|S(e_j)\|^2$ , gdzie ciąg  $\{e_1, e_2, \dots\}$  stanowi bazę ortonormalną oraz, że norma ta jest nie mniejsza od normy operatorowej, tj.  $\|S\|_{\mathcal{L}}^2 \leq \|S\|_S^2$ .

**Lemat 2.2.3** [Kokoszka et al. (2008)], [Horváth, Kokoszka]  
 Przy założeniach Twierdzenia 2.2.1 mamy

$$\mathbb{E} \left\| \widehat{\Delta} \right\|_{\mathcal{S}}^2 = N^{-1} \mathbb{E} \|X\|^2 \mathbb{E} \|\varepsilon_1\|^2.$$

*Dowód.* Zauważmy, że

$$\left\| \widehat{\Delta}(e_j) \right\|^2 = N^{-2} \sum_{n,n'=1}^N \langle X_n, e_j \rangle \langle X_{n'}, e_j \rangle \langle Y_n, Y_{n'} \rangle.$$

Stąd, przy założeniu  $H_0$ , mamy

$$\begin{aligned} \mathbb{E} \left\| \widehat{\Delta} \right\|_{\mathcal{S}}^2 &= N^{-2} \sum_{j=1}^{\infty} \sum_{n,n'=1}^N \mathbb{E} [\langle X_n, e_j \rangle \langle X_{n'}, e_j \rangle \langle \varepsilon_n, \varepsilon_{n'} \rangle] \\ &= N^{-2} \sum_{j=1}^{\infty} \sum_{n=1}^N \mathbb{E} \langle X_n, e_j \rangle^2 \mathbb{E} \|\varepsilon_n\|^2 \\ &= N^{-1} \mathbb{E} \|\varepsilon_1\|^2 \sum_{j=1}^{\infty} \langle X, e_j \rangle^2 = N^{-1} \mathbb{E} \|\varepsilon_1\|^2 \|X\|^2. \end{aligned}$$

□

**Lemat 2.2.4** [Kokoszka et al. (2008)], [Horváth, Kokoszka]

Założmy, że  $\{U_n\}_{n=1}^{\infty}$  oraz  $\{V_n\}_{n=1}^{\infty}$  są ciągami elementów losowych z przestrzeni Hilberta takich, że  $\|U_n\| \xrightarrow{P} 0$  i  $\|V_n\| = O_P(1)$ , tj.

$$\lim_{C \rightarrow \infty} \limsup_{n \rightarrow \infty} P(\|V_n\| > C) = 0.$$

Wtedy zachodzi

$$\langle U_n, V_n \rangle \xrightarrow{P} 0.$$

*Dowód.* Prawdziwość lematu wynika z analogicznej własności dla losowych ciągów liczb rzeczywistych i nierówności  $|\langle U_n, V_n \rangle| \leq \|U_n\| \|V_n\|$ .

[może lepiej przytoczyć skalarną wersję?]

□

**Lemat 2.2.5** [Kokoszka et al. (2008)], [Horváth, Kokoszka]

Przy założeniach Twierdzenia 2.2.1, dla  $k \leq p$ ,  $j \leq q$  zachodzi

$$\sqrt{N} \langle \widehat{\Delta}(\hat{v}_k), \hat{u}_j \rangle \xrightarrow{d} \eta_{kj} \sqrt{\lambda_k \gamma_j},$$

gdzie  $\eta_{kj}$  definiowane są jak w Lemacie 2.2.2.

*Dowód.* Na mocy Lematu 2.2.2, wystarczy pokazać, że

$$\sqrt{N} \langle \widehat{\Delta}(\hat{v}_k), \hat{u}_j \rangle - \sqrt{N} \langle \widehat{\Delta}(v_k), u_j \rangle \xrightarrow{P} 0. \quad (2.6)$$

Równość (2.6) wynika z nierówności trójkąta oraz z

$$\sqrt{N} \langle \widehat{\Delta}(\hat{v}_k), \hat{u}_j - u_j \rangle \xrightarrow{P} 0 \quad (2.7)$$

i

$$\sqrt{N} \langle \widehat{\Delta}(\hat{v}_k - v_k), \hat{u}_j \rangle \xrightarrow{P} 0. \quad (2.8)$$

Aby udowodnić równość (2.7), zauważmy, że z Lematu 2.2.1 mamy  $\sqrt{N}(\hat{u}_j - u_j) = O_P(1)$  oraz, na mocy Lematu 2.2.3,  $\mathbb{E}\|\hat{\Delta}(v_k)\| \leq \mathbb{E}\|\hat{\Delta}\|_{\mathcal{S}} = O(N^{-1/2})$ . Stąd równość (2.7) wynika z Lematu 2.2.4.

Aby wykorzystać takie samo uzasadnienie dla (2.8) (skorzystać z Lematu 2.2.1), zauważmy, że

$$\sqrt{N}\langle \hat{\Delta}(\hat{v}_k - v_k), \hat{u}_j \rangle = \sqrt{N}\langle \hat{v}_k - v_k, \tilde{\Delta}(\hat{u}_j) \rangle,$$

gdzie  $\tilde{\Delta}(x) = N^{-1} \sum_{n=1}^N \langle Y_n, x \rangle X_n$ . Lemat 2.2.3 stanowi, że przy założeniu  $H_0$  mamy  $\mathbb{E}\|\tilde{\Delta}\|_{\mathcal{S}} = \mathbb{E}\|\hat{\Delta}\|_{\mathcal{S}}$ , co kończy dowód.

[na pewno?]

□

Z Lematu 2.2.1,  $\hat{\lambda}_k \xrightarrow{P} \lambda_k$  oraz  $\hat{\gamma}_j \xrightarrow{P} \gamma_j$ .

**Wniosek 2.2.1** [Kokoszka et al. (2008)], [Horváth, Kokoszka]  
Przy założeniach Twierdzenia 2.2.1, dla  $j \leq q$ ,  $k \leq p$  zachodzi

$$\sqrt{N}\langle \hat{\lambda}_k^{-1/2} \hat{\gamma}_j^{-1/2} \hat{\Delta}(\hat{v}_k), \hat{u}_j \rangle \xrightarrow{d} \eta_{kj},$$

gdzie  $\eta_{kj}$  definiowane są jak w Lemacie 2.2.2.

Dowód Twierdzenia 2.2.1 [...]

**Lemat 2.2.6** [Kokoszka et al. (2008)], [Horváth, Kokoszka]  
Jeśli  $\{Y_n\}_{n \geq 1}$  są zmiennymi funkcjonalnymi o jednakowych rozkładach, to zachodzi

$$\mathbb{E}\|\hat{\Delta}\| \leq \mathbb{E}\|Y\|^2.$$

Dowód. Dla dowolnego  $u \in L^2$  takiego, że  $\|u\| \leq 1$ , mamy

$$\|\hat{\Delta}u\| \leq \frac{1}{N} \sum_{n=1}^N |\langle Y_n, u \rangle| \|Y_n\| \leq \frac{1}{N} \sum_{n=1}^N \|Y_n\|^2.$$

Co ze względu na założenie, że  $Y_n$  mają jednakowy rozkład, jest równoważne tezie lematu. □

**Twierdzenie 2.2.4** Mocne Prawo Wielkich Liczb [Bosq]

Niech  $\{X_n\}_{n \geq 1}$  będzie ciągiem zmiennych funkcjonalnych o jednakowym rozkładzie przyjmujących wartości w ośrodkowej przestrzeni Hilberta takich, że  $\mathbb{E}\|X_n\|^2 < \infty$ . Niech  $m = \mathbb{E}X_n$ , wtedy mamy

$$\frac{1}{N} \sum_{n=1}^N X_n \xrightarrow{p.n.} m.$$

**Lemat 2.2.7** [Kokoszka et al. (2008)], [Horváth, Kokoszka]  
Jeżeli spełnione jest Założenie 2.2.1, to dla dowolnych funkcji  $v, u \in L^2$

$$\langle \hat{\Delta}(v), u \rangle \xrightarrow{P} \langle \Delta(v), u \rangle.$$

Dowód. Tezę otrzymujemy korzystając z Prawa Wielkich Liczb zauważając

$$\langle \hat{\Delta}(v), u \rangle = \frac{1}{N} \sum_{n=1}^N \langle X_n, v \rangle \langle Y_n, u \rangle$$

oraz

$$\mathbb{E}[\langle X_n, v \rangle \langle Y_n, u \rangle] = \mathbb{E}[\langle \langle X_n, v \rangle Y_n, u \rangle] = \langle \Delta(v), u \rangle.$$

□



**Lemat 2.2.8** [Kokoszka et al. (2008)], [Horváth, Kokoszka]  
*Jeżeli spełnione są Założenia 2.2.1 oraz 2.2.2, to*

$$\langle \hat{\Delta}(\hat{v}_k), \hat{u}_j \rangle \xrightarrow{P} \langle \Delta(v_k), u_j \rangle, \quad \text{dla } k \leq p, \ j \leq q.$$

*Dowód.* Na mocy Lematu 2.2.7 wystarczy pokazać

$$\langle \hat{\Delta}(v_k), \hat{u}_j - u_j \rangle \xrightarrow{P} 0$$

i

$$\langle \hat{\Delta}(\hat{v}_k) - \hat{\Delta}(v_k), \hat{u}_j \rangle \xrightarrow{P} 0.$$

Relacje te wynikają z Lematów 2.2.4, 2.2.1 [na pewno?] oraz 2.2.6. □

*Dowód Twierdzenia 2.2.2.* Wprowadźmy oznaczenie

$$\hat{S}_N(p, q) = \sum_{k=1}^p \sum_{j=1}^q \hat{\lambda}_k^{-1} \hat{\gamma}_j^{-1} \langle \hat{\Delta}(\hat{v}_k), \hat{u}_j \rangle^2.$$

Na mocy Lematu 2.2.8 oraz Lematu 2.2.1 [na pewno?], zachodzi

$$\hat{S}_N(p, q) \xrightarrow{P} S(p, q) > 0.$$

Stąd,

$$\hat{T}_N(p, q) = N \hat{S}_N(p, q) \xrightarrow{P} \infty.$$

□

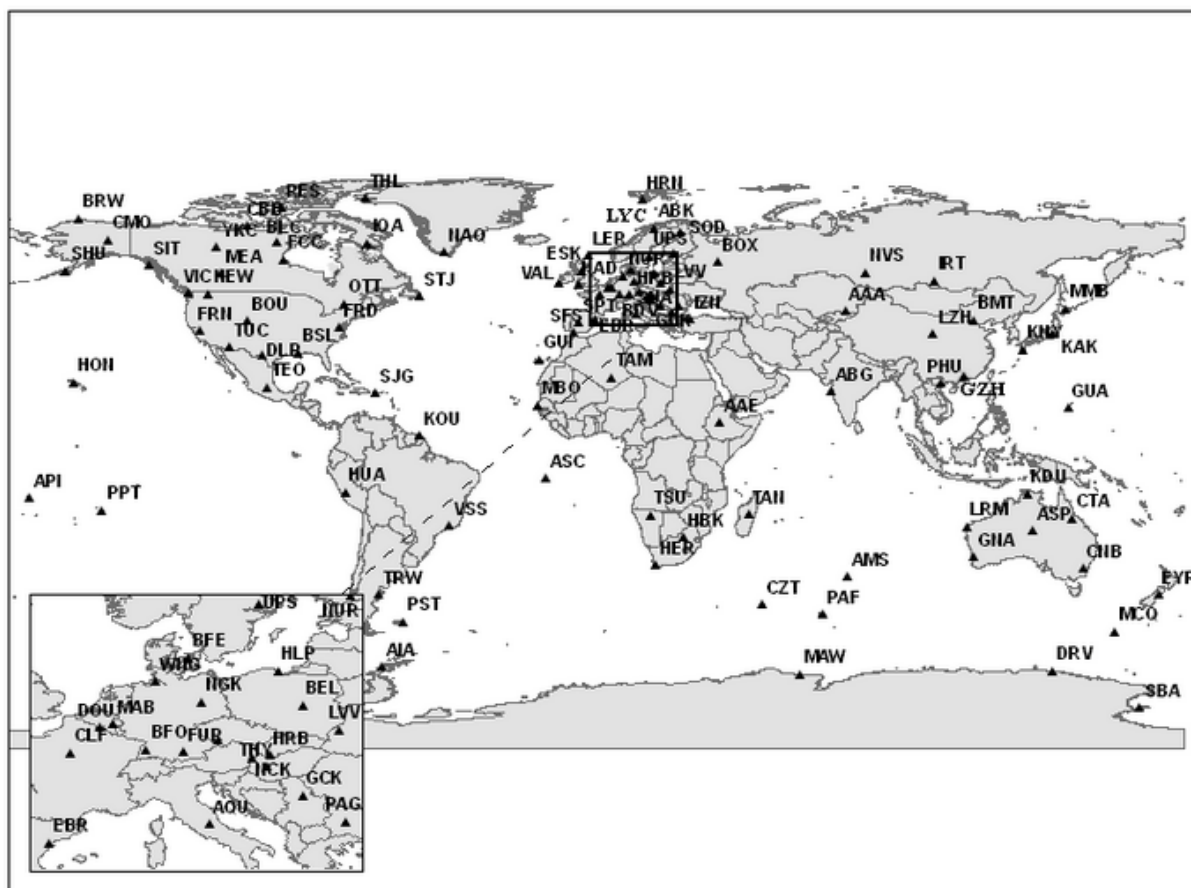
[...?]



## Rozdział 3

# Przykład zastosowania

Podobnie jak w artykule [Kokoszka et al. (2008)] oraz książce [Horváth, Kokoszka], zastosujemy przedstawiony test do modelu stworzonego na podstawie danych opisujących pola magnetyczne. Takie dane zbierane są przez stacje geofizyczne i publikowane są w ramach międzynarodowego programu INTERMAGNET i udostępniane są na stronie internetowej projektu [I]. Do programu należy obecnie 129 naziemnych obserwatoriów, w tym dwie stacje znajdujące się w Polsce (mapa stacji na Rysunku 3.1)



Rysunek 3.1: Mapa stacji geofizycznych należących do programu INTERMAGNET

Dane o polu magnetycznym zbierane są za pomocą tzw. magnetometru. Urządzenie to od-

czytuje informacje... [...]

Ze strony programu można pobrać dane dokładne: w odstępach jednosekundowych lub uproszczone: w odstępach jednodominutowych (obserwacja jest średnią z 60 sekund). W pracy wykorzystano dane uproszczone, mamy zatem 1440 punktów każdego dnia (poza sytuacjami z brakiem części danych), które wykorzystamy do stworzenia danych funkcjonalnych. Tym sposobem jeden dzień stanie się jedną obserwacją.

Magnetometer data...

Zorze polarne + substorms (?)

Korzystając z dostępnego pakietu *fda* ([R: fda])...

### 3.1. Ameryka Północna (Kanada)

W kręgu zainteresowań autorów artykułu [Kokoszka et al. (2008)] leżą dane pochodzące z obserwatoriów Ameryki Północnej, zaczniemy zatem od analizy podobnych danych.

[do opracowania: punkt po punkcie według opisu procedury testowej w rozdziale 2]

[do opracowania: kod w R!]

### 3.2. Europa (Polska)

Do programu INTERMAGNET należą także dwie polskie stacje geofizyczne: obserwatorium w Belsku oraz obserwatorium na Helu. Przeprowadzimy zatem podobną j.w. analizę dla północnej Europy.

[...]

**Dodatek A**

**Kod w R**

[KOD] [Bosq]



# Bibliografia

- [Bosq] D. Bosq, *Linear Processes in Function Spaces*. Springer, 2000.
- [Ferraty, Vieu] F. Ferraty, P. Vieu, *Nonparametric Functional Data Analysis. Theory and practice*. Springer, 2006.
- [Horváth, Kokoszka] L. Horváth, P. Kokoszka, *Interference for Functional Data with Applications*. Springer, 2012.
- [I] INTERMAGNET <http://www.intermagnet.org/index-eng.php>
- [Kokoszka et al. (2008)] P. Kokoszka, I. Maslova, J. Sojka, L. Zhu, *Testing for lack of dependence in the functional linear model*. Canadian Journal of Statistics, 2008, 36, 207-222.
- [Maslova et al. (2010)] Maslova, I., Kokoszka, P., Sojka, J. and Zhu, L., *Statistical significance testing for the association of magnetometer records at high-, mid- and low latitudes during substorm days*. Planetary and Space Science, 58 (2010), 437-445.
- [Pytlik] T. Pytlik, *Analiza funkcjonalna*. Instytut Matematyczny Uniwersytetu Wrocławskiego, 2000.
- [R: fda] J. O. Ramsay, H. Wickham, S. Graves, G. Hooker, *Package 'fda'*, wersja 2.4.4. On-line: <https://cran.r-project.org/web/packages/fda/fda.pdf>
- [Ramsay, et al. (2009)] J. O Ramsay, G. Hooker and S. Graves, *Functional Data Analysis with R and Matlab*. Springer, 2009.
- [Ramsay, Silverman] J. O. Ramsay, B. W. Silverman, *Functional Data Analysis*. Springer, 2005.
- [Wojtaszczyk] P. Wojtaszczyk, *Banach Spaces For Analysts*. Cambridge Universiti Press, 1991, 86-87.