

Адаптация метода главных компонент к данным с пропусками

Рукавишникова Анна Александровна, гр. 422

Санкт-Петербургский государственный университет
Прикладная математика и информатика
Вычислительная стохастика и статистические модели

Научный руководитель: к.ф.-м.н., доцент Алексеева Н. П.
Рецензент: д.ф.-м.н., профессор Кривулин Н. К.



Санкт-Петербург
2018 г.

Таблица: Общий вид табличных данных

	y_1	y_2	\dots	y_j	\dots	y_n
1	x_{11}	*	\dots	x_{1j}	\dots	x_{1n}
2	*	x_{22}	\dots	*	\dots	x_{2n}
.	.	.	\dots	.	\dots	.
i	x_{i1}	*	\dots	x_{ij}	\dots	x_{in}
.	.	.	\dots	.	\dots	*
m	x_{m1}	x_{m2}	\dots	*	*	x_{mn}

Причины пропусков:

- невозможность получения или обработки данных,
- искажение или сокрытие информации,
- утеря данных.

Почему данные с пропусками являются проблемой?

Стандартные статистические методы предполагают, что все переменные в указанной модели измерены для всех случаев.

Идея метода: редукция размерности данных при наименьшей потере информативности.

Определение

Пусть имеется k центрированных ($EX_i = 0$) признаков $X = (X_1, \dots, X_k)^T$ для m индивидов, где $X_i \in \mathbb{R}^m$. Тогда j -ой **главной компонентой** называется линейная комбинация

$$Y_j = A_j^T X = \sum_{i=1}^k \alpha_{ij} X_i,$$

где A_j — собственные вектора ковариационной матрицы $\Sigma = EXX^T$, соответствующие собственным числам λ_j , упорядоченным по неубыванию: $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_k$.

Необходимость использования ортогональных компонент в качестве вещественной и мнимой части реализаций комплексного марковского нормального стационарного (КМНС) процесса $x_j(t) = u_j(t) + iv_j(t)$, $j = 1, \dots, n$, $t = 1, \dots, k$, где n — число наблюдений, k — число временных точек.

Замечание

Важным свойством метода главных компонент является ортогональность (независимость) главных компонент:

$$EY_m Y_n^T = 0 \text{ при } m \neq n.$$

Структура данных с пропусками по больным туберкулёзом лёгких:

- 30 пациентов, из которых:
 - 19 пациентов с улучшениями после лечения,
 - 11 пациентов без улучшений,
- 32 признака,
- 2 временные точки,
- 23% данных отсутствуют.

Таблица: Основные признаки (измерены в 2-х временных точках)

ММР.1 ММР.8 ММР.9	Белки, участвующие в процессах воспаления при туберкулёзе лёгких
TIMP1	Вещество, регулирующее активность ММР в тканях

- Заполнить пропуски в имеющихся данных.
- Применить метод главных компонент к полученным полным данным для каждой группы пациентов в отдельности.
- Использовать полученные главные компоненты для исследования корреляционной структуры данных и анализа динамики заболевания на основе модели КМНС процесса.

Метод заполнения пропусков Predictive Mean Matching

$\Omega = (\omega_1, \dots, \omega_n)$ — наблюдения по неполным данным:
зависимой переменной $\mathcal{Y}(\Omega) = (y(\omega_1), \dots, y(\omega_n))$,
независимых переменных $\mathcal{X} = (X_1(\Omega), \dots, X_p(\Omega))$.

Модель линейной регрессии $Y = X\beta + \varepsilon$, $\varepsilon \sim \mathcal{N}(0, \sigma^2)$

- $Y = \mathcal{Y}(\Omega_1)$, $\Omega_1 \subseteq \Omega$ — полные наблюдения, $\Omega_2 = \Omega \setminus \Omega_1$,
- $X = (X_{t_1}(\Omega_1), \dots, X_{t_k}(\Omega_1))$, $(t_1, \dots, t_k) \subseteq (1, \dots, p)$.
- $\hat{\beta}$, $\hat{\sigma}^2$ по МНК

Предикторы:

- $\hat{Y}(\Omega_1) = X\hat{\beta}$,
- $Y^*(\Omega_2) = X_*\beta^*$, где $X_* = (X_{t_1}(\Omega_2), \dots, X_{t_k}(\Omega_2))$,
 $\beta^* \sim \mathcal{N}(\hat{\beta}, \text{Cov}(\hat{\beta}))$.

$$\Delta(u, v) = |\hat{Y}(u \in \Omega_1) - Y^*(v \in \Omega_2)|$$

Подбор значений

$\forall v \in \Omega_2$ подбирается $y(v) = y(u) \in \mathcal{Y}(\Omega_1)$ при условии $\Delta(u, v) \leq \Delta_*$.

- Используемая среда программирования: R.
- Используемый программный пакет: MICE (Multivariate Imputation by Chained Equations).
- Метод заполнения пропусков: Predictive Mean Matching.
- Количество полученных полных наборов данных: 5.

Цель: для двух групп пациентов исследовать корреляционную структуру данных.

Подход к решению: оценка параметров τ и η ковариационной функции комплексного марковского нормального стационарного (КМНС) процесса (Дж.Л. Дуб, 1956)

$$\mathcal{B}(t) = \sigma^2 e^{-\eta|t| - i\tau t}, \quad \eta > 0.$$

Её выборочная оценка с m реализациями в k временных точках имеет вид:

$$\hat{\mathcal{B}}(t) = \frac{1}{(k-t)m} \sum_{l=1}^{k-t} \sum_{j=1}^m x_{jl} x_{j,l+t}^*.$$

Реализация

В качестве вещественной и мнимой части независимых реализаций КМНС процесса $x_j(t) = u_j(t) + iv_j(t) = x_{jt}$ берём значения первой и второй главной компоненты, соответственно.

Теорема (А.Г. Барт, 2003, Н.П. Алексеева, 2012)

Пусть $x_j(t) = u_j(t) + iv_j(t) = x_{jt}$ — m независимых реализаций КМНС процесса в k временных точках,

$$A_1 = \sum_{l=1}^k \sum_{j=1}^m x_{jl}^* x_{jl}, \quad A_2 = \sum_{l=2}^{k-1} \sum_{j=1}^m x_{jl}^* x_{jl}.$$

Если $A_1(k-2) = A_2k$, то ОМП $\hat{\tau}, \hat{\eta}, \hat{\sigma}$ удовлетворяют соотношениям:

$$\operatorname{tg} \hat{\tau} = \frac{\operatorname{Im}(\hat{\mathcal{B}}(1)/\hat{\mathcal{B}}(0))}{\operatorname{Re}(\hat{\mathcal{B}}(1)/\hat{\mathcal{B}}(0))}, \quad \hat{\eta} = -\ln \left| \frac{\hat{\mathcal{B}}(1)}{\hat{\mathcal{B}}(0)} \right|, \quad \hat{\sigma}^2 = \hat{\mathcal{B}}(0),$$

где $\hat{\mathcal{B}}(t)$ — выборочная оценка ковариационной функции КМНС процесса.

Замечание

Главные компоненты определены с точностью до знака.

Утверждение

Пусть \mathbf{u} и \mathbf{v} — первая и вторая главные компоненты соответственно. Тогда для оценок параметров КМНС процесса $\hat{\eta} = \hat{\eta}(\mathbf{u}, \mathbf{v})$ и $\hat{\tau} = \hat{\tau}(\mathbf{u}, \mathbf{v})$ верны следующие соотношения:

- $\hat{\eta}(\pm \mathbf{u}, \pm \mathbf{v}) = \hat{\eta}(\mathbf{u}, \mathbf{v})$,
- $\hat{\eta}(\mathbf{v}, \mathbf{u}) = \hat{\eta}(\mathbf{u}, \mathbf{v})$,
- $\hat{\tau}(-\mathbf{u}, \mathbf{v}) = \hat{\tau}(\mathbf{u}, -\mathbf{v}) = -\hat{\tau}(\mathbf{u}, \mathbf{v})$,
- $\hat{\tau}(\mathbf{v}, \mathbf{u}) = -\hat{\tau}(\mathbf{u}, \mathbf{v})$.

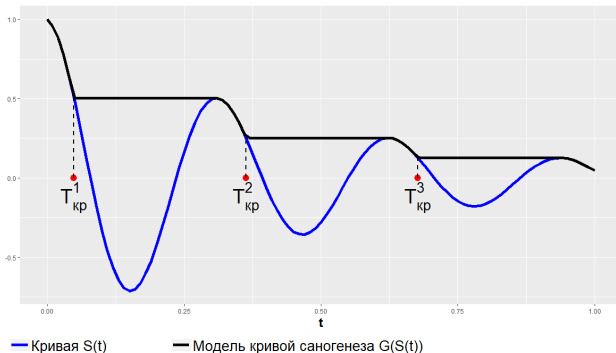
Параметрическая модель кривой саногенеза

Оценки $\hat{\tau}$ и $\hat{\eta}$ используются для построения вещественной части ковариационной функции КМНС процесса ($\sigma^2 = 1$)

$$S(t) = e^{-\eta t} \cos \tau t, \quad t > 0, \eta > 0.$$

Рассматриваем модель **кривой саногенеза** (А.Г. Барт, 2003)

$F(t) = G(S(t))$, полученную двойным правым обращением кривой $S(t)$, а именно: $F(t) = (S_{11}^-)_{11}(t)$, где $S_{11}^-(t) = \sup\{x : S(x) \geq t\}$.



— Кривая $S(t)$

— Модель кривой саногенеза $G(S(t))$

Важную роль в течении болезни играют **критические точки** $T_{кр}^j$ — моменты времени, в которые возможно изменение характера течения болезни.

Задача

Найти первую критическую точку для двух групп пациентов и выяснить, у кого она раньше.

Решение

По имеющимся оценкам параметров КМНС процесса $\hat{\tau}$ и $\hat{\eta}$ находим первые критические точки по формуле (А.Г. Барт, 2003)

$$T_{кр}^1 = \frac{2\pi - \phi}{\tau} + Q,$$

где $\operatorname{tg} \phi = \eta/\tau$, $e^{2\pi\eta/\tau} (\cos \tau Q + \frac{\eta}{\tau} \sin \tau Q) = e^{\eta Q}$.

Результаты: корреляционная структура

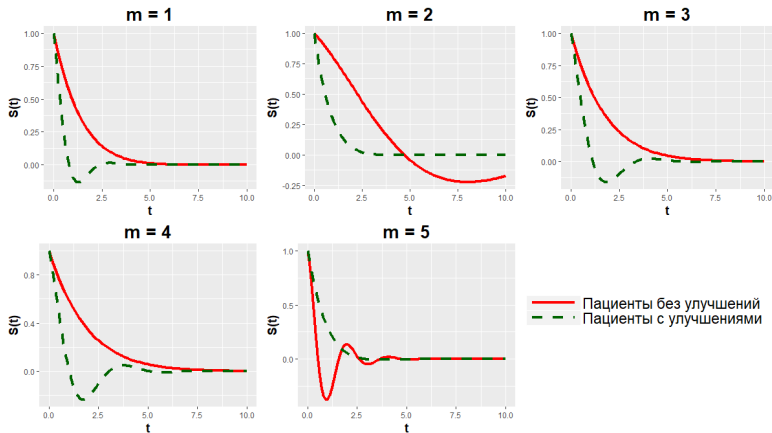


Рис.: $S(t) = e^{-\eta t} \cos \tau t$, m – номер заполнения.

Результаты по 50 заполнениям: в 33 случаях распад корреляционных связей для пациентов с улучшениями больше.

По критерию знаков p -значение = $0.016 < 0.05$.

Результаты: модели кривых саногенеза

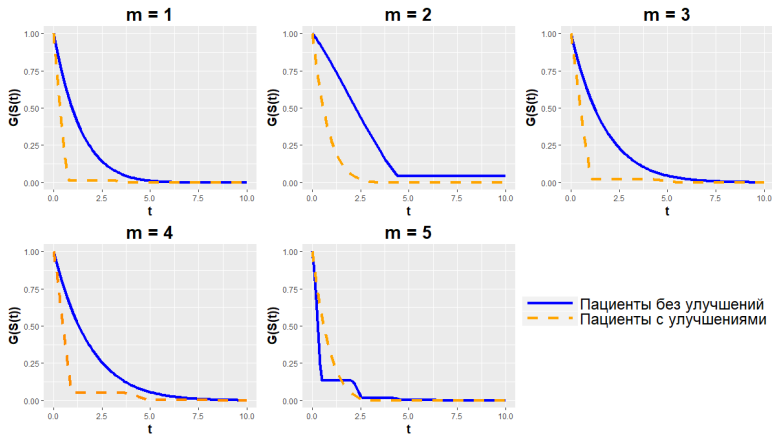


Рис.: Модель кривой саногенеза для обеих групп пациентов при разных заполнениях пропусков (m — номер заполнения).

Вывод: для пациентов без улучшений модель кривой саногенеза выше.

m	$T_{кр}^1$ для пациентов без улучшений	$T_{кр}^1$ для пациентов с улучшениями
1	6.57	0.77
2	4.41	3.40
3	10.96	1.08
4	11.18	0.94
5	0.45	2.84

Таблица: Первая критическая точка при разных заполнениях пропусков (**m** — номер заполнения).

Результаты по 50 заполнениям: в 38 случаях первая критическая точка для пациентов с улучшениями раньше. По критерию знаков p -значение = 0.0002 < 0.05.

- Реализован алгоритм заполнения пропусков в реальных медицинских данных на основе метода РММ.
- Осуществлен анализ корреляционной структуры данных на основе модели КМНС процесса.
- Доказана устойчивость оценок основных параметров КМНС процесса к знакопеременной структуре главных компонент и их перестановке.
- Показана статистическая значимость порядка критических точек для двух групп пациентов.
- Создана программа на языке программирования R, позволяющая производить вычисления от заполнения пропусков до оценки критических точек.