

Модели обнаружения зависимостей во временных рядах в задачах построения прогностических моделей

Карина Равилевна Усманова

Московский физико-технический институт
Физтех-школа прикладной математики и информатики
Факультет управления и прикладной математики
Кафедра «Интеллектуальные системы»

Научный руководитель: д.ф.-м.н. В.В.Стрижов

Выпускная квалификационная работа бакалавра

Москва 2018

Цель

Установить связи между заданными временными рядами с помощью метода сходящегося перекрестного отображения (ССМ).

Мотивация

Повысить качество прогноза временного ряда путем использования истории временных рядов, коррелирующих с ним.

- Sugihara G., May R., Ye H., Hsieh C., Deyle E., Fogarty M., and Munch S // Detecting causality in complex ecosystems. 2012.
- Sugihara G., May R. Nonlinear forecasting as a way of distinguishing chaos from measurement error in time series // Nature. 1990.
- Golyandina N. and Stepanov D. SSA-based approaches to analysis and forecast of multidimensional time series // 5th St. Petersburg workshop on simulation. 2005.

Задан временной ряд $\mathbf{x} = [x_1, \dots, x_N]$. По нему строится траекторная матрица \mathbf{H}_x :

$$\mathbf{H}_x = \begin{pmatrix} x_1 & \dots & x_{L-1} & x_L \\ x_2 & \dots & x_L & x_{L+1} \\ \vdots & \vdots & \vdots & \vdots \\ x_{N-L+1} & \dots & x_{N-1} & x_N \end{pmatrix} = \begin{pmatrix} \mathbf{x}_L \\ \mathbf{x}_{L+1} \\ \vdots \\ \mathbf{x}_N \end{pmatrix}$$

где L – ширина окна.

$$\mathbf{x}_i = (x_{i-L+1}, \dots, x_{i-1}, x_i) \in \mathbf{M}_x, \quad i = L, \dots, N,$$

где \mathbf{M}_x – траекторное пространство ряда \mathbf{x} .

Задача

Для заданных временных рядов $\mathbf{x} = [x_1, \dots, x_N]$ и $\mathbf{y} = [y_1, \dots, y_N]$ установить наличие или отсутствие связи между ними.

Решение

Считаем, что ряд \mathbf{y} зависит от ряда \mathbf{x} , если существует липшицево отображение $\varphi : \mathbf{M}_x \rightarrow \mathbf{M}_y$

$$\rho_{\mathbf{M}_y}(\varphi(\mathbf{x}_i), \varphi(\mathbf{x}_j)) \leq L \cdot \rho_{\mathbf{M}_x}(\mathbf{x}_i, \mathbf{x}_j), \quad \forall \mathbf{x}_i, \mathbf{x}_j \in \mathbf{M}_x$$

- Выбираем $\mathbf{x}_{t^*} = (x_{t^*-L+1}, \dots, x_{t^*-1}, x_{t^*})$.
- Пусть $\mathbf{x}_{t_1}, \dots, \mathbf{x}_{t_k}$ – k ближайших соседей вектора \mathbf{x}_{t^*} в пространстве \mathbf{M}_x . Тогда $\mathbf{y}_{t^*}, \mathbf{y}_{t_1}, \dots, \mathbf{y}_{t_k}$ – строки матрицы \mathbf{H}_y , соответствующие индексам t_1, \dots, t_k .

$$S(\mathbf{x}, \mathbf{y}) = \frac{\text{dist}(\mathbf{x})}{\text{dist}(\mathbf{y})}, \quad \text{dist}(\mathbf{x}) = \frac{1}{k} \sum_{i=1}^k \|\mathbf{x}_{t^*} - \mathbf{x}_{t_i}\|_2$$

Если $S(\mathbf{x}, \mathbf{y})$ меньше некоторого порога s , то ряд \mathbf{y} зависит от ряда \mathbf{x} .

Построение проекции

- Сингулярное разложение траекторной матрицы:

$$\mathbf{H}_x = \mathbf{U}_x \mathbf{\Lambda}_x \mathbf{V}_x$$

- Выберем \mathcal{T}_x – некоторый набор индексов компонент ряда \mathbf{x}
- $\mathbf{M}_{\mathcal{T}_x} \subset \mathbf{M}_x$ – траекторное подпространство
- Проекция ряда \mathbf{x} в подпространство $\mathbf{M}_{\mathcal{T}_x}$, описывается траекторной матрицей

$$\mathbf{P}_{\mathcal{T}_x} = \mathbf{U}_x \tilde{\mathbf{\Lambda}}_x \mathbf{V}_x$$

$$S(\mathbf{x}, \mathbf{y}, \mathcal{T}_x, \mathcal{T}_y) = \frac{\text{dist}(\mathbf{x}, \mathcal{T}_x)}{\text{dist}(\mathbf{y}, \mathcal{T}_y)}, \quad \text{dist}(\mathbf{x}, \mathcal{T}_x) = \frac{1}{k} \sum_{i=1}^k \|\mathbf{x}_{t^*} - \mathbf{x}_{t_i}\|_2$$

Задача поиска подпространств $\mathbf{M}_{\mathcal{T}_x}$ и $\mathbf{M}_{\mathcal{T}_y}$ эквивалентна поиску номеров главных компонент $(\mathcal{T}_x, \mathcal{T}_y)$

$$(\mathcal{T}_x, \mathcal{T}_y) = \arg \max_{\mathcal{T}_x, \mathcal{T}_y} S(\mathbf{x}, \mathbf{y}, \mathcal{T}_x, \mathcal{T}_y),$$

$$|\mathcal{T}_x| \rightarrow \min$$

$$|\mathcal{T}_y| \rightarrow \min$$

$$\mathbf{x} = \sin t + 2 \sin \frac{t}{2} + \sigma_x^2 \varepsilon, \quad \sigma_x^2 = 0.3, \quad \varepsilon \in \mathcal{N}(\mathbf{0}, \mathbf{I})$$

$$\mathbf{y} = \sin(2t + 5) + \sigma_y^2 \varepsilon, \quad \sigma_y^2 = 0.25, \quad \varepsilon \in \mathcal{N}(\mathbf{0}, \mathbf{I})$$

Эксперимент, часть 1

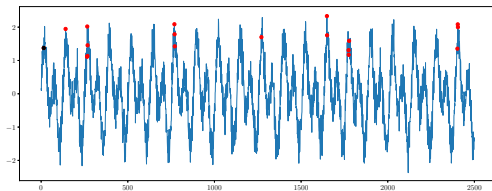


Рис.: Ближайшие соседи точки x_{15}

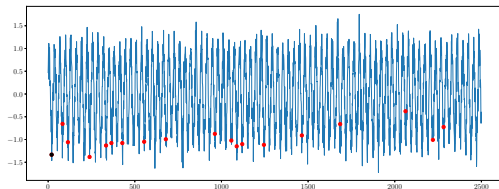


Рис.: Ближайшие соседи точки y_{15}

Ближайшие соседи на фазовых диаграммах

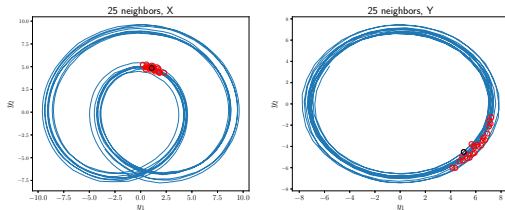


Рис.: Определение ближайших соседей по ряду x

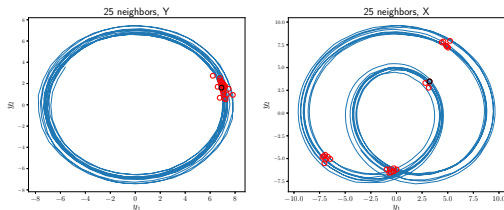


Рис.: Определение ближайших соседей по ряду y

Эксперимент, часть 2

Эксперимент проводился на данных потребления электроэнергии и температуры в течение года.

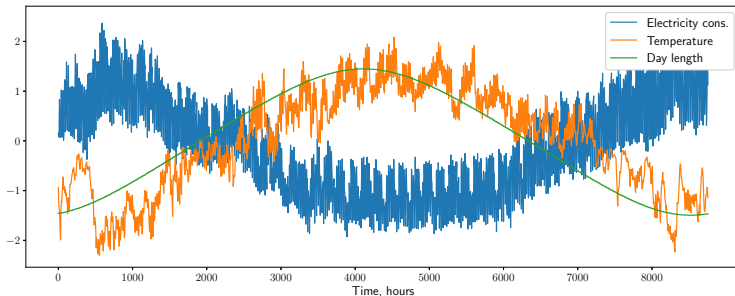


Рис.: Нормированные ряды потребления электроэнергии, температуры и длины светового дня

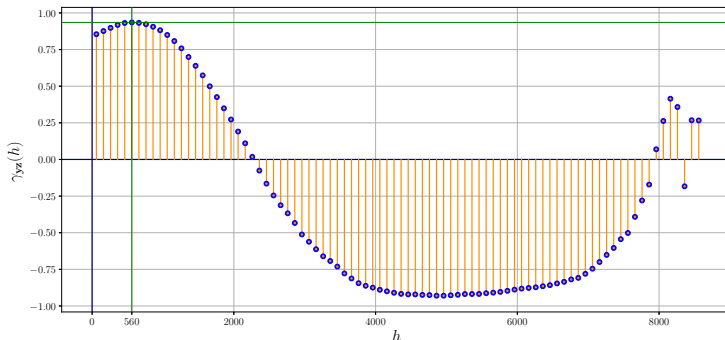


Рис.: Кросс-корреляционная диаграмма для ряда температуры и длины светового дня

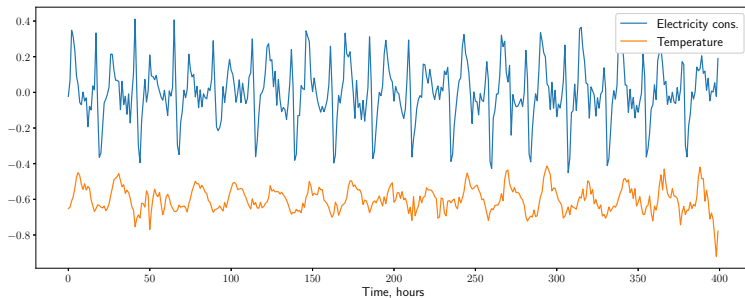


Рис.: Продифференцированные и нормированные ряды потребления электроэнергии и температуры

Ближайшие соседи на фазовых траекториях

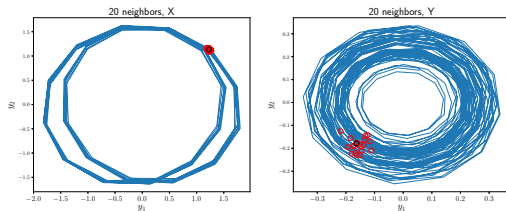


Рис.: Определение ближайших соседей по ряду x

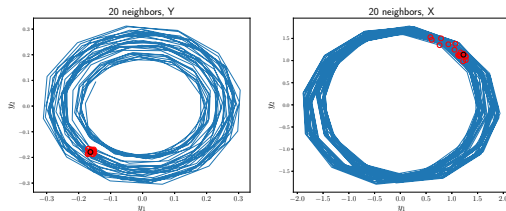


Рис.: Определение ближайших соседей по ряду y

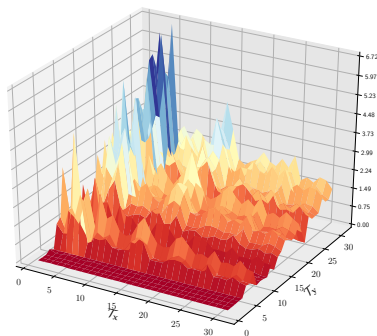


Рис.: Значения функционала $S(x, y, T_x, T_y)$ для различных наборов компонент (T_x, T_y)

- С помощью метода ССМ исследованы связи между рядами потребления электроэнергии и температуры.
- Исследована связь между проекциями этих рядов в различные подпространства.

К. Р. Усманова, С. П. Кудияров, Р. В. Мартышкин, А. А. Замковой, В. В. Стрижов // Анализ зависимостей между показателями при прогнозировании объема грузоперевозок // Системы и средства информатики, 2018.