

Погружение временных рядов с высокой волатильностью в метрическое пространство

Алтай Эйнуллаев Эльшан оглы

Московский физико-технический институт

Курс: Автоматизация научных исследований
(практика, В. В. Стрижов)/Группа 105

Эксперт: д.ф-м.н. В.В.Стрижов

Консультант: К.Яковлев

2024

Цель исследования

Решается задача прогнозирования набора сильно коррелированных временных рядов с высокой волатильностью. Предлагается прогнозировать с помощью матрицы попарных расстояний между временными рядами набора. Исследуются влияние различных способов вычисления попарного расстояния на точность прогноза.

Постановка задачи

Пусть $X = \{\mathbf{x} = [x_1, \dots, x_m]^T | x_i \in R\}$ — множество временных рядов, заданных своей реализацией. Обозначим через $\mathbf{Y} \in R^{n \times m}$ заданный набор из n временных рядов:

$$\mathbf{Y} = [\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}]^T. \quad (1)$$

Через $\mathbf{Y}_t \in R^{n \times t}$ обозначим $t < m$ первых столбцов \mathbf{Y} :

$$\mathbf{Y}_t = [\mathbf{x}_{1:t}^{(1)}, \dots, \mathbf{x}_{1:t}^{(n)}]^T. \quad (2)$$

Таким образом, по известной \mathbf{Y}_t выполняется прогноз значений набора временных рядов в момент времени $t + 1$:

$$\mathbf{Y}_t \rightarrow \mathbf{x}_{t+1} = [x_{t+1}^{(1)}, \dots, x_{t+1}^{(n)}]^T. \quad (3)$$

Определим функцию расстояния между временными рядами:
 $d : X \times X \rightarrow R$, удовлетворяющую условиям Мерсера.

$$d_t(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) = d_t(i, j) \quad (4)$$

Таким образом, в каждый момент времени t набору временных рядов \mathbf{Y}_t поставлена в соответствие матрица попарных расстояний $\Sigma_t \in \mathcal{S}_n^+$ (симметричная, неотрицательно определенная матрица). С помощью Score-Based Generative Models и LSTM прогнозируем матрицу попарных расстояний $\hat{\Sigma}_{t+1}$. Далее возвращаем прогноз в исходное пространство:

$$\mathbf{x}_{t+1} = \arg \min_{\mathbf{x} \in R^n} \|\Sigma_{t+1} - \hat{\Sigma}_{t+1}\|. \quad (5)$$

Выбор функции попарных расстояний

Набор временных рядов в момент времени t :

$$\mathbf{X}_t = \begin{pmatrix} x_{t-L+1}^{(1)} & x_{t-L+2}^{(1)} & \cdots & x_t^{(1)} \\ x_{t-L+1}^{(2)} & x_{t-L+2}^{(2)} & \cdots & x_t^{(2)} \\ \cdots & \cdots & \cdots & \cdots \\ x_{t-L+1}^{(n)} & x_{t-L+2}^{(n)} & \cdots & x_t^{(n)} \end{pmatrix} \quad (6)$$

i -й столбец матрицы \mathbf{X}_t — \mathbf{y}^i , i -я строка — $\mathbf{x}^{(i)}$.

$$\Sigma_t^1 = \frac{1}{L} \sum_{i=1}^L \mathbf{y}^i \mathbf{y}^{iT} \quad (7)$$

Введем вектор $\mathbf{m} \in R^n$ следующим образом:

$$\mathbf{m}_j = \frac{1}{L} \sum_{i=1}^L y_j^i. \quad (8)$$

Выбор функции попарных расстояний

Второй способ:

$$\Sigma_t^2 = \frac{1}{L} \sum_{i=1}^L (\mathbf{y}^i - \mathbf{m})(\mathbf{y}^i - \mathbf{m})^T \quad (9)$$

Для каждой строки $\mathbf{x}^{(j)}$, $j \in 1, \dots, n$ матрицы \mathbf{X} введем величину:

$$\sigma_j = \sqrt{\frac{1}{L} \sum_{i=1}^L (\mathbf{x}_i^{(j)} - \mathbf{m}_i)^2}. \quad (10)$$

Теперь введем матрицу Σ_t^3 следующим образом:

$$(\Sigma_t^3)_{ij} = \frac{(\Sigma_t^2)_{ij}}{\sigma_i \sigma_j}. \quad (11)$$

Вычислительный эксперимент

Эксперимент состоял в прогнозировании матрицы попарных расстояний для каждого способа ее вычисления.

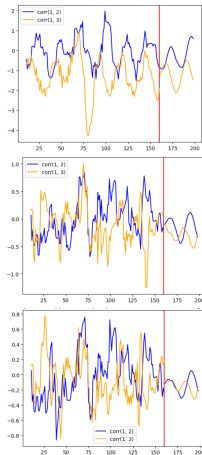


Рис.: Прогнозирование Σ_t^i на синтетическом наборе синусов

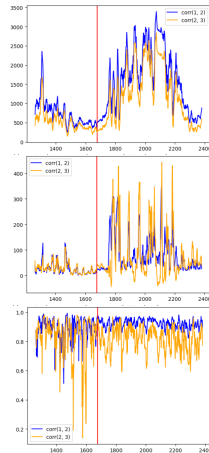


Рис.: Прогнозирование Σ_t^i на наборе цен на электричество

Анализ ошибки

Вычислим ошибку MAE каждого из прогнозов. В случае синтетического набора: $MAE_1 = 0.55$, $MAE_2 = 0.17$, $MAE_3 = 0.17$. На реальном наборе: $MAE_1 = 895.2$, $MAE_2 = 152.8$, $MAE_3 = 0.11$.

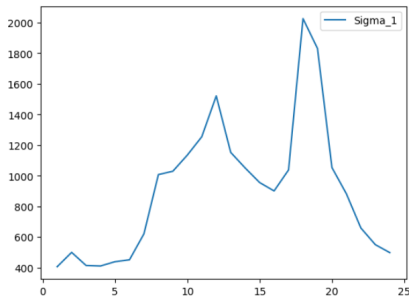


Рис.: Средняя MAE прогноза попарных расстояний для каждого из рядов при 1-ом способе подсчета Σ

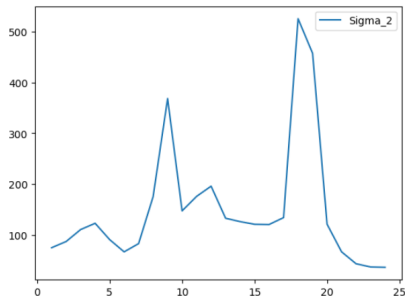


Рис.: Средняя MAE прогноза попарных расстояний для каждого из рядов при 2-ом способе подсчета Σ

Наибольшая ошибка на временных рядах под номером 18, 19 связана с высокой дисперсией соответствующих рядов.

Отказываясь от их прогнозирования, получаем существенное улучшение прогноза:

$$MAE_1 = 785.7 \quad (12)$$

$$MAE_2 = 94.1 \quad (13)$$

Перечислите ваши результаты

- ▶ Прогнозирования набора высоко коррелированных временных рядов с высокой волатильностью.
- ▶ Проведен анализ различных способов вычисления матрицы попарных расстояний и их прогноза.
- ▶ Осталось провести эксперимент с Score-Based Gen Models и оценить точность возвращения прогноза в исходное пространство.