

Генеративный причинно-следственный подход к анализу данных ИМК

Владимиров Э.А.

Московский физико-технический институт

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

2023

Причинно-следственный анализ

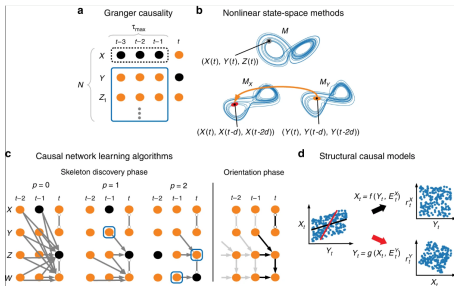
Проблема

- ▶ Традиционные методы (корреляция, линейная регрессия) неадекватны для сложных нелинейных связей
- ▶ Данные имеют высокую размерность, что усложняет поиск причинно-следственных связей
- ▶ Зависимости между переменными могут изменяться во времени

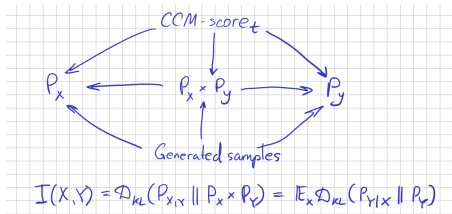
Решение

Построить устойчивую и интерпретируемую форму вероятностного анализа причинного влияния $\mathbf{X} \rightarrow \mathbf{Y}$

Различные подходы к поиску связей



Существующие подходы



Предлагаемый подход

Постановка задачи

Пусть $\mathbf{X}(t) = \{X_1(t), X_2(t), \dots, X_{n_x}(t)\}$ и $\mathbf{Y}(t) = \{Y_1(t), Y_2(t), \dots, Y_{n_y}(t)\}$ — два набора многомерных временных рядов, наблюдаемых в моменты времени $t = 1, \dots, T$.

Необходимо определить направленные причинные связи:

1. $X_i(t - \tau) \rightarrow Y_j(t)$ для $i = 1, \dots, n_x, j = 1, \dots, n_y$, и лагов $\tau \geq 0$,
2. $Y_j(t - \tau) \rightarrow X_i(t)$ для $i = 1, \dots, n_x, j = 1, \dots, n_y$, и лагов $\tau \geq 0$.

Предполагаем, что многомерные временные ряды $\mathbf{X}(t)$ и $\mathbf{Y}(t)$ генерируются следующим образом:

$$X_i(t) = f_i(\text{Pa}_{X_i}(t), \varepsilon_{X_i}(t)),$$

$$Y_j(t) = g_j(\text{Pa}_{Y_j}(t), \varepsilon_{Y_j}(t)),$$

где:

- ▶ $\text{Pa}_{X_i}(t) \subseteq \{Y_1(t - \tau), \dots, Y_{n_y}(t - \tau)\}$ — множество родителей переменной $X_i(t)$ из \mathbf{Y} ,
- ▶ $\text{Pa}_{Y_j}(t) \subseteq \{X_1(t - \tau), \dots, X_{n_x}(t - \tau)\}$ — множество родителей переменной $Y_j(t)$ из \mathbf{X} ,
- ▶ f_i и g_j — детерминированные функции, описывающие зависимость,
- ▶ $\varepsilon_{X_i}(t)$ и $\varepsilon_{Y_j}(t)$ — шумовые компоненты.

Оптимизационная задача:

$$\min_{G_{XY}, G_{YX}} \mathcal{L}(\mathbf{X}, \mathbf{Y} \mid G_{XY}, G_{YX}) + \lambda_1 \mathcal{R}(G_{XY}, G_{YX}) + \lambda_2 \mathcal{T}(G_{XY}, G_{YX}),$$

где:

- ▶ G_{XY} — граф зависимостей $X_i \rightarrow Y_j$,
- ▶ G_{YX} — граф зависимостей $Y_j \rightarrow X_i$,
- ▶ $\mathcal{L}(\mathbf{X}, \mathbf{Y} \mid G_{XY}, G_{YX})$ — правдоподобие наблюдаемых данных с учетом графов G_{XY} и G_{YX} ,
- ▶ $\mathcal{R}(G_{XY}, G_{YX})$ — регуляризатор, штрафующий за сложность графов,
- ▶ $\mathcal{T}(G_{XY}, G_{YX})$ — штраф за избыточную изменчивость графов во времени.

Independent Component Analysis

Предположим, что $\mathbf{X}(t)$ образуется из нескольких скрытых источников $\mathbf{S}(t) \in \mathbb{R}^{d_s}$:

$$\mathbf{X}(t) = A\mathbf{S}(t), \quad A \in \mathbb{R}^{d_x \times d_s}.$$

Каждая компонента $S_k(t)$ предполагается статистически независимой от других:

$$p(\mathbf{S}) = \prod_{k=1}^{d_s} p(S_k).$$

Задача оптимизации: Найти обратную матрицу \hat{A}^{-1} , дающую

$$\hat{\mathbf{S}}(t) = \hat{A}^{-1} \mathbf{X}(t),$$

чтобы минимизировать взаимную информацию между компонентами $\hat{S}_k(t)$.

$$\text{MI}(\hat{\mathbf{S}}(t)) \approx \sum_{k=1}^{d_s} H(S_k) - H\left(\sum_k S_k\right),$$

Convergent Cross Mapping

Теневое вложение (delay embedding):

$$M_{X,t} = (X_t, X_{t-\tau}, \dots, X_{t-(E-1)\tau}) \in \mathbb{R}^E,$$

где E — размерность вложения, τ — временной лаг. Аналогично задаётся $M_{Y,t} = (Y_t, Y_{t-\tau}, \dots)$.

Реконструкция:

$$\hat{Y}_t = \sum_{i=1}^k w_i Y_{n_i},$$

здесь n_i — индексы ближайших соседей точки $M_{X,t}$ в пространстве M_X , а w_i — веса, зависящие от расстояния до $M_{X,t}$.

Критерий причинности:

$$\rho_{X \rightarrow Y} = \text{corr}(\{\hat{Y}_t\}, \{Y_t\}).$$

Если при увеличении размера “библиотеки” (количества доступных точек) значение $\rho_{X \rightarrow Y}$ *возрастает*, считается, что $\mathbf{X}(t)$ действительно влияет на $\mathbf{Y}(t)$.

Probabilistic CCM

Идея: Вместо единственного прогноза \hat{Y}_t рассматривается *полное условное распределение*

$$p_L(Y_t | M_{X,t}),$$

оценённое по выборке размера L . Ближайшие соседи в пространстве M_X позволяют построить *вероятностную аппроксимацию* (например, ядерным методом):

$$p_L(y | M_{X,t}) = \frac{1}{Z_t} \sum_{i \in N_L(t)} K(y - Y_{n_i}),$$

где $K(\cdot)$ — ядро (например, гауссово), $N_L(t)$ — множество соседей точки $M_{X,t}$, а Z_t — нормировочная константа.

Оценка причинности как MI:

$$I_L(X \rightarrow Y) = \mathbb{E}_{M_{X,t}} D_{\text{KL}}(p_L(Y_t | \hat{M}_{X,t}) \parallel p(Y_t)),$$

Сходящееся свойство: При $L \rightarrow \infty$ (при достаточно плотном покрытии пространства)

$$p_L(Y_t | M_{X,t}) \rightarrow p(Y_t | M_{X,t}),$$

Предлагаемый метод

1. Независимый анализ компонент (ICA).

Для исходных ЭЭГ-данных $\mathbf{X}_{\text{raw}}(t) \in \mathbb{R}^{d_x}$ получаем независимые компоненты:

$$\hat{\mathbf{S}}(t) = \hat{A}^{-1} \mathbf{X}_{\text{raw}}(t).$$

2. Построение эмбедингов.

Для каждого времени t формируем вектор:

$$M_{X,t} = (\hat{\mathbf{S}}(t), \hat{\mathbf{S}}(t - \tau), \dots, \hat{\mathbf{S}}(t - (E - 1)\tau)),$$

3. Оценка причинно-следственных связей.

В полученном пространстве $(M_{X,t}, M_{Y,t})$ определяем меру влияния $\mathbf{X} \rightarrow \mathbf{Y}$, вычисляя:

$$\gamma(t) = \text{Prob} - \text{CCM}(M_{X,t}, M_{Y,t}).$$

Результат — временной ряд $\{\gamma(t)\}$, отражающий динамику влияния $X \rightarrow Y$.

Риманова постановка задачи

Причинно-следственная связь — вероятность push-forward-a или диффеоморфизма.

Основная проблема — построение многообразий

Возможные варианты:

- ▶ PINN
- ▶ Riemannian space of covariance matrices

Вычислительный эксперимент на данных ЭЭГ - ИИМ

Данные

У 25 участников были записаны показания ЭЭГ, ИИМ, МРТ во время игры в настольный теннис. С каждым участником было сыграно 4 сессии, длительность каждой из них составляет 7-10 минут.

Human Player



Ball Machine



Block 1		Block 2		Block 3		Block 4	
Machine Rally	Cooperative	Machine Serve	Competitive	Cooperative	Machine Serve	Competitive	Machine Rally
2:30 2:30 2:30	7:30						

15 min.