

Учебник по временным рядам: начало

Винни-Пух

28 января 2023 г.

Одномерные временные ряды

Стационарные процессы

Из курса математического анализа мы знаем разницу между рядами и последовательностями. В последовательности числа записаны одно за другим, скажем, через запятую,

$$5, 8, -3, 2, 4, 5, \dots$$

А ряд — это бесконечная сумма чисел, например,

$$0.9 + 0.09 + 0.009 + 0.0009 + \dots$$

Настала пора дать первое определение и раскрыть заговор рептилоидов!

Определение 1. Временной ряд — это последовательность случайных величин.

Индекс временного ряда может быть любым, но чаще всего мы работаем с тремя случаями. Бесконечный в обе стороны индекс,

$$\dots, y_{-2}, y_{-1}, y_0, y_1, y_2, \dots$$

бесконечный в одну сторону,

$$y_1, y_2, y_3, y_4, \dots$$

либо конечный,

$$y_1, y_2, y_3, y_4, \dots, y_T.$$

Чтобы отличать весь временной ряд от одной конкретной случайной величины, мы будем использовать обозначения:

y_t — одна конкретная случайная величина;

$(y_t) = y_1, y_2, y_3, y_4, \dots$ — вся последовательность случайных величин.

Если контекст требует, то можно проявить больше аккуратности и указать возможные значения индекса, например, $(y_t)_{t=1}^{\infty}$.

Начнём с самого простого временного ряда — белого шума.

Определение 2. Ряд (u_t) называется белым шумом (white noise), если он удовлетворяет трём свойствам:

Да, да, всё верно, временной ряд — это не ряд, ноль — чётное число, единица — не простое, бульёные кубики — не кубики, московские диаметры — не диаметры, а Деда Мороза не существует.

или занудства

- a) Нулевое математическое ожидание, $E(u_t) = 0$ для любого t .
- б) Постоянная дисперсия, $\text{Var}(u_t) = \sigma_u^2$ для любого t .
- в) Нулевая ковариация, $\text{Var}(u_t, u_s) = 0$ для любых $t \neq s$.

Заметим, что случайные величины в белом шуме вполне могут быть зависимы. Например,

...

Определение 3. Ряд (y_t) называется слабо стационарным (weakly stationary), или просто стационарным, если он удовлетворяет трём свойствам:

- a) Постоянное математическое ожидание, $E(y_t) = \mu$ для любого t .
- б) Постоянная дисперсия, $\text{Var}(y_t) = \gamma_0$ для любого t .
- в) Ковариация двух величин зависит только от их удалённости по времени друг от друга, $\text{Var}(y_t, y_s) = \gamma_{t-s}$ для любых t и s .

Из третьего условия на ковариацию $\text{Var}(y_t, y_s) = \gamma_{t-s}$ следует постоянство дисперсии, достаточно подставить $t = s$ и увидеть, что $\text{Var}(y_t, y_t) = \text{Var}(y_t) = \gamma_0$. Мы выписали второе свойство отдельно от третьего, чтобы лучше его выделить.

Определение 4. Ряд (y_t) называется процессом скользящего среднего порядка q (moving average of order q), если он представим в виде:

$$y_t = u_t + b_1 u_{t-1} + \dots + b_q u_{t-q},$$

где (u_t) — белый шум. Обозначаем такие процессы мы так: $y_t \sim MA(q)$.

Процесс скользящего среднего — это статистическая модель. Название скользящего среднего имеет одну из простых процедур сглаживания ряда.

Сглаживание ряда

При сглаживании ряда мы из исходного ряда (y_t) получаем новый ряд (\tilde{y}_t) с меньшей изменчивостью. Количество наблюдений при этом может как немного поменяться, так и сохраниться, в зависимости от конкретного алгоритма.

Определение 5. Взятие скользящего среднего с шириной окна $h = 3$ — алгоритм сглаживания с формулой

$$\tilde{y}_t = \frac{y_{t-1} + y_t + y_{t+1}}{3}.$$

Для краткости можно использовать обозначение $\tilde{y} = MA(y)$.

Как выглядит скользящее среднее с другой нечётной шириной окна читатель может попробовать догадаться сам, к примеру, выписав формулу для скользящего среднего с шириной окна $h = 5$.

LOESS

Обычная регрессия (ordinary least squares) проводит одну линию $\hat{y}_t = \hat{\beta}_1 + \hat{\beta}_2 x_t$. В роли x_t может быть просто само время t .

Вспомним целевую функцию обычной регрессии

$$Q(\hat{\beta}_1, \hat{\beta}_2) = \sum_{t=1}^T (y_t - \hat{y}_t)^2 = \sum_{t=1}^T (y_t - (\hat{\beta}_1 + \hat{\beta}_2 x_t))^2.$$

В результате минимизации

$$Q(\hat{\beta}_1, \hat{\beta}_2) \rightarrow Q(\hat{\beta}_1, \hat{\beta}_2)_{\hat{\beta}_1, \hat{\beta}_2}$$

получается единственное значение $\hat{\beta}_1$ и $\hat{\beta}_2$.

LOESS = LOcal regrESSion = ЛОкальная регрессия

LOESS проводит свою линию регрессии для каждого x . Целевая функция теперь зависит от абсциссы точки x , в которой мы строим регрессию,

$$Q(\hat{\beta}_1, \hat{\beta}_2) = \sum_{t=1}^T K(x_t, x)(y_t - \hat{y}_t)^2 = \sum_{t=1}^T K(x_t, x)(y_t - (\hat{\beta}_1 + \hat{\beta}_2 x_t))^2.$$

Функция весов $K(x_t, x)$ должна давать большой положительный вес точкам x_t рядом с точкой x и маленький положительный, или даже нулевой, вес точкам x_t далеко от точки x .

Например, в качестве функции весов $K(x_t, x)$ можно использовать

$$K(x_t, x) = \exp\left(-\frac{(x_t - x)^2}{h^2}\right).$$

При $h \rightarrow \infty$ мы получим обычные оценки метода наименьших квадратов $\hat{\beta}_1(x) = \hat{\beta}_1^{\text{OLS}}$, $\hat{\beta}_2(x) = \hat{\beta}_2^{\text{OLS}}$.

Функция весов может быть и такой

$$K(x_t, x) = \begin{cases} 1, & \text{если } x_t \text{ — это один из пяти ближайших соседей } x, \\ 0, & \text{иначе} \end{cases}.$$

Оптимизируем мы по прежнему по двум переменным,

$$Q(\hat{\beta}_1, \hat{\beta}_2, x) \rightarrow Q(\hat{\beta}_1, \hat{\beta}_2)_{\hat{\beta}_1, \hat{\beta}_2}$$

Задача оптимизации — выпуклая, есть решение в явном виде. Только теперь получаются оптимальные коэффициенты, зависящие от точки x , $\hat{\beta}_1(x)$ и $\hat{\beta}_2(x)$.

Для получения сглаженного значения \tilde{y}_t мы берём $x_t = t$ и

$$\tilde{y}_t = \hat{\beta}_1(t) + \hat{\beta}_2(t)t.$$

STL

Изложим упрощённый вариант STL-алгоритма для месячных данных без выбросов.

a) Положим $\text{trend}_t = 0$ и $\text{season}_t = 0$.

б) Детрендируем исходный ряд,

$$D_t = y_t - \text{trend}_t.$$

в) Разрежем детрендированный ряд (D_t) на двенадцать подрядов по месяцам.

$$(D_t) \rightarrow (D_t^{jan}), (D_t^{feb}), \dots, (D_t^{dec}).$$

г) Сгладим каждый подряд с помощью LOESS:

$$C^{jan} = LOESS(D^{jan}), \dots, C^{dec} = LOESS(D^{dec})$$

д) Соберём двенадцать подрядов в один ряд

$$(C_t^{jan}), (C_t^{feb}), \dots, (C_t^{dec}) \rightarrow (C_t)$$

е) Сильно сгладим собранный ряд

$$L = LOESS(MA(MA(C))).$$

ж) Обновим сезонную составляющую

$$\text{season}_t = C_t - L_t.$$

з) Обновим тренд

$$\text{trend}_t = y_t - \text{season}_t.$$

Далее перейдём к шагу 2 и пройдём шаги 2-8 ещё раз.

DLT

Вспомним ETS(A, Ad, A) модель в структурной форме.

$$\begin{cases} y_t = (\ell_{t-1} + \phi b_{t-1}) + s_{t-12} + u_t \\ s_t = s_{t-12} + \gamma u_t \\ \ell_t = \ell_{t-1} + \phi b_{t-1} + \alpha u_t \\ b_t = \phi b_{t-1} + \beta u_t \\ u_t \sim \text{Normal}(\text{loc} = 0, \text{scale} = \sigma). \end{cases}$$

Стандартной ссылкой по ETS моделям является [hyndman2018forecasting](#).

Чтобы перейти к DLT модели сделаем ряд обобщений:

- Добавим в наблюдаемый процесс y_t регрессионную составляющую:

$$r_t = \beta_1 x_{t1} + \beta_2 x_{t2}.$$

- Добавим в наблюдаемый процесс глобальный тренд, например, линейный:

$$g_t = \delta_1 + \delta_2 t.$$

- Перейдем от нормального распределения ошибки к распределению Стьюдента,

$$u_t \sim \text{Student}(\text{df} = \nu, \text{loc} = 0, \text{scale} = \sigma).$$

Кроме того, в описании DLT модели по сравнению с ETS моделью почему-то сдвинут индекс у сезонной составляющей. То есть величина, называемая s_t у Хиндмана в **hyndman2018forecasting**, в статье **ng2020orbit** названа s_{t+12} . Поэтому уравнение на сезонную составляющую принимает вид

$$s_{t+12} = s_t + \gamma u_t.$$

С учётом новых составляющих и сдвига индекса у сезонности уравнение на наблюдаемый y_t примет вид

$$y_t = g_t + (\ell_{t-1} + \phi b_{t-1}) + s_t + r_t + u_t.$$

DLT модель

Наблюдаемый ряд y_t раскладывается в сумму составляющих: глобальный тренд, локальное отклонение от глобального тренда, сезонная составляющая, регрессионная составляющая, ошибка.

$$y_t = g_t + (\ell_{t-1} + \phi b_{t-1}) + s_t + r_t + u_t$$

Глобальный тренд g_t может быть задан по-разному. Например, линейно

$$g_t = \gamma_1 + \gamma_2 t.$$

Сезонная составляющая плавно меняется во времени,

$$s_{t+12} = s_t + \gamma u_t.$$

Скорость локального тренда плавно меняется,

$$b_t = \phi b_{t-1} + \beta u_t.$$

Локальный тренд

$$\ell_t = \ell_{t-1} + \phi b_{t-1} + \alpha u_t.$$

Регрессионная составляющая на примере двух регрессоров,

$$r_t = \beta_1 x_{t1} + \beta_2 x_{t2}.$$

Ошибка,

$$u_t \sim \text{Student}(df = \nu, \text{loc} = 0, \text{scale} = \sigma).$$

Одной системой,

$$\begin{cases} y_t = g_t + (\ell_{t-1} + \phi b_{t-1}) + s_t + r_t + u_t \\ g_t = \gamma_1 + \gamma_2 t \\ s_{t+12} = s_t + \gamma u_t \\ r_t = \beta_1 x_{t1} + \beta_2 x_{t2} \\ \ell_t = \ell_{t-1} + \phi b_{t-1} + \alpha u_t \\ b_t = \phi b_{t-1} + \beta u_t \\ u_t \sim \text{Student}(df = \nu, \text{loc} = 0, \text{scale} = \sigma). \end{cases}$$

Используя формулу для ℓ_t можно записать y_t также в виде

$$y_t = g_t + \ell_t + s_t + r_t + (1 - \alpha)u_t.$$

Рекуррентные соотношения

Можно элегантно отказаться от u_t в уравнениях на s_{t+12} , ℓ_t и b_t .

Это полезно для описания модели на вероятностных языках программирования, будь то stan, numpguro или что-то ещё.

Выразим ошибку u_t из формулы для локального тренда и подставим в формулу для скорости роста локального тренда :

$$b_t = \phi b_{t-1} + \frac{\beta}{\alpha}(\ell_t - \ell_{t-1} - \phi b_{t-1}).$$

Перегруппируем и увидим, что скорость роста локального тренда b_t является средневзвешенным,

$$b_t = \frac{\beta}{\alpha}(\ell_t - \ell_{t-1}) + \left(1 - \frac{\beta}{\alpha}\right)\phi b_{t-1}.$$

Можно определить $\rho_b = \frac{\beta}{\alpha}$, и тогда

$$b_t = \rho_b(\ell_t - \ell_{t-1}) + (1 - \rho_b)\phi b_{t-1}.$$

В коде пакета orbit на stan соответствующая строка имеет вид

```
b[ t ] = slp_sm * (1[ t ] - 1[ t - 1 ]) + (1 - slp_sm) * DAMPED_FACTOR * b[ t - 1];
```

Теперь выразим ошибку u_t из и подставим в .

$$\ell_t = \ell_{t-1} + \phi b_{t-1} + \alpha(y_t - g_t - \ell_{t-1} - \phi b_{t-1} - s_t - r_t).$$

Перегруппируем и снова получаем вид средневзвешенного.

$$\ell_t = \alpha(y_t - g_t - s_t - r_t) + (1 - \alpha)(\ell_{t-1} + \phi b_{t-1}).$$

Смотрим на исходный код модели в stan,

```
lt_sum[t] = l[t-1] + DAMPED_FACTOR * b[t-1];
l[t] = lev_sm * (RESPONSE[t] - gt_sum[t] - s_t - r[t]) + (1 - lev_sm) * lt_sum[t];
```

На этот раз выразим ошибку u_t из и подставим в .

$$s_{t+12} = s_t + \frac{\gamma}{1-\alpha}(y_t - g_t - \ell_t - s_t - r_t).$$

Перегруппируем и получаем вид средневзвешенного,

$$s_{t+12} = \frac{\gamma}{1-\alpha}(y_t - g_t - \ell_t - r_t) + \left(1 - \frac{\gamma}{1-\alpha}\right)s_t.$$

При обозначении $\rho_s = \frac{\gamma}{1-\alpha}$ получаем

$$s_{t+12} = \rho_s(y_t - g_t - \ell_t - r_t) + (1 - \rho_s)s_t.$$

Соответствующий фрагмент кода в stan,

```
s[t + SEASONALITY] = sea_sm * (RESPONSE[t] - gt_sum[t] - l[t]
- r[t]) + (1 - sea_sm) * s_t;
```

Замечаем, что в статье **ng2020orbit** в описании DLT модели есть пара описок. Пропущено ϕ перед b_{t-1} в рекуррентной формуле для ℓ_t . Пропущено g_t в рекуррентной формуле для s_{t+m} .

Начальные условия

$$b_1 = 0$$

$$g_1 =$$

$$r_1 =$$

$$s_{12} = -(s_1 + s_2 + \dots + s_{11})$$

$$\ell_1 = y_1 - g_1 - s_1 - r_1$$

Априорные распределения

При $t \in \{1, 2, \dots, 11\}$,

$$s_t \sim \text{Normal}(\text{loc} = 0, \text{scale} = \sigma_s)$$

```
for ( i in 1:(SEASONALITY - 1))
  init_sea[ i ] ~ normal(0, SEASONALITY_SD);
```

$$\beta_j \sim \text{Normal}(\text{loc} = \mu_j, \text{scale} = \sigma_j),$$

где μ_j, σ_j — гиперпараметры, по умолчанию равные $\mu_j = 0$ и $\sigma_j = 1$.

$$\sigma \sim \text{HalfCauchy}(\text{loc} = 0, \text{scale} = \gamma_0),$$

где $\gamma_0 \dots$

Тесты на прогнозную силу

Тест Диболда-Мариано

Предпосылки:

Два прогноза, \hat{y}_t^A и \hat{y}_t^B . Разница произвольных метрик качества,

$$d_t = (\hat{y}_t^A - y_t)^2 - (\hat{y}_t^B - y_t)^2.$$

Процесс (d_t) стационарный. Другими словами $E(d_t) = \mu$, $\text{Var}(d_t, d_{t-k}) = \gamma_k$, в частности, $\text{Var}(d_t) = \gamma_0$.

Гипотезы:

$$H_0: E(d_t) = 0;$$

$$H_a: E(d_t) \neq 0;$$

Тестовая статистика при верной H_0 :

$$DM = \frac{\bar{d} - 0}{se(\bar{d})} \rightarrow \mathcal{N}(0; 1)$$

Трудность возникает только в оценке $se(\bar{d})$, так как значения d_t коррелированы.

Как правило оценивают регрессию вектора d_t на константу и используют робастную стандартную ошибку se_{HAC} .

$$\hat{d}_t = \hat{\beta}_1, \quad DM = \frac{\hat{\beta}_1 - 0}{se_{HAC}(\hat{\beta}_1)}.$$

В качестве альтернативного подхода можно дополнительно предположить, что (d_t) описывается стационарным $ARMA(p, q)$ процессом с небольшими p и q и рассчитать $se(\bar{d})$ в рамках этого предположения.

RC и SPA тесты

RC (Reality Check) тест Уайта и SPA (Superior Predictive Ability) тест Хансена обобщают тест Диболда-Мариано на случай сравнения множества прогнозов против одного эталонного.

Для обоих тестов используется стационарный бутстрэп. Опишем стационарный бутстрэп на примере более простой задачи.

Кратко про стационарный бутстрэп

Представим себе, что у нас есть ряд y_1, \dots, y_T , и мы хотим построить доверительный интервал для $\rho = \text{Corr}(y_t, y_{t-1})$ с помощью бутстрэпа.

Если использовать обычный бутстрэп, который из исходной выборки (y_t) много раз делает случайную выборку с повторениями, то структура временного ряда будет разрушаться при создании бутстрэп-выборок, и оценка корреляции по бутстрэп-выборкам будет каждый раз примерно нулевой.

Алгоритм стационарного бутстрэпа пытается решить эту проблему. На входе у нас временной ряд y_1, \dots, y_T . На выходе мы хотим получить бутстрэп копию этого ряда той же длины y_1^*, \dots, y_T^* .

1. Выберем параметр вероятности p . О правилах выбора чуть позже.
2. Выберем случайный момент времени $s \in \{1, \dots, T\}$ и запишем y_s очередным элементом в бутстрэп копию.
3. С вероятностью p вернемся к шагу 2, с вероятностью $1 - p$ пойдём дальше.
4. Увеличим s на 1, запишем y_s очередным элементом в бутстрэп копию и перейдем к подкидыванию монетки на шаге 3.

Алгоритм продолжается до тех пор, пока не наберем T наблюдений в бутстрэп-копию ряда.

Теперь мы можем построить бутстрэп-доверительный интервал для корреляции. Например, с помощью перцентильного бутстрэпа.

По исходному ряду создаем 10000 бутстрэп-копий ряда. По каждой бутстрэп-копии считаем оценку корреляции. Удаляем по 2.5% самых больших и самых маленьких оценок корреляции. Полученные края и будут границами доверительного интервала.

Про выбор p .

...

Возвращаемся к тестам.

Решения

Источники мудрости

Источники мудрости, кои автор подборки постарался не замутить.
Смело направляйте к ним верблюдов своего любопытства!

Список литературы

- (<https://math.stackexchange.com/users/169207/stumped>), stumped (б.г.).
- Expected value of sock pairs.* Mathematics Stack Exchange. URL:<https://math.stackexchange.com/q/894509> (version: 2014-08-11). eprint: <https://math.stackexchange.com/q/894509>. URL: <https://math.stackexchange.com/q/894509>.
- Blom, Gunnar (1994). *Problems and Snapshots from the World of Probability*. Springer. Задачи с решениями и небольшими исследованиями.
- Bruss, F Thomas (2000). «Sum the odds to one and stop». В: *Annals of Probability*, с. 1384—1391. URL: <https://projecteuclid.org/euclid.aop/1019160340>. Простейшая стратегия для игр, где надо остановиться на последнем успехе!
- Friedlen, DM и Doug Pryor (1990). «E3265». В: *The American Mathematical Monthly* 97.3, с. 242—244.
- Gravner, Janko (б.г.). *Twenty problems in probability*. URL: <https://www.math.ucdavis.edu/~gravner/MAT135A/resources/chpr.pdf>.
- Li, Shuo-Yen Robert (1980). «A martingale approach to the study of occurrence of sequence patterns in repeated experiments». В: *the Annals of Probability*, с. 1171—1176. URL: <https://projecteuclid.org/DPubS?service=UI&version=1.0&verb=Display&handle=euclid.aop/1176994578>.
- Ruggles, Richard и Henry Brodie (1947). «An empirical approach to economic intelligence in World War II». В: *Journal of the American Statistical Association* 42.237, с. 72—91. Подробности про то, как захватив все-го два танка, можно оценить ежемесячный выпуск с точностью в пару процентов.
- Vanderbei, Robert (2011). «Postdoc variant of the secretary problem». В: URL: <https://vanderbei.princeton.edu/tex/PostdocProblem/PostdocProb.pdf>. Подружка Разборчивой невесты хочет второго красавца!
- Wilf, Herbert S (2013). *generatingfunctionology*. Elsevier. URL: <https://www.math.upenn.edu/~wilf/DownldGF.html>. Шикарная книжка про производящие функции.
- Winkler, Peter (2002). «Games people don't play». В: *Puzzler's Tribute: a Feast for the Mind*, с. 301—313. URL: <http://www.teorver.ru/newkatalog/1193689162.pdf>. Несколько красивейших задач с решениями!
- Гусейн-Заде, С.М. (2003). *Разборчивая невеста*. МЦНМО. Задача в изложении для девятиклассников, <http://www.mccme.ru/mmmf-lectures/books/books/book.25.pdf>.
- Колмогоровские студенческие олимпиады по теории вероятностей (б.г.). URL: <http://new.math.msu.su/department/probab/olimpia/olimpia.htm>.
- Свойства случайных перестановок (б.г.). URL: https://en.wikipedia.org/wiki/Random_permutation_statistics.
- Секей, Гabor (1990). *Парадоксы в теории вероятностей и математической статистике*. Москва, Мир. Парадоксы с решениями и историей.

Феллер, Вильям (2010). *Введение в теорию вероятностей и ее приложения*. URSS: Либроком. Очень старый и очень классный учебник по теории вероятностей в двух томах.