

МСПС | Временные ряды |

Опр. Временной ряд - послед. са. величин, зависящих от n (здесь n - время, зависящих от времени).

$\{z_n\}_{n \in \mathbb{N}}$ В дальнейшем - y_n

(1)

Задача прогнозирования: найти ф-цию $f_T: \mathbb{R}^T \rightarrow \mathbb{R}$ $z_{T+d} \approx f_T(z_1, \dots, z_T, d)$

Простейшие методы типа средних регрессии на время, как правило, не помогают, потому что наблюдения ав. ^{сильно} зависящими от близких наблюдений.

Опр. Автокорреляция (оценка ковариан. ф-ции), тогда (ACF)

$$\gamma_k = \gamma_{t, t+k} = \frac{\sum_{t=1}^{T-k} (y_t - \bar{y})(y_{t+k} - \bar{y})}{\sum_{t=1}^T (y_t - \bar{y})^2}, \quad \bar{y} = \frac{1}{T} \sum_{t=1}^T y_t$$

в идеальном случае, то набл. образуют стат. ряд.

Тренд - плавное долгосрочное изменение ур. ряда

Сезонность - циклич. изм. ряда с постоян. периодом.

Ошибка - непрогнозируем. случ. компонента ряда.

Цикл - описывает повторяющиеся, но непериодич. колебания.

Есть алгоритм, который разлагает временной ряд на 3 компонента:

Тренд + циклич. составляющая, сезон. составл. и ошибка: STL-алгоритм (в R - функция `stl`), алгоритм ав. робастности, есть продолжение

идей, алгоритма LOWESS и его более ранней версии Cozss

(сначала с помощью непараметрич. регрессии оцек. тренд, затем циклич. составляющая, а потом ошибки корректируются с помощью процедуры, похожей на ту, что применялось в LOWESS). Написать это самому

тяжело, но если хочется, то читайте оригинал. статью Кливлинга STL: a seasonal-trend decomposition procedure based on Loess.

Сооб., предсказывать с помощью этой штуки тоже можно, ведь непериодич. компонента не меняется, надо лишь оцекить тренд.

Можно на временной ряд посмотреть с другой стороны.

Опр. Врем. ряд (y_1, \dots, y_T) стационарен в уз. смысле, если $\forall z \forall t_1, \dots, t_k$

$(y_{t_1}, \dots, y_{t_k}) \stackrel{d}{=} (y_{t_1+z}, \dots, y_{t_k+z})$. В широком смысле мощность - константа и ковар. ф-ция зависит только от разности: $\gamma(t-s) = \text{cov}(y_t, y_s) = \text{cov}(y_{t+h}, y_{s+h})$

Ряды с трендом и сезонностью нестационарны. Для рядов с непериодич. циклами модель стат. ряда использовать можно.

Как избавиться от тренда и ~~сезонности~~ ^{сезонности} в ряду? А использовать

"дифференцирование": $y'_t = y_t - y_{t-1}$ (это от постоян. тренда), $y''_t = y'_t - y'_{t-1}$

y'_t помогает от лин. тренда; пусть сезон имеет длину s , тогда преобразование $y'_t = y_t - y_{t-s}$ наз. сезонным дифференцированием.

Сезонное дифференцир. лучше делать первым - после него ряд уже может оказаться стационарным (т.е. дальше пометью, как решать задачу прогноза)

Далее мы будем строить разные модели, но требования к остаткам $\hat{\varepsilon}_t = y_t - \hat{y}_t$, \hat{y}_t - наше предсказание, будут одинаковыми:

- 1) Нормальность (и проверяем крит. Стюдента, Уилкоксона и др.) крит. Дики-Фуллера
- 2) Стационарность, т.е. независимость от времени (визуал. анализу, крит. KPSS)
- 3) Неавтокоррелируемость - отсут. неустойчивой зависимости от пред. наблюдений.

Для остатков - отсутствие кластеров положит. и отриц. значений.
Проверяется с помощью коррелограммы, Q-крит. Льюнга-Бокса.
Хелпательное, но не обязатель. св-ва:

- 4) Нормальность (у нас: qq-plot и др.), 5) Гомоскедастичность (Уайт, Бреши-Пэган)

Критерий KPSS: $y_t = \beta_0 + \beta_1 t + z_t + \varepsilon_t$, z_t - сп. блуждание $= \sum_{i=1}^t \eta_i$, $\eta_i \sim iid(0, 1)$, ε_t - стат. $\varepsilon_t = \varepsilon_{t-1} \cdot 0.6^2 \neq 0$

Статистика $KPSS(\varepsilon) = \frac{\sum_{t=1}^T \varepsilon_t^2}{\frac{1}{2} \sum_{t=1}^T \varepsilon_t^2}$, где $S_t = \sum_{i=1}^t \varepsilon_i$; $\hat{\sigma}_\varepsilon^2 = \frac{1}{T} \sum_{i=1}^T \varepsilon_i^2$, ε_t - остатки при регрессии y_t по t и тренду.

(если ε_t не одинаково распр., то $\hat{\sigma}_\varepsilon^2$ заменяется на другую статистику $\hat{\sigma}_\varepsilon^2 = \frac{1}{T} \sum_{t=1}^T \varepsilon_t^2 + \frac{2}{T} \sum_{\tau=1}^{[T]H} (1 - \frac{\tau}{T}) \cdot \sum_{t=\tau+1}^T \varepsilon_t \varepsilon_{t-\tau}$, подробнее см. монографию Maddala, Kim. Unit roots... (1998)). При H_0 статистика KPSS имеет табл. распределение.

Q-критерий Льюнга-Бокса

Пусть $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)$ - ряд ошибок прогноза. $H_0: \varepsilon_1 = \dots = \varepsilon_L = 0$, $H_1: H_0$ неверна
Статистика $Q = T(T+2) \sum_{\tau=1}^L \frac{\hat{\gamma}_\tau^2}{T-\tau}$, при верной H_0 $Q \sim \chi^2_{L-k}$, k - число настраиваемых пар-ров модели ряда.

Коррелограмма - это график автокорреляцион. ф-ции.

Теперь к авторегрессии:

Опр. Процесс авторегрессии порядка p (обозн. $AR(p)$) - это $X_t = a_1 X_{t-1} + a_2 X_{t-2} + \dots + a_p X_{t-p} + \varepsilon_t$, $a_p \neq 0$, где ε_t - проц. Белого шума, не корр. с X_t . $DE\varepsilon = \sigma_\varepsilon^2$

Введём оператор запаздывания $L: LX_t = X_{t-1}$, тогда соотношение $AR(p)$ можно записать так: $a(L)X_t = \varepsilon_t$, где $a(L) = 1 - (a_1 L + a_2 L^2 + \dots + a_p L^p)$

Теорема (Условие стат. $AR(p)$) $a(L)X_t = \varepsilon_t$ стационарен \Leftrightarrow корни $a(z) = 0$ лежат вне $|z| \leq 1$ на комплексной плоскости.

Опр. Процесс скользящего среднего порядка q ($MA(q)$) - это $X_t = \varepsilon_t + b_1 \varepsilon_{t-1} + \dots + b_q \varepsilon_{t-q}$, $b_q \neq 0$, где ε_t - проц. Белого шума с $DE\varepsilon = \sigma_\varepsilon^2$.

Такая модель всегда авт. стат. процессом

Утв. \forall стат. процесс $AR(p)$, задаваемый соотнош. $a(L)(X_t - \mu) = \varepsilon_t$, можно представить в виде проц. $MA(\infty)$. Если X_t имеет $MA(q)$ -представл. $X_t - \mu = b(L)\varepsilon_t = \sum_{j=0}^q b_j \varepsilon_{t-j}$, $b_0 = 1$, и все корни ур-ния $b(z) = 0$ по модулю > 1 (в компл. плоскости) (это наз. условие обратимости), то \exists представл. X_t в виде $AR(\infty)$ модели.

Опр. X_t - стационарный процесс авторегрессии, если $EX_t = 0$ и

$$a(L)X_t = b(L)\varepsilon_t, \text{ т.е. } X_t = \sum_{j=1}^p a_j X_{t-j} + \sum_{j=0}^q b_j \varepsilon_{t-j}, \quad a_p \neq 0, b_q \neq 0.$$

Теорема (разложение Вольфа) \forall стац. в широком смысле ^{детерм.} процесс X_t может быть представлен в виде $X_t = \sum_{j=0}^{+\infty} c_j \varepsilon_{t-j} + \sum_{i=1}^{+\infty} a_i X_{t-i}$, где

ε_t - проц. белого шума, Z_t - стац., детерм. сл. проц., ^{зависит} с ε_t .

Выбор: \forall стац. ряд можно приблизить моделью ARMA(p, q) сколь угодно точно.

Опр. Ряд опис. моделью ARIMA(p, d, q), если ряд его разностей (т.е. дифференцированный) $\nabla^d X_t = (1-L)^d X_t$ описывается моделью ARMA(p, q):

$$a(L) \nabla^d X_t = b(L)\varepsilon_t$$

Можно ещё сезон. ~~ком.~~ добавить - будет SARIMA.

Теперь к выбору коэф. p, d, q и оценке пар-ров модели:

Определим PACF (частот. автокоррел. стац. ряда X_t)

$$\varphi_h = \begin{cases} r(X_{t+h}, X_t) & h=1 \\ r(X_{t+h} - X_{t+h}^{h-1}, X_t - X_t^{h-1}) & h \geq 2 \end{cases}$$

где X_t^{h-1} - регрессия X_t на $X_{t+1}, \dots, X_{t+h-1}$, т.е. $X_t^{h-1} = \beta_1 X_{t+1} + \dots + \beta_{h-1} X_{t+h-1}$ и, наоборот, $X_{t+h}^{h-1} = \alpha_1 X_{t+h-1} + \dots + \alpha_{h-1} X_{t+1}$.

Оценки пар-ров модели:

- 1) При выборе p и q можно использовать ACF и PACF:
 - в модели ARIMA(p, d, 0) ACF экспоненц. затухает или имеет сезон. вид, а PACF значимо отлич. от 0 при $h=p$.
 - в модели ARIMA(0, d, q), наоборот, PACF эксп. затухает или имеет сезон. вид, а ACF значимо отлич. от 0 при $\tau=q$.
- 2) d выбирается так, чтобы ряд был похож на стационарный
- 3) Пусть знаем распр. ε (капр., нормальное) - при заданных p, d, q коэф. модел. оцен. методом макс. пр-я.

Информационные критерии (чем меньше, тем лучше модель)

• $AIC_c = -2L + \frac{2(p+q+1)(p+q+2)}{T-p-q-2}$, где L - ^{логарифм.} ф-ция пр-я

• $BIC = -2L + (\log T - 2)(p+q+1)$. Если модель - AR(p), то лучше заменить $-2L$ на $\ln \sigma_\varepsilon^2$.

В R существует 2 функ. ф-ции: авто.arima и forecast.

И, наконец, как же получить прогноз \hat{X}_{T+h} ?

Каждо сначала оценить последовательно $\hat{X}_{T+1}, \hat{X}_{T+2}$ и т.д., каждый раз обновляя время на 1, коэф. при X_i и ε_i оставляя теми же, новые ошибки ε_{T+1}, \dots приравнивая к 0, а старые заменяя на остатки. Стоит ещё обратить внимание на метод SSA ("Тусеница").