

Содержание.

1. Чем хороша стационарность?
2. Какие бывают стационарные ряды?
3. $MA(q)$ процесс
4. $AR(p)$ процесс
5. $ARMA$ процесс
6. Наконец - зачем действительно нужна стационарность
7. $ARIMA$ процесс
8. $SARIMA$ процесс
9. Как подбирать гиперпараметры модели
10. Современный подход к прогнозированию - переход к задаче регрессии.

Чем хороша стационарность?

Классический пример - продолжите последовательность?

1, 2, 4, 8, 16, ?

Чем хороша стационарность?

Правильный ответ - 42!

Чем хороша стационарность?

Если ряд не постоянен во времени, мы не можем делать никаких осмысленных прогнозов.

Стационарность дает нам гарантии того, что:

$$\text{Var}(Y_t) = \text{const}$$

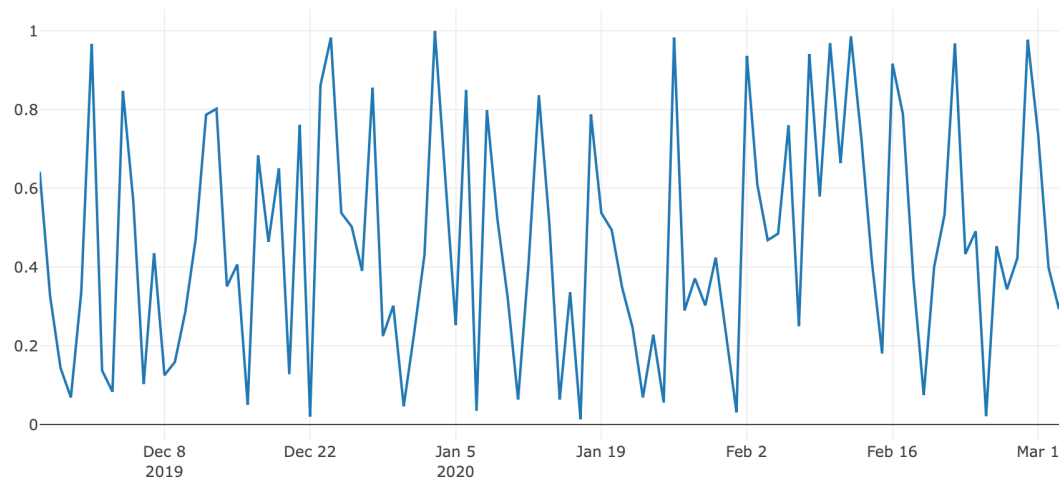
$$E(Y_t) = \text{const}$$

$$\text{Cov}(Y_t, Y_{t-k}) = \text{const}$$

Какие бывают стационарные ряды

Самый простой тип - белый шум. $Y_t = \epsilon_t, \epsilon_t \sim N(0, \sigma^2)$

Всегда можем сказать, что этот ряд имеет матожидание 0 и дисперсию σ^2



MA(q) процесс

Следующий тип стационарный рядов - MA(q) процесс (Moving Average)

$$y_t = \alpha + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}$$

То есть каждая точка моделируется как линейная комбинация q предыдущих шумовых компонент

Большинство MA(q) процессов будут стационарными

*Примечание о шуме

Здесь и далее будет предполагаться, что шум это гауссов (нормально распределенный) шум с нулевым матожиданием и дисперсией σ^2

$$e_t \sim N(0, \sigma^2)$$

AR(p) - процесс

Данный процесс представляет собой зависимость каждой точки ряда от P предыдущих точек

$$y_t = \alpha + \phi_1 y_{t-1} + \phi_2 y_{t-2} + \dots + \phi_p y_{t-p} + \varepsilon_t$$

Не любой AR процесс - стационарен

Для того, чтобы данный процесс был стационарен, необходимо выполнение следующего условия - необходимо, чтобы коэффициенты ϕ лежали на единичном круге, т.е., например

в AR(1) необходимо $-1 < \phi_1 < 1$;

в AR(2) необходимо $-1 < \phi_2 < 1$, $\phi_1 + \phi_2 < 1$, $\phi_2 - \phi_1 < 1$

ARMA процесс

Комбинация AR(p) и MA(q) процессов называется ARMA(p, q)

$$ARMA(p, q): y_t = \alpha + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q}$$

Истина - зачем нужна стационарность

По теореме Вольда - любой стационарный ряд с любой наперед заданной точностью может быть смоделирован моделью $ARMA(p, q)$!

Таким образом, сделав ряд стационарным, мы можем подобрать какой модели он соответствует.

Выбрать эту модель для прогнозирования.

И восстановить исходный ряд обратными преобразованиями.

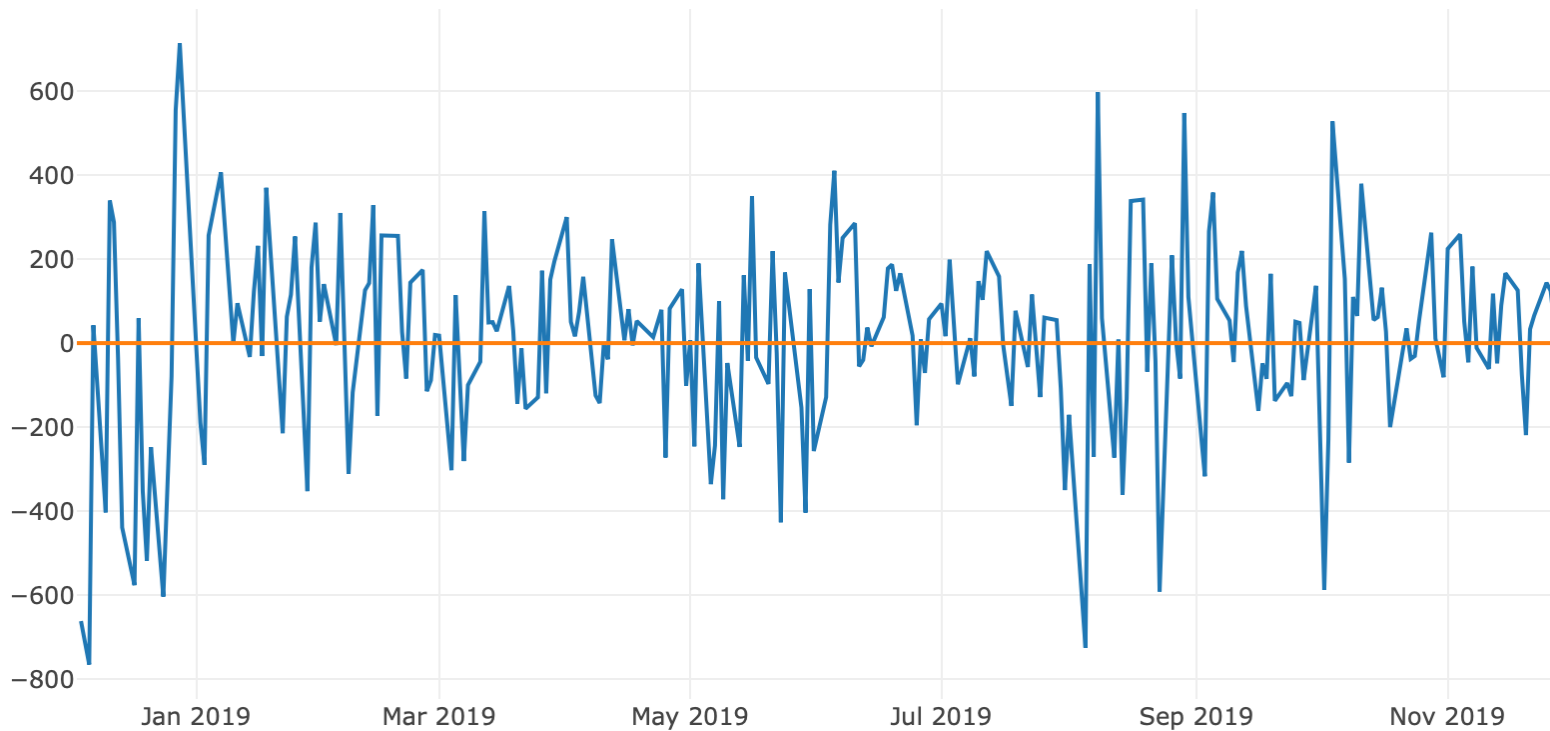
Как это правильно понимать?

Это что, мы теперь сможем предсказывать белый шум?

Нет, разумеется.

Теорема Вольда лишь говорит, что мы сможем понять какой модели $ARMA(p, q)$ ряд соответствует. То есть, говоря простыми словами, для, например, белого шума, ответ будет таков: ваш ряд описывается моделью $ARMA(0, 0)$.

Как это правильно понимать?



Модель ARIMA

Если ARMA работает для всех стационарных рядов, то почему не сделать модель, которая работает для всех рядов, что можно сделать стационарными, продифференцировав?

Ряд описывается моделью $ARIMA(p, d, q)$, если d раз продифференцированный ряд описывается моделью $ARMA(p, q)$.

Как это правильно понимать?

Или индекс Доу-Джонса, допустим.



SARMA

Окей, моделью ARIMA мы можем смоделировать все стационарные ряды, а также все не стационарные которые можно сделать стационарными дифференцированием.

Проблема - сезонность не всегда можно убрать дифференцированием.

Решение - добавить в модель ARMA сезонные компоненты.

$$\begin{aligned} y_t = & \alpha + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q} \\ & + \phi_S y_{t-S} + \phi_{2S} y_{t-2S} + \dots + \phi_{PS} y_{t-PS} \\ & + \theta_S \varepsilon_{t-S} + \theta_{2S} \varepsilon_{t-2S} + \dots + \theta_{PS} \varepsilon_{t-PS} \end{aligned}$$

Возьмем все вместе и получим SARIMA

Ряд описывается моделью SARIMA(p, d, q)(P, D, Q), если d раз обычно и D раз сезонно продифференцированный ряд описывается моделью SARMA(p, q)(P, Q)

$$\begin{aligned} y_t = & \alpha + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \varepsilon_t + \theta_1 \varepsilon_{t-1} + \theta_2 \varepsilon_{t-2} + \dots + \theta_q \varepsilon_{t-q} \\ & + \phi_S y_{t-S} + \phi_{2S} y_{t-2S} + \dots + \phi_{PS} y_{t-PS} \\ & + \theta_S \varepsilon_{t-S} + \theta_{2S} \varepsilon_{t-2S} + \dots + \theta_{PS} \varepsilon_{t-QS} \end{aligned}$$

Как подбирать параметры модели?

Как подобрать α, ϕ, θ (при условии зафиксированных гиперпараметров)?

По методу наименьших квадратов, так же как и в линейной регрессии.

Как подбирать параметры модели?

Порядки дифференцирования d , D ?

Пока ряд не станет стационарным (по критерию Дики-Фуллера, например)

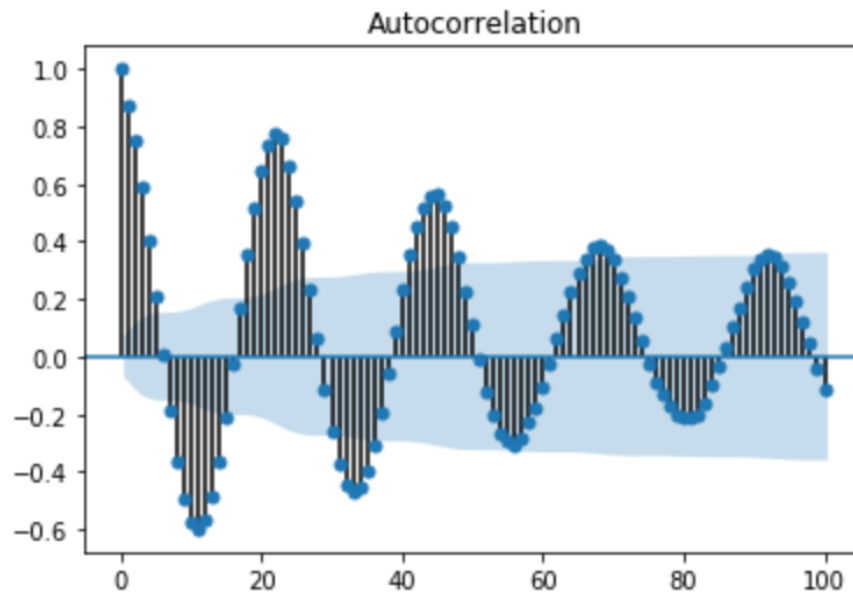
Как подбирать параметры модели?

q , Q ?

По графику автокорреляционной функции.

Q - последний значимый сезонный лаг

q - последний значимый не сезонный лаг



Как подбирать параметры модели?

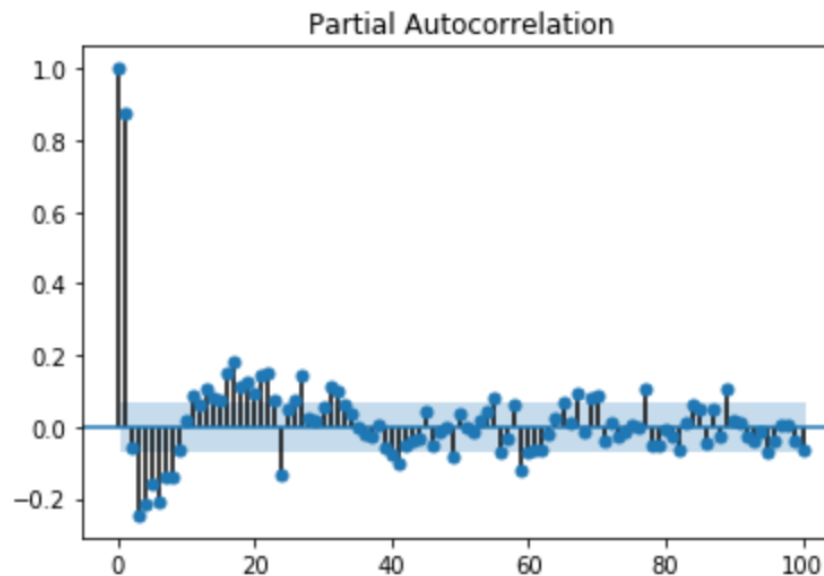
p , P ?

По графику частичной автокорреляционной функции.

*Частичная автокорреляция - автокорреляция после снятия регрессии на промежуточные значения.

p - последний значимый сезонный лаг

P - последний значимый не сезонный лаг



Критерий Акаике

Выбор всех предыдущих параметров был приблизительным, что использовать для сравнения разных параметров?

Метод максимального правдоподобия не подходит, так как приводит к переобучению модели, просто выбирая максимальные значения параметров.

Поэтому используется информационный критерий Акаике.

$$AIC = 2k - 2\ln(L)$$

Алгоритм

