

# Анализ временных рядов

## часть 2

Алёна Ким

НИУ ВШЭ

9 ноября, 2018

# Автокорреляционная функция (ACF)

**Временной ряд**  $Y^T: y_1, \dots, y_T, \dots, y_t \in \mathbb{R}$  - значения признака, измеренные через постоянные временные интервалы.

**Автокорреляция:**

$$r_\tau = r_{y_t t_{t+\tau}} = \frac{\sum_{t=1}^{T-\tau} (y_t - \bar{y})(y_{t+\tau} - \bar{y})}{\sum_{t=1}^T (y_t - \bar{y})^2}, \bar{y} = \frac{1}{T} \sum_{t=1}^T y_t.$$

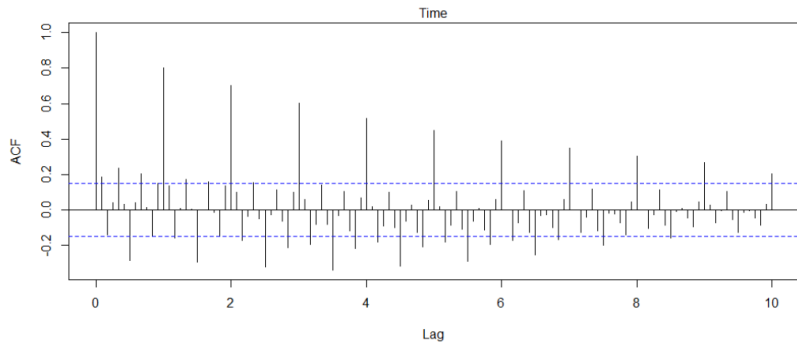
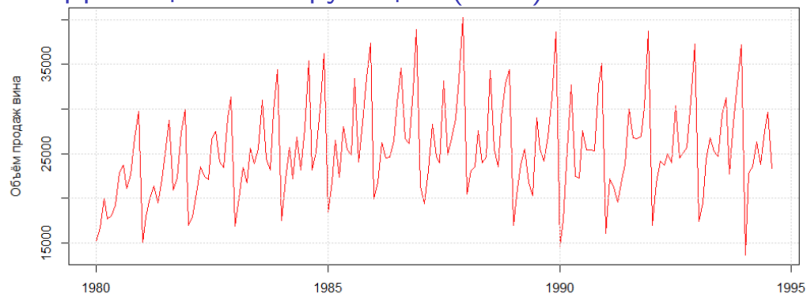
$r_\tau \in [-1, 1]$ ,  $\tau$  - лаг автокорреляции

Нулевая гипотеза:  $H_0 : r_\tau = 0$ ;

Альтернативная гипотеза:  $H_1 : r_\tau \neq 0$ ;

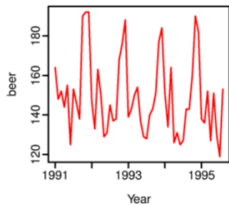
Статистика:  $T(Y^T) = \frac{r_\tau \sqrt{T-\tau-2}}{\sqrt{1-r_\tau^2}}$ ;

# Автокорреляционная функция (ACF)



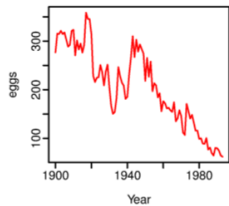
# Нестационарность

## ► СЕЗОННОСТЬ



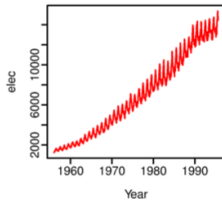
# Нестационарность

- ▶ тренд



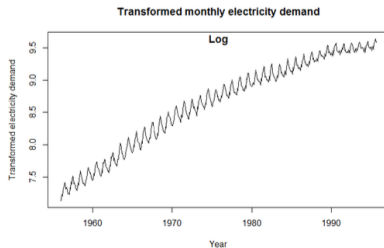
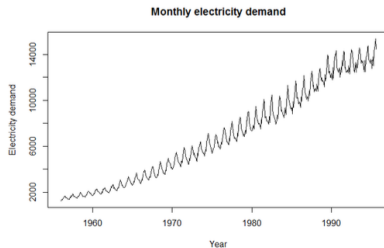
# Нестационарность

- ▶ меняющаяся дисперсия



# Стабилизация дисперсии

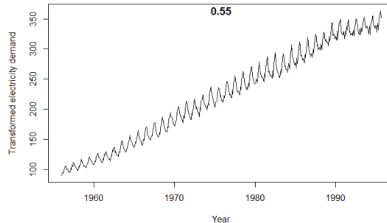
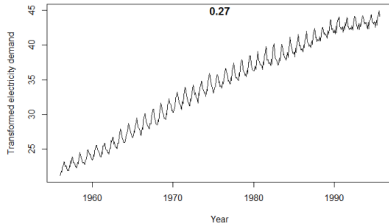
Логарифмирование как стабилизирующее преобразование:



# Преобразования Бокса-Кокса

$$y'_t = \begin{cases} \ln y_t, & \lambda = 0, \\ \frac{y_t^\lambda - 1}{\lambda}, & \lambda \neq 0, \end{cases}$$

Параметр  $\lambda$  выбирается так, чтобы минимизировать дисперсию или максимизировать правдоподобие





# Преобразования Бокса-Кокса

$$\hat{y}_t = \begin{cases} \exp(\hat{y}'_t), \lambda = 0, \\ (\lambda \hat{y}'_t + 1)^{\frac{1}{\lambda}}, \lambda \neq 0, \end{cases}$$

Если какие-то  $y_t < 0$ , преобразования Бокса-Кокса невозможны

# Дифференцирование ряда

Временной ряд  $Y^T: y_1, \dots, y_T$

Дифференцирование:

$$y_1, \dots, y_T \rightarrow y'_2, \dots, y'_T$$

$$y'_t = y_t - y_{t-1}$$

Зачем?

- ▶ стабилизирует значение ряда
- ▶ избавляет от сезонности/тренда

Неоднократное дифференцирование

$$y_1, \dots, y_T \rightarrow y'_2, \dots, y'_T \rightarrow y''_3, \dots, y''_T$$

# Сезонное дифференцирование ряда

Переход к попарным разностям значений в соседних сезонах

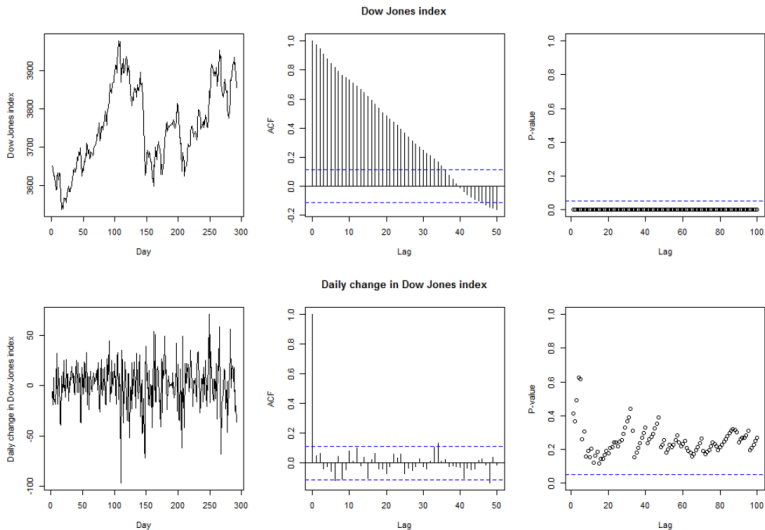
$$y_1, \dots, y_T \rightarrow y'_{s+1}, \dots, y'_T$$

$$y'_t = y_t - y_{t-s}$$

Обычное и сезонное дифференцирование можно выполнять в любом порядке.

Если ряд кажется сезонным, сначала лучше выполнить сезонное дифференцирование. После этого ряд может уже быть стационарным

# Дифференцирование ряда



Критерий KPSS:

Для исходного ряда  $p\text{-value} < 0.01$

Для ряда первых разностей  $p\text{-value} > 0.1$

# Критерий KPSS (Kwiatkowski-Philips-Schmidt-Shin)

ряд ошибок прогноза:  $\epsilon^T = \epsilon_1, \dots, \epsilon_T$ ;

нулевая гипотеза:  $H_0$  : ряд  $\epsilon^T$  стационарен;

альтернативная гипотеза:  $H_1$  : ряд  $\epsilon^T$  описывается моделью вида  $\epsilon_t = \alpha \epsilon_{t-1}$ ;

статистика:  $KPSS(\epsilon^T) = \frac{1}{T^2 \lambda^2} \sum_{i=1}^T (\sum_{t=1}^i \epsilon_t)^2$ ;

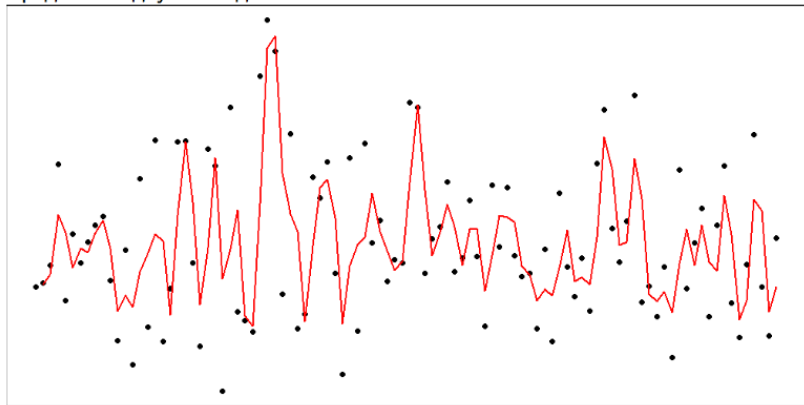
# Скользящее среднее

Пусть у нас есть независимый одинаково распределённый во времени шум  $\epsilon_t$ :



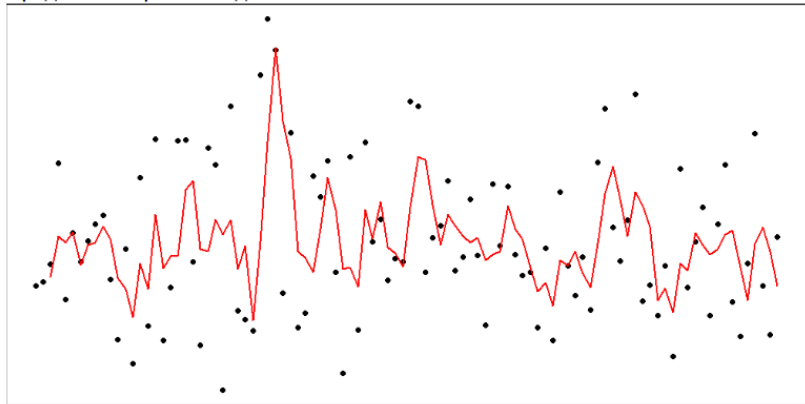
# Скользящее среднее

Среднее по двум соседним точкам:



# Скользящее среднее

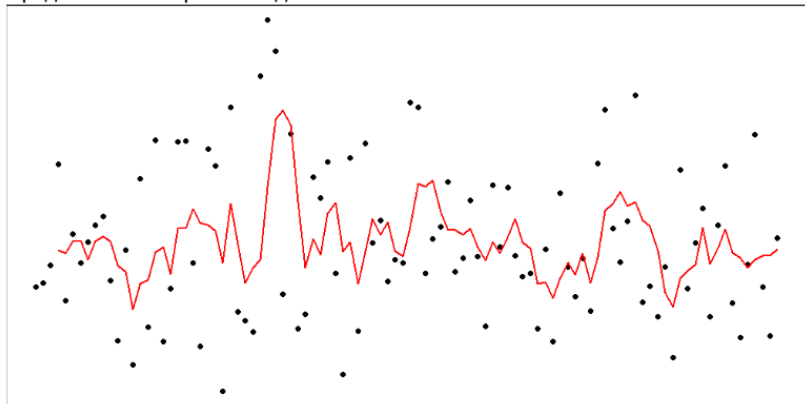
Среднее по трём соседним точкам:





# Скользящее среднее

Среднее по четырём соседним точкам:



## AR(p), MA(q)

$$AR(p) : y_t = \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \epsilon_t,$$

$$MA(q) : y_t = \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q},$$

где  $y_t$  – стационарный ряд с нулевым средним,  $\phi_t, \theta_t$  – константы,  $\epsilon_t$  – гауссов белый шум с нулевым средним и постоянной дисперсией  $\sigma_\epsilon^2$

Если среднее равно  $\mu$  :

$$AR(p) : y_t = \alpha + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \epsilon_t,$$

$$\alpha = \mu(1 - \phi_1 - \dots - \phi_p)$$

$$MA(q) : y_t = \mu + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q},$$

## ARMA (Autoregressive moving average)

$$ARMA(p, q) : y_t = \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q},$$

где  $y_t$  – стационарный ряд с нулевым средним,  $\phi_t, \theta_t$  – константы,  $\epsilon_t$  – гауссов белый шум с нулевым средним и постоянной дисперсией  $\sigma_\epsilon^2$

Если среднее равно  $\mu$ :

$$y_t = \alpha + \phi_1 y_{t-1} + \dots + \phi_p y_{t-p} + \epsilon_t + \theta_1 \epsilon_{t-1} + \dots + \theta_q \epsilon_{t-q},$$

$$\alpha = \mu(1 - \phi_1 - \dots - \phi_p)$$

Другой способ записи:

$$\phi(B)y_t = \theta(B)\epsilon_t,$$

где  $B$  – разностный оператор ( $By_t = y_{t-1}$ )

$$\phi(B)y_t = (1 - \phi_1 B - \phi_2 B^2 - \dots - \phi_p B^p)y_t =$$

$$= \epsilon_t(1 + \theta_1 B + \dots + \theta_q B^q) = \theta(B)\epsilon_t$$

# Теорема Вольда

Любой стационарный ряд может быть аппроксимирован моделью  $ARMA(p, q)$  с любой точностью

# ARIMA (Autoregressive integrated moving average)

Ряд описывается моделью  $ARIMA(p, d, q)$ , если ряд его разностей

$$\nabla^d y_t = (1 - B)^d y_t$$

описывается моделью  $ARMA(p, q)$

$$\phi(B)\nabla^d y_t = \theta(B)\epsilon_t$$

## Подбор коэффициентов

- ▶  $\alpha, \theta, \epsilon$ 
  - если все остальные параметры фиксированы, коэффициенты регрессии подбираются методом наименьших квадратов
  - чтобы найти коэффициенты  $\theta$ , шумовая компонента предварительно оценивается с помощью остатков авторегрессии
  - если шум белый, то МНК даёт оценки максимального правдоподобия
- ▶  $d$ 
  - порядки дифференцирования подбираются так, чтобы ряд стал стационарным
  - чем меньше раз продифференцируем, тем меньше будет дисперсия итогового прогноза
- ▶  $p, q$ 
  - начальные приближения можно выбрать с помощью автокорреляций

# Частичная автокорреляционная функция (PACF)

**Частичная автокорреляция:** стационарного ряда  $y_t$  – автокорреляция остатков авторегрессии предыдущего порядка

$$\phi_h = \begin{cases} r(y_{t+1}, y_t), & h = 1, \\ r(\epsilon_{t+h}, \epsilon_t), & h > 1, \end{cases}$$

$p$  – номер последнего лага при котором PACF значима

$q$  – номер последнего лага при котором ACF значима

# Построение ARIMA

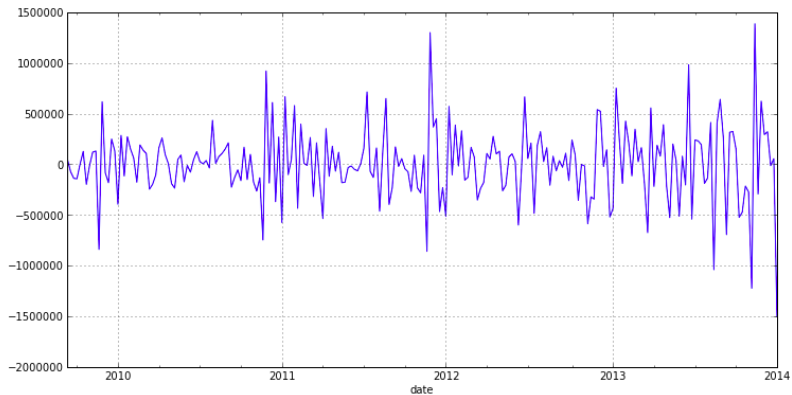




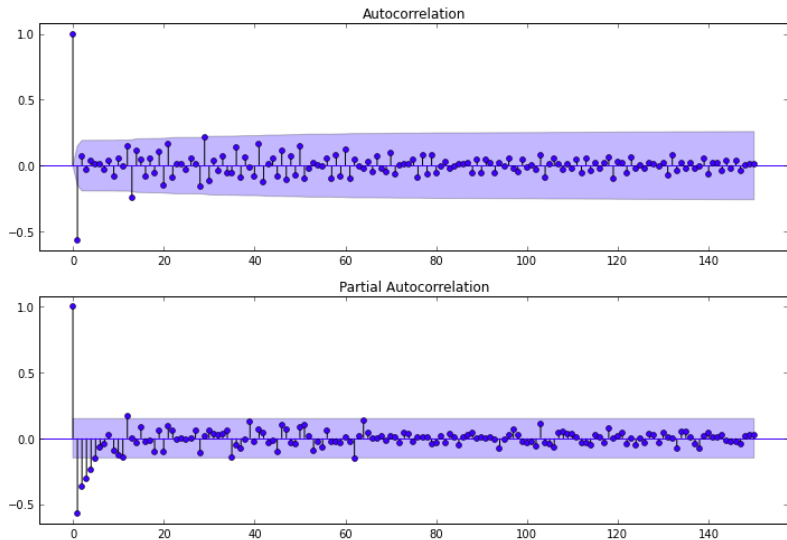
# Построение ARIMA



# Построение ARIMA

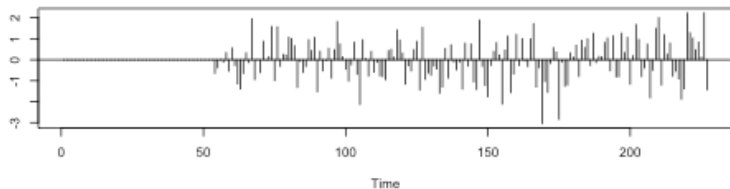


# Построение ARIMA

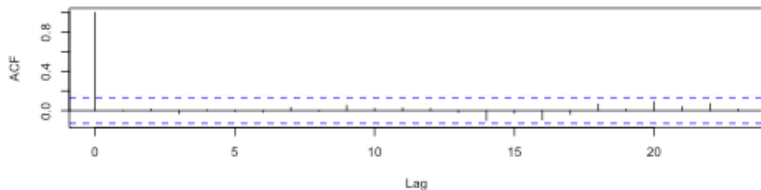


# Построение ARIMA

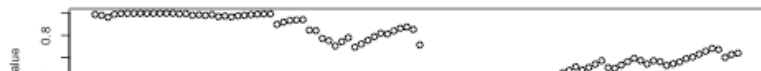
Standardized Residuals



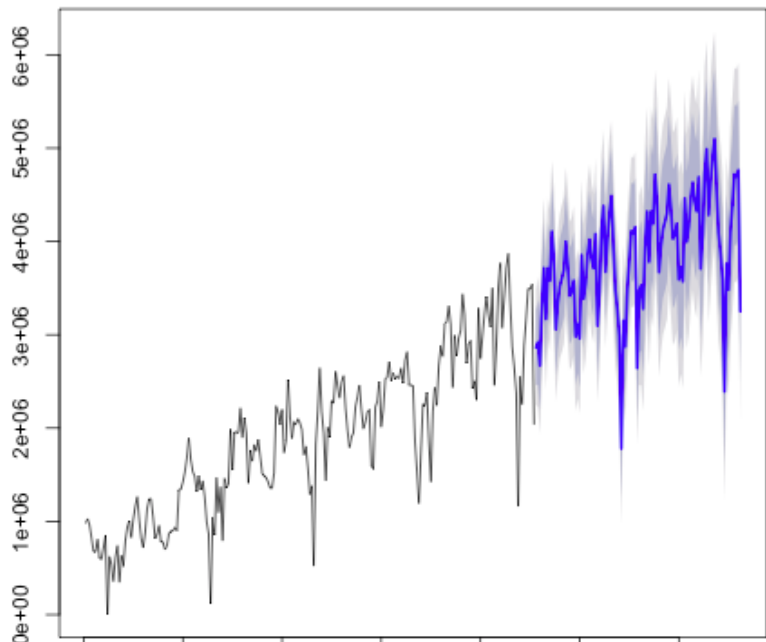
ACF of Residuals



p values for Ljung-Box statistic



## Построение ARIMA(4,1,13)



# Дискретное преобразование Фурье

## DFT, Discrete Fourier Transform

$$X_k = \sum_{n=0}^{N-1} x_n e^{-\frac{2\pi i}{N} kn} =$$

$$\sum_{n=0}^{N-1} x_n \cdot \left[ \cos\left(\frac{2\pi kn}{N}\right) - i \cdot \sin\left(\frac{2\pi kn}{N}\right) \right],$$

$$(k = 0, \dots, N - 1)$$

### Ограничения:

При анализе с помощью преобразования Фурье, мы исходим из предположения, что он периодический на текущем временном интервале и состоит из элементарных синусоид.

# Дискретное преобразование Фурье



- ▶ Получаем комплексные числа, содержащий информацию об амплитудном и фазовом спектрах анализируемого кадра. Причём спектры также являются дискретными с шагом (частота дискретизации)/(N отсчётов).
- ▶ Разложение в базис в некотором пространстве периодических функций. Коэффициенты - веса этих функций.
- ▶ Выполняется за  $O(N \log(N))$  - БПФ.

- [1].  
<http://www.machinelearning.ru/wiki/images/3/31/Psad<sub>t</sub>s<sub>a</sub>rima<sub>2</sub>017.pdf>
- [2]. <https://habr.com/post/210530/>
- [3]. <https://ru.wikipedia.org>