

(w07)

как выбрать модель?

Идея!

кросс-валидация.

Модели = { Модель А, Модель Б }

$$MSE_M = E((y_{T+1} - \hat{y}_{T+1}^M)^2)$$

теоретическая  
MSE

CV с разгущенным окном

$e_{F+1} = y_{F+1} - \hat{y}_{F+1|F}$   
 $e_{F+2} = y_{F+2} - \hat{y}_{F+2|F+1}$   
 $e_{F+3} = y_{F+3} - \hat{y}_{F+3|F+2}$   
 $e_{F+4} \dots$

$$\hat{MSE}_M = \sum_{h=1}^H e_{F+h}^2 / H$$

Идеальный способ!

- \* независимый (если много моделей/разов)
- \* хотелся много наблюдений!

(хотим: \*  $H$  помеченные, чтобы  
модель обучать на более длинном  
ряду \*  $H$  побольше, чтобы  
точно оценивать MSE  
всего).

\* Способы, когда нужна CI горга:

- \* формульные статистические тесты
- \* информационные критерии (AIC)
- \* максимальные правдоподобия.

Автоматическая сезонная ARIMA.

год - месячный / кварт-ый

①. Смотрим по ряду «сезонности»  
(максимум)

Если  $\text{сез-сдв} > \text{порог} = 0.8$

то переходим от  $y_t \rightarrow \Delta_s y_t =$   
 $= y_t - y_{t-k}$  \* может и другое !!

②. Формальный статистический тест \*\*  
KPSS

$H_0: (y_t) \sim \text{стационарен}$

т.е. после  $k$ го шага  
возможно это уже  
от некого блага  $\Delta$

$H_1: (y_t)$  не стационар, но  
есть тренд  $(\Delta y_t)$  стационар

Если  $H_0$  отвергается то делаем  
обычную разницу

$(y_t) \rightarrow (\Delta y_t)$

т.е.  $y_t$  после возможного  
преобразования не имеет  
тренда и дисперсия постоянна

перв 3-м марку [мы рассматриваем]  
 перв будем сразу сдв прогресс.

Шаг 3  $y_t \rightarrow \begin{pmatrix} d & D \\ \Delta & \Delta_S \end{pmatrix} y_t$   $d \in \{0, 1, 2\}$   
 $D \in \{0, 1\}$

оцениваем несколько возможных  $d$  и  $D$   
 SARIMA модель (сразу) для  $\Delta \Delta_S y_t$

например

$$\begin{aligned} y_t &\sim \text{SARIMA}(0, d, 1) (0, D, 1) \\ y_t &\sim \text{SARIMA}(1, d, 0) (0, D, 1) \end{aligned}$$

$\vdots$

$$y_t \sim \text{SARIMA}(2, d, 0) (2, D, 0)$$

модель не является  
 оптимальной

\*возможно надо  
 переопределить модель

$$\text{SARIMA}(p, d, q) (P, D, Q)$$

неог.
сезон

$$\begin{aligned} p+q &\leq 3 \\ P+Q &\leq 3 \end{aligned}$$

выбирается модель с мин AIC.

$$AIC = 2p - 2 \ln L$$

$p$  - число параметров  
 свободных параметров  
 в модели

$\ln L$  - логарифм. максим.  
 лн правды

Detail

that

time series

$$(y_t) \rightarrow \boxed{STL} \rightarrow (\text{trend}) + (\text{seas}) + (\text{rem})$$

$$\underline{F_{seas}} = \max\left(0, 1 - \frac{\text{stoe}(\text{rem})}{\text{stoe}(\text{seas} + \text{rem})}\right)$$

$\text{stoe}(a)$  =  $\text{выборочная дисперсия}$   $\text{по } a$

sample  $\rightarrow$

$$= \frac{\sum (a_i - \bar{a})^2}{n-1}$$

# Шаг 2. Говорим про KPSS тест.

(идея)

$y_1, \dots, y_T$  — зависимые наблюдения

интуитивно: несут меньше информации  
чем независимые наблюдения  
способы циприя?

- \* информация Фишера
- \* гетеросkedастическая дисперсия (long run)
- \* эргодическое число наблюдений

идеальный мир  
равных независимых точек

$y_1, \dots, y_T \sim \text{незав}$  с  $E(y_t) = \mu$  и  $\text{Var}(y_t) = \frac{\lambda^2}{T}$

95% доверительный интервал для  $\mu$ :  $\left[ \bar{y} - 1.96 \cdot \sqrt{\frac{\lambda^2}{T}}, \bar{y} + 1.96 \cdot \sqrt{\frac{\lambda^2}{T}} \right]$

$$\text{Var}(\bar{y}) = \frac{\lambda^2}{T}$$

$$\text{Var}(y_t) = \lambda^2$$

$y_1, \dots, y_T$  — стационарный процесс

$$\text{Var}(y_t) = \gamma_0$$

$$\text{Var}(\bar{y}) = \frac{\text{Var}(y_1 + \dots + y_T)}{T^2} =$$

$$= \frac{T \cdot \gamma_0 + 2(T-1) \cdot \gamma_1 + \dots + 2 \cdot \gamma_{T-1}}{T^2}$$

$$\gamma_k = \text{Cov}(y_t, y_{t+k})$$

$$\neq \frac{\gamma_0}{T}$$

[def] [гетеросkedастическая дисперсия  $\frac{\lambda^2}{T}$ ] — т.е. такое, что  
[стационарный процесс]  $\text{Var}(\bar{y}) = \frac{\lambda^2}{T} + o\left(\frac{1}{T}\right)$

def  $T_{eff}$  - среднее число наблюдений  
 сур-ср у формулы:  
 $Var(\bar{y}) = \frac{\sigma^2}{T_{eff}}$

Ynp.  $y_t \approx \underbrace{u_t}_{\sim N(0;16)} + \underbrace{u_{t-1}}_{\text{д. шум}}$   
 а)  $Var(y_t)$ ? б) полное ли смешение?  
 в)  $\lim_{T \rightarrow \infty} \frac{T_{eff}}{T}$  ? =  $\frac{1}{2}$

$$Var(y_t) = 32$$

$$Var\left(\frac{y_1 + y_2 + \dots + y_T}{T}\right) = \frac{T \cdot 32 + 2 \cdot (T-1) \cdot 16 + 0}{T^2} =$$

$$= \frac{64}{T} - \frac{32}{T^2} = \frac{64}{T} + o\left(\frac{1}{T}\right)$$

$$1^2 = 64$$

KPSS тест с константой

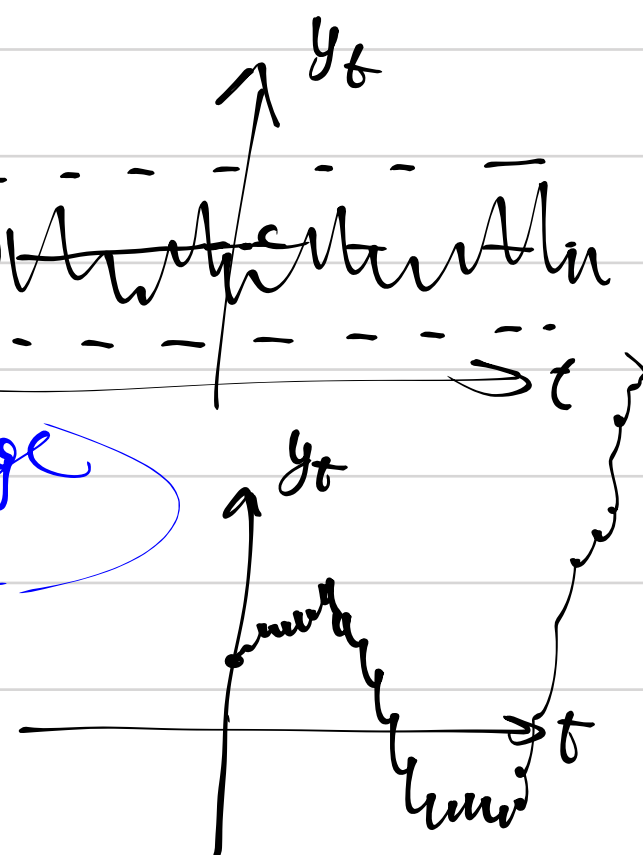
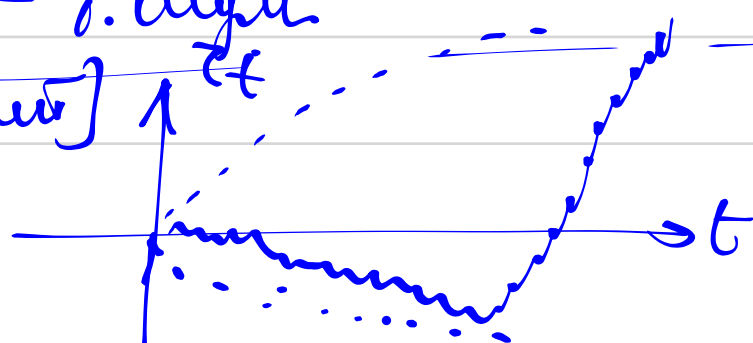
Предположим:  $y_t = \alpha + z_t + x_t$ , где

$(x_t) \sim$  стационарный процесс с  $E(x_t) = 0$ .  
 независимый с  $z_t$

$H_0: z_t \equiv 0$ .  $(y_t)$  - случайн.

$H_1: z_t = u_1 + u_2 + u_3 + \dots + u_t$ , где  
 $(u_t)$  - д. шум

$[z_t - \text{не стационарный}]$





формула процедура:

$$\hat{c} = \bar{y} = \frac{y_1 + y_2 + \dots + y_T}{T}$$

можем

$$KPSS = \frac{\sum_{t=1}^T S_t^2}{T^2 \cdot \hat{\lambda}^2}$$

$$S_t = y_1 + y_2 + \dots + y_t = t \cdot \bar{y}$$

↑ накопленная  
сумма ошибок  
прогноза приваивной

$\hat{\lambda}^2$  — оценка  
долгосрочной  
дисперсии.

$KPSS \xrightarrow[T \rightarrow \infty]{\text{К. Верна}} \text{свободное распределение } [KPSS^c]$

---