

СРС: Моделирование тенденции развития временного ряда

Решение типового примера

На сайте Федеральной службы государственной статистики (www.gks.ru) выберем ряд динамики некоторого показателя для анализа и моделирования тенденции.

Например, в разделе «Официальная статистика/Публикации/Каталог публикаций/Статистические сборники/Социально-экономические показатели Российской Федерации в 1991 - 2011гг. (приложение к статистическому сборнику «Российский статистический ежегодник. 2012»)» скачиваем файл «Социально-экономические показатели Российской Федерации в 1991-2011гг. (0,9 Мб)» в формате .xls. На листе «Раз.2» книги выбираем показатель «Выбросы загрязняющих веществ в атмосферный воздух от автотранспорта), млн.т.»:

Год	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	
y_t	17,3	22	19	13,5	11	11	11,3	11,8	12,2	13,5	
Год	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011
y_t	14,2	14,4	14,8	15,3	15,4	14,7	14,7	13,6	13,5	13,1	13,3

График представлен на рис. 1:

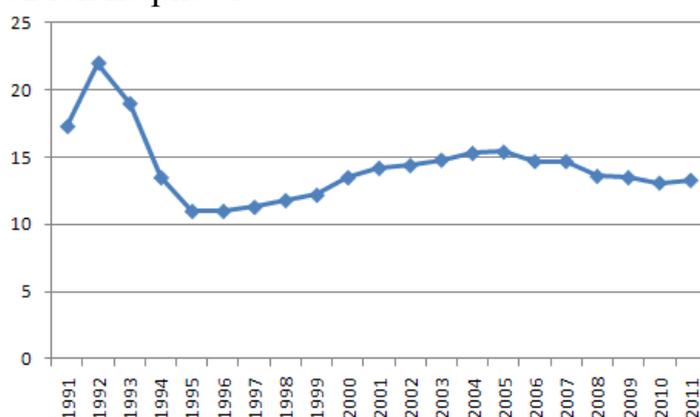


Рис. 1. Выбросы загрязняющих веществ в атмосферный воздух от автотранспорта, млн.т., 1991-2011 гг.

При выборе показателя учитывали следующее:

- длина временного ряда (чем больше, тем лучше);
- сопоставимость уровней временного ряда (показатели, выраженные в денежных единицах лучше не использовать, т.к. за 20 лет в РФ темпы инфляции существенно менялись);
- не тривиальность поведения показателя (значение показателя со временем заметно меняется и изменение более сложное, чем небольшой постоянный рост или снижение);
- субъективный интерес исследователя.

Выполним следующие расчеты:

- 1) Проверим наличие тенденции у временного ряда.
- 2) Проведем аналитическое выравнивание временного ряда.
- 3) Проверим адекватность и точность модели.

Проверим наличие тенденции у временного ряда с помощью критерия серий, основанного на медиане выборки

1. Из исходного ряда с уровнями y_1, y_2, \dots, y_n образуем ранжированный (вариационный) ряд y'_1, y'_2, \dots, y'_n :

11; 11; 11,3; 11,8; 12,2; 13,1; 13,3; 13,5; 13,5; 13,5; 13,6; 14,2; 14,4; 14,7; 14,7; 14,8; 15,3; 15,4; 17,3; 19; 22

2. Определяем медиану (Me) этого вариационного ряда. В случае нечетного значения длины ряда $n = 2m + 1$: $Me = y'_{m+1}$. Если объем ряда четный $n = 2m$, то $Me = (y'_m + y'_{m+1})/2$.

Так как объем ряда $n = 21$ – нечетный, то $Me = y'_{m+1} = y'_{11} = 13,6$.

3. Построим последовательность δ_i из плюсов и минусов по следующему правилу:

$$\delta_i = \begin{cases} +, & \text{если } y_t > Me, & t = \overline{1, n}; \\ -, & \text{если } y_t < Me, & t = \overline{1, n}. \end{cases}$$

Если значение уровня исходного ряда y_t равно медиане, то это значение пропускается.

Год	1991	1992	1993	1994	1995	1996	1997	1998	1999	2000	
y_t	17,3	22	19	13,5	11	11	11,3	11,8	12,2	13,5	
δ_i	+	+	+	-	-	-	-	-	-	-	
Год	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011
y_t	14,2	14,4	14,8	15,3	15,4	14,7	14,7	13,6	13,5	13,1	13,3
δ_i	+	+	+	+	+	+	+		-	-	-

4. Подсчитываем $\nu(n)$ – число серий в совокупности δ_i , где под серией понимается последовательность подряд идущих плюсов или минусов. Один плюс или один минус тоже будет считаться серией:

$$\nu(n) = 4.$$

Определяем $\tau_{max}(n)$ – протяженность самой длинной серии:

$$\tau_{max}(n) = 7.$$

5. Проверка гипотезы основывается на том, что при условии случайности ряда (при отсутствии систематической составляющей) протяженность самой длинной серии не должны быть слишком большой, а общее число серий – слишком маленьким. Поэтому для того, чтобы не была отвергнута гипотеза о случайности исходного ряда (об отсутствии систематической составляющей), должны выполняться следующие неравенства:

$$\nu(n) > \left[\frac{1}{2} (n + 1 - 1,96\sqrt{n - 1}) \right];$$

$$\tau_{max}(n) < [1,43 \ln(n + 1)],$$

где n – длина временного ряда;

$\nu(n)$ – число серий;

$\tau_{max}(n)$ – число подряд идущих плюсов или минусов в самой длинной серии.

Квадратные скобки в правой части неравенств означают целую часть числа (целая часть числа x , т.е. $[x]$ – это целое число, ближайшее к x и не превосходящее его).

Если хотя бы одно из записанных неравенств нарушается, то гипотеза о случайности исходного ряда отвергается с вероятностью ошибки α заключенной между 0,05 и 0,0975 (следовательно, подтверждается наличие зависящей от времени неслучайной составляющей).

Проверим выполнение записанных неравенств для нашего примера:

$$\begin{aligned}v(n) = 4 &> \left[\frac{1}{2} (21 + 1 - 1,96\sqrt{21 - 1}) \right]; \\4 &> \left[\frac{1}{2} (22 - 1,96\sqrt{20}) \right]; \\4 &> [6,62] = 6.\end{aligned}$$

Получили неверное неравенство, значит гипотезу о случайности исходного ряда отвергаем на уровне значимости $\alpha = 0,0975$. Принимаем гипотезу о наличии тенденции в данном временном ряду.

Для полноты проведения расчета проверим выполнение второго неравенства:

$$\begin{aligned}\tau_{max}(n) = 7 &< [1,43 \ln(21 + 1)], \\7 &< [4,42] = 4.\end{aligned}$$

Получили неверное неравенство, что подтверждает сделанный выше вывод.

Проведем аналитическое выравнивание временного ряда средствами MS Excel

Построим график временного ряда (рис. 1).

Выполним правый щелчок мыши по графику (по самому графику или по точкам, изображающим уровни ряда) и в появившемся контекстном меню выберем пункт «Добавить линию тренда». Появится окно диалога, изображенное на рис. 2.

Для повышения информативности необходимо установить флажки:

- показывать уравнение на диаграмме;
- поместить на диаграмму величину достоверности аппроксимации (R^2) (т.е. коэффициент детерминации).

Последовательно добавим следующие тренды:

- линейный;
- экспоненциальный;
- логарифмический;
- полиномиальный (степень равна 2);
- полиномиальный (степень равна 3);
- степенной.

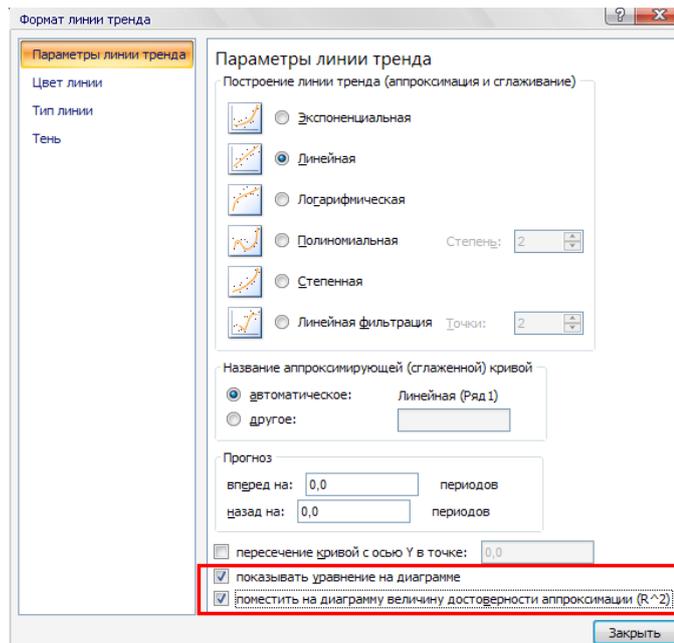


Рис. 2. Окно диалога «Формат линии тренда»

Результаты представлены на рис.3-8:

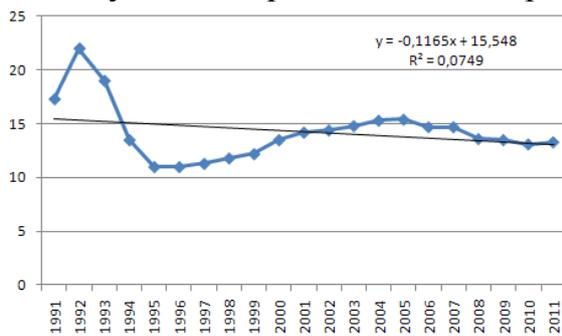


Рис. 3. Линейный тренд

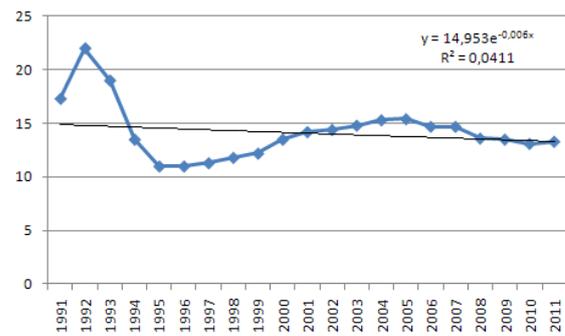


Рис. 4. Экспоненциальный тренд

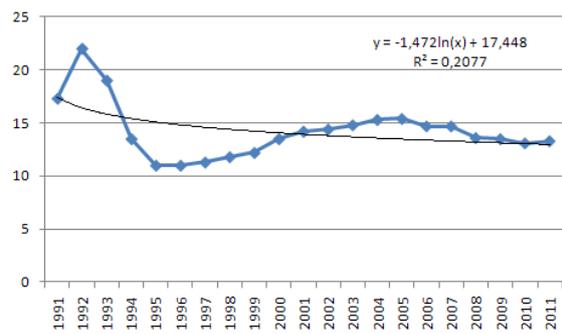


Рис. 5. Логарифмический тренд

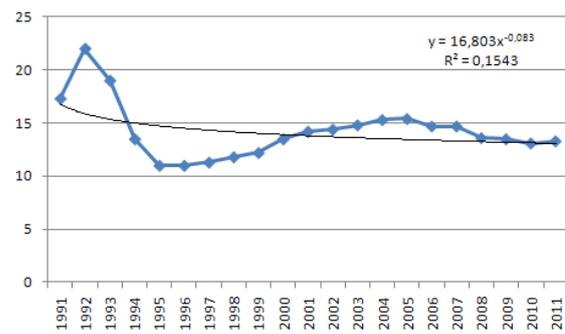


Рис. 6. Степенной тренд

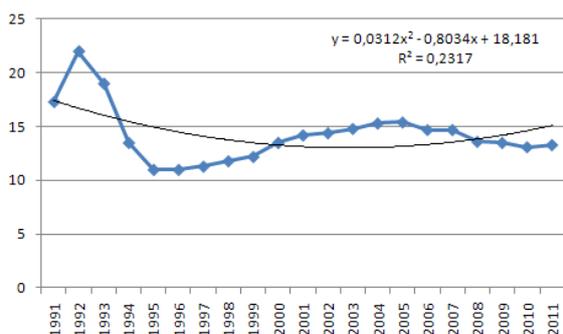


Рис. 7. Полиномиальный тренд (степень 2)

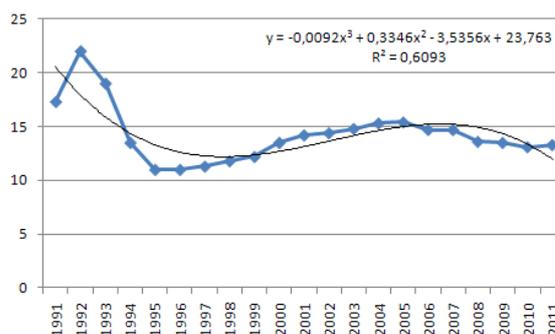


Рис. 8. Полиномиальный тренд (степень 3)

Наиболее высокий коэффициент детерминации получен для полиномиального тренда (степень 3): $R^2 = 0,61$, т.е. 61% вариации изучаемого показателя объясняется изменением года – ростом показателя времени. Причем коэффициенты детерминации для всех остальных рассмотренных видов тренда существенно меньше: для параболического тренда – 0,23, а это меньше более чем в два раза. В связи с этим принцип простоты использовать неуместно (брать более простой вид тренда с несущественно меньшим коэффициентом детерминации).

На данном этапе будем считать, что моделью изучаемого временного ряда является кубический полином:

$$\hat{y}_t = -0,0092t^3 + 0,3346t^2 - 3,5356t + 23,763.$$

Замечания:

1) Увеличение степени полинома всегда приводит к увеличению коэффициента детерминации. При этом на практике очень редко используются полиномы 4-й и 5-й степеней и выше как модели тенденции, т.к. коэффициенты перед переменными в высоких степенях, как правило, статистически не значимо отличаются от нуля. Даже в выбранной нами модели коэффициент перед t^3 близок к нулю, что наводит на мысли о его статистически незначимом отличии от нуля, а квадратный трехчлен как модель тенденции имеет весьма низкий коэффициент детерминации.

2) На рис. 8. Видно, что сглаживающая кривая в конце рассматриваемого промежутка заметно направляется вниз, а сам показатель несколько стабилизировался и не собирается быстро уменьшаться. Это еще один довод, мешающий в реальной ситуации рассматривать в качестве модели тенденции ряда кубический полином.

3) В связи с этим, представляется полезным рассматривать ряд не с 1991 года, а с 1995 года. Можно предположить, что с этого времени произошли структурные изменения и неверно строить одну модель для всего промежутка времени. Длина временного ряда в этом случае пострадает незначительно.

Проверим адекватность и точность выбранной модели

Для этого проведем анализ ряда остатков: $e_t = y_t - \hat{y}_t$, где

$$\hat{y}_t = -0,0092t^3 + 0,3346t^2 - 3,5356t + 23,763.$$

Для проверки адекватности модели необходимо проверить:

- 1) случайность остатков;
- 2) соответствие остатков нормальному закону распределения;
- 3) независимость ряда остатков.

Случайность остатков проверим с помощью проверки гипотезы о существовании тренда, воспользуемся *критерием «восходящих и нисходящих серий»*. Если ряд остатков не содержит тренда, значит он является случайным.

Построим последовательность знаков – плюсов и минусов:

$$\delta_i = \begin{cases} +, & \text{если } y_{t+1} - y_t > 0, & t = 1, 2, \dots, n - 1 \\ -, & \text{если } y_{t+1} - y_t < 0, & t = 1, 2, \dots, n - 1 \end{cases}$$

В случае, когда последующее наблюдение окажется равным предыдущему, учитывается только одно наблюдение. Таким образом, элементы этой последовательности принимают значение «+», если последующее значение уровня ряда больше предыдущего, и «-» если меньше. Общее число знаков «+», «-» заранее неизвестно.

t	y_t	\hat{y}_t	e_t	$e_{t+1} - e_t$	δ_i
1	17,3	20,55	-3,25		
2	22	17,96	4,04	7,30	+
3	19	15,92	3,08	-0,96	-
4	13,5	14,39	-0,89	-3,97	-
5	11	13,30	-2,30	-1,41	-
6	11	12,61	-1,61	0,69	+
7	11,3	12,25	-0,95	0,65	+
8	11,8	12,18	-0,38	0,57	+
9	12,2	12,34	-0,14	0,24	+
10	13,5	12,67	0,83	0,97	+
11	14,2	13,11	1,09	0,25	+
12	14,4	13,62	0,78	-0,31	-
13	14,8	14,14	0,66	-0,11	-
14	15,3	14,60	0,70	0,03	+
15	15,4	14,96	0,44	-0,26	-
16	14,7	15,17	-0,47	-0,90	-
17	14,7	15,16	-0,46	0,01	+
18	13,6	14,88	-1,28	-0,82	-
19	13,5	14,27	-0,77	0,50	+
20	13,1	13,29	-0,19	0,58	+
21	13,3	11,87	1,43	1,62	+

Подсчитываем $\nu(n)$ – число серий в совокупности δ_i , где под серией понимается последовательность подряд идущих плюсов или минусов. Один плюс или один минус тоже будет считаться серией:

$$v(n) = 9.$$

Определяем $\tau_{max}(n)$ – протяженность самой длинной серии:

$$\tau_{max}(n) = 6.$$

Проверим выполнение неравенств:

$$v(n) > \left[\frac{1}{3}(2n - 1) - 1,96 \sqrt{\frac{16n - 29}{90}} \right],$$

$$\tau_{max}(n) < \tau_0(n),$$

где $\tau_0(n)$ – табличное значение, зависящее от n – длины временного ряда. При $n \leq 26$, $\tau_0(n) = 5$.

Проверим выполнение записанных неравенств для нашего примера:

$$v(n) = 9 > \left[\frac{1}{3}(2n - 1) - 1,96 \sqrt{\frac{16n - 29}{90}} \right];$$

$$9 > \left[\frac{1}{3}(40 - 1) - 1,96 \sqrt{\frac{320 - 29}{90}} \right];$$

$$9 > [9,48] = 9 - \text{не верно.}$$

$$\tau_{max}(n) = 6 < \tau_0(n) = 5 - \text{не верно.}$$

Таким образом, оба неравенства не выполняются, значит гипотеза о случайности ряда остатков отвергается и подтверждается наличие зависящей от времени неслучайной составляющей.

Вывод: случайность ряда остатков не подтвердилась.

Проверим соответствие остатков нормальному закону распределения.

При нормальном распределении показатели асимметрии и эксцесса равны нулю.

Мы предполагаем, что отклонения от тренда представляют собой выборку из некоторой генеральной совокупности, поэтому можно определить выборочные характеристики асимметрии (A) и эксцесса (Э):

$$A = \frac{\frac{1}{n} \sum_{t=1}^n e_t^3}{\sqrt{\left(\frac{1}{n} \sum_{t=1}^n e_t^2\right)^3}}, \quad \text{Э} = \frac{\frac{1}{n} \sum_{t=1}^n e_t^4}{\sqrt{\left(\frac{1}{n} \sum_{t=1}^n e_t^2\right)^3}} - 3.$$

Если одновременно выполняются следующие неравенства:

$$|A| < 1,5 \cdot \sqrt{\frac{6(n-2)}{(n+1)(n+3)}}, \quad (1)$$

$$\left| \text{Э} + \frac{6}{n+1} \right| < 1,5 \cdot \sqrt{\frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}}, \quad (2)$$

то гипотеза о нормальном характере распределения случайной компоненты не отвергается.

Если выполняется хотя бы одно из неравенств

$$|A| \geq 2 \cdot \sqrt{\frac{6(n-2)}{(n+1)(n+3)}}, \quad (3)$$

$$\left| \varepsilon + \frac{6}{n+1} \right| \geq 2 \cdot \sqrt{\frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}}, \quad (4)$$

то гипотеза о нормальном характере распределения отвергается.

Выполним необходимые расчеты в таблице:

t	e_t	e_t^2	e_t^3	e_t^4
1	-3,25	10,58	-34,42	111,95
2	4,04	16,35	66,11	267,29
3	3,08	9,49	29,24	90,09
4	-0,89	0,78	-0,69	0,61
5	-2,30	5,29	-12,17	27,98
6	-1,61	2,59	-4,16	6,68
7	-0,95	0,91	-0,87	0,83
8	-0,38	0,15	-0,06	0,02
9	-0,14	0,02	0,00	0,00
10	0,83	0,69	0,58	0,48
11	1,09	1,18	1,29	1,40
12	0,78	0,61	0,47	0,37
13	0,66	0,44	0,29	0,20
14	0,70	0,49	0,34	0,24
15	0,44	0,19	0,08	0,04
16	-0,47	0,22	-0,10	0,05
17	-0,46	0,21	-0,10	0,04
18	-1,28	1,63	-2,09	2,67
19	-0,77	0,60	-0,46	0,36
20	-0,19	0,04	-0,01	0,00
21	1,43	2,04	2,91	4,15
Сумма	0,36	54,49	46,19	515,45
Среднее	0,02	2,59	2,20	24,55

$$A = \frac{\frac{1}{n} \sum_{t=1}^n e_t^3}{\sqrt{\left(\frac{1}{n} \sum_{t=1}^n e_t^2\right)^3}} = \frac{2,2}{\sqrt{2,59^3}} = 0,53,$$

$$\varepsilon = \frac{\frac{1}{n} \sum_{t=1}^n e_t^4}{\sqrt{\left(\frac{1}{n} \sum_{t=1}^n e_t^2\right)^3}} - 3 = \frac{24,55}{\sqrt{2,59^3}} - 3 = 2,87.$$

$$\sqrt{\frac{6(n-2)}{(n+1)(n+3)}} = \sqrt{\frac{108}{21 \cdot 23}} = 0,47,$$

$$\sqrt{\frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}} = \sqrt{\frac{24 \cdot 21 \cdot 19 \cdot 18}{484 \cdot 24 \cdot 26}} = 0,76.$$

Т.к. выполняется неравенство (4):

$$\left| \vartheta + \frac{6}{n+1} \right| = \left| 2,87 + \frac{6}{22} \right| = 3,14 \geq 2 \cdot \sqrt{\frac{24n(n-2)(n-3)}{(n+1)^2(n+3)(n+5)}} = 2 \cdot 0,76 = 1,52.$$

то гипотеза о нормальном характере распределения ряда остатков отвергается

Вывод: соответствие ряда остатков нормальному закону распределения не подтвердилось.

Проверим независимость ряда остатков.

Если ряд остатков не обладает свойством независимости, то существует автокорреляция остатков. Для проверки воспользуемся критерием Дарбина-Уотсона. Критическая статистика рассчитывается по формуле:

$$d = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2}.$$

Расчетное значение статистики d сравнивают с табличными (граничными) значениями d_u (индекс верхней границы) и d_l (индекс нижней границы). Граничные значения d_u и d_l , зависящие от числа наблюдений n , числа объясняющих переменных в модели k' , уровня значимости α , находят по таблице.

В нашем случае: $n = 21$, $k' = 1$ (т.к. в модели один фактор – время), $\alpha = 0,05$, по таблице: $d_u = 1,42$, $d_l = 1,22$.

Алгоритм выявления автокорреляции остатков на основе критерия Дарбина-Уотсона следующий.

Выдвигается гипотеза H_0 об отсутствии автокорреляции остатков.

Затем проверяют в какую область попало расчетное значение статистики d :



Если фактическое значение критерия Дарбина-Уотсона попадает в зону неопределенности, то на практике предполагают существование автокорреляции остатков и отклоняют гипотезу H_0 .

Рассчитаем фактическое значение критерия:

t	e_t	e_t^2	$e_t - e_{t-1}$	$(e_t - e_{t-1})^2$
1	-3,25	10,58		
2	4,04	16,35	7,30	53,23
3	3,08	9,49	-0,96	0,93
4	-0,89	0,78	-3,97	15,73
5	-2,30	5,29	-1,41	2,00
6	-1,61	2,59	0,69	0,48
7	-0,95	0,91	0,65	0,43
8	-0,38	0,15	0,57	0,33
9	-0,14	0,02	0,24	0,06
10	0,83	0,69	0,97	0,94
11	1,09	1,18	0,25	0,06
12	0,78	0,61	-0,31	0,09
13	0,66	0,44	-0,11	0,01
14	0,70	0,49	0,03	0,00
15	0,44	0,19	-0,26	0,07
16	-0,47	0,22	-0,90	0,82
17	-0,46	0,21	0,01	0,00
18	-1,28	1,63	-0,82	0,67
19	-0,77	0,60	0,50	0,25
20	-0,19	0,04	0,58	0,34
21	1,43	2,04	1,62	2,62
Сумма	0,36	54,49		79,08

$$d = \frac{\sum_{t=2}^n (e_t - e_{t-1})^2}{\sum_{t=1}^n e_t^2} = \frac{79,08}{54,49} = 1,45.$$

Так как:

$$d_u = 1,42 < d = 1,45 < 4 - d_u = 2,58,$$

то нет оснований отклонять H_0 (автокорреляция остатков отсутствует).

Вывод: ряд остатков обладает свойством независимости.

Общий вывод: выбранная модель не является адекватной в полной мере, т.к.:

- случайность ряда остатков не подтвердилась;
- соответствие ряда остатков нормальному закону распределения не подтвердилось.

Проверим точность модели.

Характеристика MAPE (Mean Absolute Percentage Error) используется для сравнения точности прогнозов. При этом обычно этот показатель интерпретируется следующим образом:

MAPE < 10% – высокая точность модели;

10% < MAPE < 20% – точность можно признать хорошей;

20% < MAPE < 50% – удовлетворительная точность модели.

$$MAPE = |\bar{\delta}| = \frac{1}{n} \sum_{t=1}^n \left| \frac{e_t}{y_t} \right| \cdot 100\%.$$

При этом n – число уровней временного ряда, для которых определялось прогнозное значение.

Расчеты выполним в таблице:

t	y_t	\hat{y}_t	e_t	$\frac{e_t}{y_t}$	$\left \frac{e_t}{y_t} \right $
1	17,3	20,55			
2	22	17,96			
3	19	15,92			
4	13,5	14,39			
5	11	13,30			
6	11	12,61			
7	11,3	12,25			
8	11,8	12,18			
9	12,2	12,34			
10	13,5	12,67			
11	14,2	13,11			
12	14,4	13,62	0,78	0,05	0,05
13	14,8	14,14	0,66	0,04	0,04
14	15,3	14,60	0,70	0,05	0,05
15	15,4	14,96	0,44	0,03	0,03
16	14,7	15,17	-0,47	-0,03	0,03
17	14,7	15,16	-0,46	-0,03	0,03
18	13,6	14,88	-1,28	-0,09	0,09
19	13,5	14,27	-0,77	-0,06	0,06
20	13,1	13,29	-0,19	-0,01	0,01
21	13,3	11,87	1,43	0,11	0,11
Сумма					0,51
Среднее					0,05

$$MAPE = |\bar{\delta}| = \frac{1}{n} \sum_{t=1}^n \left| \frac{e_t}{y_t} \right| \cdot 100\% = 0,05 \cdot 100\% = 5\%.$$

Вывод: т.к. MAPE < 10%, то у модели высокая точность прогнозирования.

Примечание: при оценке точности модели рассматривали не весь ряд, а последние 10 значений, для которых фактические значения сравнивались с теоретическими.